

PERVASIVE ALTERNATIVE RNA EDITING IN *TRYPANOSOMA BRUCEI*

By

Laura Elizabeth Kirby

A DISSERTATION

Submitted to  
Michigan State University  
in partial fulfillment of the requirements  
for the degree of

Microbiology and Molecular Genetics—Doctor of Philosophy

2019

## ABSTRACT

### PERVASIVE ALTERNATIVE RNA EDITING IN *TRYPANOSOMA BRUCEI*

By

Laura Elizabeth Kirby

*Trypanosoma brucei* is a single celled eukaryote that utilizes a complex RNA editing system to render many of its mitochondrial genes translatable. Editing of these genes requires multiple small RNAs called guide RNAs to direct the insertion and deletion of uridines. These gRNAs act sequentially, each generating the anchor binding site for the next gRNA. This sequential dependence should render the process quite fragile, and mutations in the gRNAs should not be tolerated. In the examination of the gRNA transcriptome of *T. brucei*, many gRNAs were identified that are capable of generating alternative mRNA sequences, and potentially disrupting the editing process. In this work, the effects of alternative editing are characterized. This analysis revealed the role of gRNAs in developmental regulation of gene expression, showing a correlation between the abundance of the initiating gRNAs across two different points in the life cycle of *T. brucei* and their expression. This study also revealed the existence of mitochondrial dual-coding genes, which provide protection for genetic material that is not under selection at all points of the life cycle of *T. brucei*. The examination of these dual-coding genes showed that RNA editing patterns can shift between cell lines and under different energetic conditions. Examining the gRNAs involved in these editing pathways revealed that there is a high amount of mismatching base pairs that are tolerated for editing to function, and that gRNA abundance is not a reliable predictor for editing preference. Finally, a

reexamination of the gRNA transcriptome revealed that many gRNAs are still unidentified and most likely are generating new alternatively edited sequences.

## **ACKNOWLEDGEMENTS**

I would first like to acknowledge my family, without whose support I could not have completed the work I have done. I am forever grateful for their continuous encouragement and constant faith in me.

I gratefully acknowledge Dr. Donna Koslowsky, who has been an excellent mentor. She not only provided her extensive insights on the field of RNA editing, but also aided and advised me in my professional development. She had shaped me as a researcher and as a person, and I was very fortunate in being able to work with her.

I would also like to thank the small undergraduate army I have had the privilege to work with. They taught me a great deal about mentoring and leadership, and their many hours greatly aided me in the completion of this work.

For their assistance in my computational education, I would like to thank Dr. Yanni Sun and Dr. Arend Hintze. Without their assistance, my research would not have been possible.

For serving as my graduate committee, I would like to thank Dr. Shannon Manning, Dr. Charles Hoogstraten, Dr. Chris Adami and Dr. Yanni Sun. Their many contributions have shaped and refined my research.

I would also like to thank Dr. Cori Fata-Hartley, for serving as my teaching mentor and helping me design and conduct an education research project.

I would like to thank the friends I have gained during my time at Michigan State who have supported me and been wonderful colleagues: Alexis Weber, Sandy Olenic, Ahrom Kim, and Shreya Saha.

Finally, I would like to thank the Department of Microbiology and Molecular Genetics, the College of Natural Science, the Elenor L. Gilmore Endowment, the Frank Peabody Microbiology Student Research Fund, the Russell B. DuVall Endowment, the Berttina Wentworth Fellowship, and the Marvis A. Richardson Endowed Fellowship for their support of my research.

## TABLE OF CONTENTS

LIST OF TABLES .....	ix
LIST OF FIGURES.....	x
CHAPTER 1: INTRODUCTION .....	1
Kinetoplastids.....	1
Kinetoplastid RNA Editing .....	1
<i>Trypanosoma brucei</i> .....	2
<i>Trypanosoma vivax</i> .....	6
<i>Trypanosoma cruzi</i> .....	6
<i>Leishmania spp.</i> .....	7
<i>Phytomonas spp.</i> .....	7
Procyclic gRNA Transcriptome .....	8
Evolution and retention of RNA editing in kinetoplastids.....	10
Dual-coding and dual-function genes.....	13
Project Summary.....	15
CHAPTER 2: ANALYSIS OF THE <i>TRYPANOSOMA BRUCEI</i> EATRO 164 BLOODSTREAM GUIDE RNA TRANSCRIPTOME .....	18
Abstract.....	18
Author Summary.....	19
Introduction .....	19
Materials and Methods.....	21
Results.....	23
Discussion.....	37
Accession Numbers.....	40
Acknowledgments.....	40
CHAPTER 3: MITOCHONDRIAL DUAL-CODING GENES IN <i>TRYPANOSOMA BRUCEI</i> .....	42
Abstract.....	42
Author Summary.....	43
Introduction .....	43
Materials and Methods.....	46
Results .....	49
Discussion.....	60
Acknowledgments.....	65
CHAPTER 4: ANALYSIS OF THREE PAN-EDITED MRNAS REVEALS DUAL-CODING GENES AND COMPLEX MULTIPATH EDITING.....	66
Abstract.....	66

Introduction .....	67
Materials and Methods.....	70
Results.....	73
Discussion.....	93
Acknowledgements.....	99
 CHAPTER 5: CLUSTER CLASSIFICATION OF UNKNOWN GRNAS REVEALS THE ROBUSTNESS OF THE RNA EDITING SYSTEM .....	100
Abstract.....	100
Introduction .....	101
Materials and Methods.....	105
Results.....	108
Discussion.....	129
Acknowledgements.....	134
 CHAPTER 6: SUMMARY AND DISCUSSION .....	135
Introduction .....	135
Summary of Chapter 2 .....	136
Summary of Chapter 3 .....	137
Summary of Chapter 4 .....	138
Summary of Chapter 5 .....	139
Genetic Integrity .....	140
Developmental Regulation .....	144
Protein Diversity .....	146
Editing Efficiency.....	147
Future Work .....	149
Conclusion.....	151
 APPENDICES .....	152
APPENDIX A. Quantification of the number of identified bloodstream and procyclic gRNA transcripts that cover a respective nucleotide in the fully edited mRNA.....	153
APPENDIX B. Alignment of the mitochondrial fully edited mRNAs and the most abundant gRNAs required for full coverage identified in the bloodstream and procyclic life cycle stages .....	158
APPENDIX C. All gRNA major classes pulled for ATPase 6 in the EATRO 164 procyclic and bloodstream transcriptomes .....	187
APPENDIX D. Identified CR3 mRNA and gRNA transcripts.....	215
APPENDIX E. ND7 5'-most gRNA populations and the predicted mRNA sequences generated.....	217
APPENDIX F. RPS12 5'-most gRNA populations and the predicted mRNA sequences generated.....	219
APPENDIX G. Alignments of <i>T. brucei</i> and <i>T. vivax</i> edited mRNAs .....	220

APPENDIX H. Alignments of protein sequences of pan-edited dual-coding genes in <i>L. tarentolae</i> , <i>L. amazonensis</i> , <i>P. serpens</i> , and <i>Perkinsela CCAP1560/4</i> with <i>T. brucei</i> and <i>T. vivax</i> sequences .....	227
APPENDIX I. RPS12 gRNA Alignments for TREU 667 SDM79 and EATRO 164 SDM79 cells , and all editing variants.....	234
APPENDIX J. gRNAs identified to edit the RPS12 mRNAs of found in both TREU 667 and EATRO 164 gRNA transcriptomes .....	241
APPENDIX K. ND7 gRNA Alignments for TREU 667 SDM79 and EATRO 164 SDM79 cells, and all editing variants.....	244
APPENDIX L. gRNAs identified to edit the ND7 5' mRNAs of found in both TREU 667 and EATRO 164 gRNA transcriptomes .....	251
APPENDIX M. Predicted ND7 protein sequences .....	253
APPENDIX N. CR3 gRNA Alignments for TREU 667 SDM79, and all editing variants.....	254
APPENDIX O. CR3 gRNA Alignments for EATRO 164, and all editing variants .....	261
APPENDIX P. gRNAs identified to edit the CR3 mRNAs of found in both TREU 667 and EATRO 164 gRNA transcriptomes .....	270
REFERENCES .....	274

## LIST OF TABLES

Table 1. Differences in mitochondrial transcript abundance, polyadenylation and the extent of RNA editing in two life cycle stages of <i>T. brucei</i> .....	5
Table 2. Number of gRNA transcripts in procyclic and bloodstream major classes and ratio of procyclic transcripts to bloodstream transcripts for each gene.....	25
Table 3. Summary of the gRNA data coverage for each gene. ....	26
Table 4. Most common gRNA transcription start sites in procyclic and bloodstream data..	29
Table 5. Identified gaps or weak overlaps (less than 6 nucleotides) between populations of gRNAs observed in both data sets. ....	31
Table 6. Summary of populations found in both data sets that have more reads in the bloodstream data set than in the procyclic data set. ....	32
Table 7. Editing efficiencies of RPS12, ND7, and CR3. ....	76
Table 8. Editing efficiency for each RPS12 gRNA population. ....	80
Table 9. Editing efficiency for each ND7 5'domain gRNA population .....	84
Table 10. Editing efficiencies by block level of CR3 .....	85
Table 11. Summary of ACORNS Results .....	110
Table 12. Cluster summary .....	110
Table 13. Cluster size summary .....	110
Table 14. gRNA population analyses for RPS12. ....	118
Table 15. gRNA population analysis for ND7 5' .....	124
Table 16. CR3 gRNA population analysis .....	127

## LIST OF FIGURES

Figure 1. The abundance of the initiating gRNA of all edited mRNAs in each stage.....	25
Figure 2. The frequency of nt variations versus nucleotide position in the gRNA. ....	27
Figure 3. Comparing the number of non-complementary nucleotides 5' of the anchoring region or 3' of the guiding region in procyclic and bloodstream gRNAs.....	28
Figure 4. Length of gRNA complementarity to fully edited mRNAs for both bloodstream and procyclic gRNAs.....	28
Figure 5. The percentage of different nucleotide overlaps found between adjacent gRNAs... .....	30
Figure 6. Alignment of conventional ATPase 6 protein sequence to hypothetical proteins generated by the 11U alternative edited mRNA and the 4U alternatively edited mRNA. ...	33
Figure 7. Editing sites 420–489 of COIII aligned with the gRNAs identified for that region in the procyclic and bloodstream data sets.....	34
Figure 8. Alternative editing of the 5' end of pan-edited genes results in access to different reading frames. ....	52
Figure 9. Positions of stop codons on all RFs of the edited genes in <i>T. brucei</i> .....	53
Figure 10. Mutational frequencies in mitochondrialy encoded genes categorized by effect on amino acid sequence. ....	55
Figure 11. Percent conservation of editing patterns between <i>T. brucei</i> and <i>T. vivax</i> .....	56
Figure 12. Principal component analysis of frequency of amino acid mutation types and editing conservation between <i>T. brucei</i> and <i>T. vivax</i> pan-edited transcripts. ....	58
Figure 13. Amino acid sequences of ARFs of dual-coding genes.....	60
Figure 14. Observed RPS12 editing pathways in the TREU 667 cell line and the EATRO 164 cell line grown in SDM79 and SDM80.....	79
Figure 15. Alignment of RPS12 proteins from <i>T. brucei</i> , <i>T. vivax</i> , <i>Leishmania tarentolae</i> , <i>Leishmania donovani</i> , and <i>Leishmania amazonensis</i> . .....	80

Figure 16. Regions with poor gRNA coverage and functionally conserved residues in RPS12 .	81
Figure 17. Observed ND7 5' editing pathways in the TREU 667 cell line and the EATRO 164 cell line grown in SDM79 and SDM80.....	82
Figure 18. Regions with poor gRNA coverage and functionally conserved residues in ND7 5'	84
Figure 19. Observed CR3 editing pathways in the TREU 667 cell line .....	87
Figure 20. Four different 3' end sequences found in the TREU 667 transcriptome for the CR3 transcript and CR3 protein sequences.....	88
Figure 21. Observed CR3 editing pathways in the EATRO 164 cell line grown in SDM79 and SDM80.....	90
Figure 22. Alignment of CR3 predicted protein variants from the EATRO 164 cell line.....	91
Figure 23. Predicted secondary structures of most abundant CR3 predicted proteins .....	91
Figure 24. Frequencies of early total deletions of DNA encoded uridines in partially edited ND7 and RPS12 transcripts .....	93
Figure 25. Example clusters of related gRNAs generated by ACORNS from the EATRO 164 PC gRNA transcriptome.....	112
Figure 26. Observed RPS12 editing pathways in the TREU 667 cell line the EATRO cell line....	117
Figure 27. Analysis of functionality and abundance of productive gRNAs populations that edit RPS12 in TREU 667 cells.....	119
Figure 28. Analysis of functionality and abundance of productive gRNAs that edit RPS12 in EATRO 164 cells.....	120
Figure 29. Observed ND7 5' editing pathways in TREU 667 cell line and the EATRO 164 cell line.....	122
Figure 30. Analysis of functionality and abundance of productive gRNAs that edit that edit ND7 5' in TREU 667 cells .....	123
Figure 31. Analysis of functionality and abundance of productive gRNAs that edit ND7 5' in EATRO 164 cells.....	123

Figure 32. Observed CR3 editing pathways in TREU 667 and EATRO 164 cell lines .....	126
Figure 33. Analysis of functionality and abundance of gRNA subpopulations that edit CR3 in TREU 667 cells .....	128
Figure 34. Analysis of functionality and abundance of gRNA subpopulations that edit CR3 in EATRO 164 cells.....	128

# CHAPTER 1: INTRODUCTION

## Kinetoplastids

*Trypanosoma brucei* is a member of the Kinetoplastea, a group of protozoans characterized by a large network of DNA in their mitochondria known as the kinetoplast that is physically attached to the flagellum [1]. While not all kinetoplastids are parasites, the group encompasses some of the most successful parasites in existence, inhabiting an incredibly wide range of hosts from plants to invertebrates to vertebrates [2,3]. The dixenous members cycle between two distinct hosts and can encounter different environments with distinct metabolic constraints. These environmental shifts require rapid and extensive changes in gene expression. This is particularly interesting considering the kinetoplastids' bizarre and complicated use of RNA editing for their mitochondrial gene expression.

## Kinetoplastid RNA Editing

RNA editing is one of several unique genetic features found in the mitochondria of these parasites. RNA editing creates open reading frames in “cryptogenes” by insertion and deletion of uridylate residues at specific sites within the mRNA. The U-insertions/deletions are directed by small guide RNAs (gRNAs) and can repair frameshifts, generate start and stop codons and more than double the size of the transcript [4].

The kinetoplast DNA (kDNA) consists of two types of DNA molecules, maxicircles and minicircles. Maxicircles are large circular DNA molecules that contain the genes for two ribosomal RNAs, 12S and 9S, and the protein coding genes [5]. While some of the protein-coding genes do not require RNA editing prior to translation, most require extensive editing

before they can be translated [6,7]. The sequence changes are guided by small complementary RNA molecules (the gRNAs) that are encoded on the minicircles [8]. Minicircles make up the bulk of the kinetoplastid network with each minicircle encoding 1–5 gRNAs.

This effectively means that the genetic information for the edited mitochondrial mRNAs is dispersed between the mRNA cryptogenes on the maxicircles and as many as 10,000 gRNA encoding minicircles. In *T. brucei*, the extensive editing of a single transcript can require more than 40 gRNAs and hundreds of editing events [9]. The gRNAs act as templates for the large multi-subunit protein complex known as the editosome [4,6]. The editosome cleaves the mRNA, inserts or deletes the correct number of uridines and then re-ligates the mRNA in an energy intensive process. This is repeated until the mRNA is complementary to the small gRNA. The initiating gRNA interacts with the 3' end of the pre-edited transcript and generates the anchor binding region for the next gRNA. In fact, all subsequent gRNAs anchor to the edited sequence created by the preceding gRNA. Editing proceeds from the 3' end to the 5' end of the mRNA transcript with the terminating gRNA either creating the start codon or bringing an existing start codon into frame. Because each gRNA directs editing that generate the anchor region for the next gRNA, the RNA editing process is sequentially dependent on correct editing by each gRNA. As a result, the process is incredibly fragile.

### **Trypanosoma brucei**

*Trypanosoma brucei* is the causative agent of Human African trypanosomiasis (HAT) and one agent of Animal African Trypanosomiasis (AAT). Each year, 10,000 new cases of HAT are reported, and 3 million cattle are killed, severely impacting the lives and livelihood of those in infected areas [10,11]. The trypanosomes live in two distinct environments: the animal host

and the insect vector, the tsetse fly. These environments are distinct in temperature and nutrient composition, providing a challenge to *T. brucei* as it cycles between hosts. While in the mammalian host, *T. brucei* lives entirely extracellularly in the bloodstream. It is frequently subject to attacks by the host's adaptive immune system, and the population evades these attacks through antigenic variation [12]. This part of the life cycle can be quite long, with the longest known infection lasting 29 years [13]. In the bloodstream, the bulk of the trypanosome population exists in the actively dividing slender form. The slender form is optimized to utilize its glucose rich environment, using glycolysis to generate energy [14,15]. During this stage of the life cycle, the mitochondrion is down-regulated, lacking both Krebs cycle enzymes and a functional electron transport chain (ETC) [16]. While the activity of the mitochondrion is relatively low during the bloodstream stage (BS), expression of the mitochondrial genome is still essential [17,18]. Once the population reaches an optimum density, a small portion of the population transitions into stumpy form trypanosomes. The stumpy form is nondividing and appears to be transitional, activating mitochondrial genes in preparation for uptake in a blood meal by its tsetse fly vector and subsequent transfer to a harsher environment [14]. Successful transition to the fly vector requires activation of the ETC, and ATP synthesis via oxidative phosphorylation. Once inside the tsetse fly, the parasite utilizes proline as its primary energy source while residing and actively dividing in the midgut [19–21]. This stage of the life cycle is followed by a dramatic bottleneck when the trypanosomes transition from the midgut to the salivary glands of the tsetse fly, with as few as 1-5 trypanosomes completing the transition [22,23]. From the salivary glands, trypanosomes are then refluxed into their next mammalian

host during a bloodmeal. In order to adapt to these sudden changes in environment, *T. brucei* must vastly alter its gene expression, most notably, in its mitochondria.

The 22 kb maxicircle of *T. brucei* encodes several genes involved in the mitochondrial ETC and oxidative phosphorylation, NADH dehydrogenase (ND) subunits 1-5 and 7-9, cytochrome oxidase (CO) subunits I-III, cytochrome b (CYb), ATP synthase subunit 6 (A6), as well as genes encoding the ribosomal protein small subunit 12 (RPS12), 12S and 9S rRNAs, and some genes with unknown functions: C-rich regions (CR) 3 and 4, and Maxicircle unidentified reading frames (Murf) 2 and 5 [5]. Twelve of these genes require some amount of RNA editing to be translatable, with some requiring only one or two gRNAs (COII, CYb, MurfII), and others requiring editing across the span of the transcript (ND3, ND7, ND8, ND9, COIII, A6, RPS12, CR3, and CR4) [4,6,7].

Distinct differences in mitochondrial transcript abundance, polyadenylation and the extent of RNA editing are observed during the complex life cycle (Table 1). The pattern of differential RNA editing observed is especially interesting. For example, the CYb and COII mRNAs are edited during the insect stage, but are primarily unedited in bloodstream forms [24,25]. In contrast, editing of the NADH dehydrogenase subunit transcripts (ND3, ND7, ND8 and ND9) and RPS12 appears to occur preferentially in bloodstream forms [5,26–30]. Other transcripts, COIII and A6 are edited equally in both life cycle stages [31,32].

Aside from the genes encoded on the maxicircle, there are the minicircle gRNAs. The minicircles range in abundance from 5,000–10,000 present in each network, and are ~1kb in size, with each minicircle encoding 2–5 gRNAs. In *T. brucei*, there are more than 200 different minicircle sequence classes (~1200 gRNAs) [8,33]. While the minicircles make up a bulk of the

**Table 1. Differences in mitochondrial transcript abundance, polyadenylation and the extent of RNA editing in two life cycle stages of *T. brucei*.**

Gene	No. of uridines		Edited size (nt)	Stage Edited	Relative level of mature RNA			PolyA tail length		Number of PC major classes <sup>g</sup>	References
	Added	Deleted			Long Slender	Short Stumpy	Procyclic	Bloodstream	Procyclic		
12S	2-17 (tail)	0	1149	N.D.	0.04	1.3	1.0	N.A.	N.A.	N.A. <sup>j</sup>	[24,34,35]
9S	7(tail)	0	611	N.D.	0.07	1.4	1.0	N.A.	N.A.	N.A.	[24,34,35]
CYb	34	0	1,151	P <sup>a</sup>	~0	0.5	1.0	Short (UE <sup>h</sup> & E <sup>i</sup> )	Short (E) & Long (E)	11	[24,25,36,37]
A6	448	28	821	P/BS <sup>b</sup>	1.0	N.D. <sup>f</sup>	1.0	Short (E) & Long (E)	Short (E) & Long (E)	81	[31,36]
COI	0	0	1,647	NE <sup>c</sup>	0.07	0.4	1.0	Short	Short & Long	N.A.	[24,37,38]
COII	4	0	663	P	~0	0.5	1.0	Short	Short (UE) & Long (E)	N.A. <sup>k</sup>	[24,37,39]
COIII	547	41	969	P/BS	1.0	N.D.	1.0	Short (E)	Short(E) & Long(E)	151	[32]
ND1	0	0	960	NE	~1	N.D.	~1	Short & Long	Short & Long	N.A.	[38,40,41]
ND2	0	0	1,343	NE	>1.0	N.D.	1.0	Short & Long	Short & Long	N.A.	[37,42,43]
ND3	210	13	452	P/BS	>1.0	N.D.	1.0	Short (E)	Short (E)	34	[26]
ND4	0	0	1,314	NE	~1.0	N.D.	~1.0	Short & Long	Short & Long	N.A.	[37,38,44,45]
ND5	0	0	1,779	NE	0.5	0.8	1.0	N.D.	N.D.	N.A.	[24,38]
ND7	553	89	1,238	5'P/BS, 3'BS <sup>d</sup>	~10	N.D.	1.0	Short(E)	Short(UE)	129	[5,27]
ND8	259	46	574	BS <sup>e</sup>	~20	N.D.	1.0	Short (E) & Long (E)	Short (E)	70	[5,28,37]
ND9	345	20	649	BS	>1.0	N.D.	1.0	Short (PE) & Long (E)	Short (PE)	39	[29]
RPS12	132	28	325	BS	>1.0	N.D.	1.0	Short(UE & E) & Long (E)	Short(UE & E) & Long (E)	50	[30]
Murf 2	26	4	1,111	P/BS	~1	~1	~1	Short & Long	Short & Long	1	[40,46]
Murf 5	N.D.	N.D.	N.D.	N.D.	N.D.	N.D.	N.D.	N.D.	N.D.	N.A.	N.D.
CR3	148	13	299	BS	>1.0	N.D.	1.0	N.D.	N.D.	37	[47]
CR4	325	40	567	BS	1.0	N.D.	~0	Short (E) & Long (E)	Short (UE)	41	[48]

<sup>a</sup>P, transcript is edited only in the procyclic (insect) developmental stage.

<sup>b</sup>P/BS, transcript is edited in both bloodstream and procyclic stages.

<sup>c</sup>NE, never edited, editing of these transcripts has not been reported.

<sup>d</sup>The ND7 transcript is differentially edited in the procyclic and bloodstream stages.

<sup>e</sup>BS These transcripts are only fully edited in the bloodstream developmental stage [49].

<sup>f</sup>N.D. Values have not yet been determined.

<sup>g</sup>All data comes from EATRO 164 procyclic gRNA transcriptome previously published [9].

<sup>h</sup>UE, unedited, the transcripts which carried these tails were typically unedited.

<sup>i</sup>E, edited, the transcripts carrying these tails were typically edited.

<sup>j</sup>N.A., Not applicable.

<sup>k</sup>COII is a cis-edited transcript. Poly-A tails listed as short are between 10 and 50 nts long and tails listed as long are between 150 and 200 nts long.

mitochondrial DNA, early studies using both Northern blot and primer extension analyses on a limited number of gRNAs indicate that gRNAs are present in both insect and bloodstream forms, suggesting that the regulation of RNA editing is not at the level of gRNA availability [28,50,51].

### **Trypanosoma vivax**

Like *T. brucei*, *Trypanosoma vivax* is a causative agent of AAT. Its kinetoplast DNA is also very similar. The maxicircles in *T. vivax* possess the same genes as the maxicircle of *T. brucei*, and the genes that are edited in *T. brucei* are also edited in *T. vivax* to the same extent [52]. The minicircles of *T. vivax* vary significantly in size, from 300-1100 bp, encoding 1-3 gRNAs [52]. The life cycle of *T. vivax* is highly similar to that of *T. brucei*. In its mammalian host, it lives extracellularly in the bloodstream, primarily metabolizing glucose. When it is taken up by the tsetse fly in a bloodmeal, it initially resides in the proventriculus and foregut. From there, cells migrate to the proboscis and propagate, preparing to be deposited into the next mammalian host [53].

### **Trypanosoma cruzi**

*T. cruzi* is known for causing Chagas disease in South and Central America and is carried between hosts by triatomine bugs. Like *T. brucei* and *T. vivax*, it has highly similar kDNA, with the maxicircle containing the same genes in the same order. The edited genes of *T. brucei* and *T. vivax* are also edited in *T. cruzi* to the same extent [54]. Like *T. brucei*, *T. cruzi* spends most of its insect stage in the nutrient depleted midgut of its host, metabolizing amino acids for survival [55,56]. Once cells have propagated, they migrate to the hindgut and are excreted from the fly.

Transmission to the mammalian host occurs by contact with a wound or mucous membrane.

Once inside the mammalian host, *T. cruzi* invades many different types of nucleated cells by using the microtubule cytoskeleton of the host cell to recruit lysosomes to create vacuolar compartments where *T. cruzi* resides [57]. Once inside the vacuole, *T. cruzi* metabolizes glucose as its primary energy source [58,59].

### ***Leishmania spp.***

*Leishmania spp.* are found on almost every continent in the world, and infect 700,000 to 1.2 million people annually [60]. *Leishmania spp.* are transmitted by the phlebotomine sand fly. Unlike the trypanosomes, while in the midgut of the sand fly, *Leishmania spp.* primarily metabolize glucose, because of the fly's frequent sap meals. *Leishmania spp.* are transmitted to their mammalian hosts by the bite of the sand fly. Once inside the host, they are phagocytosed by host cells. Inside macrophages, they replicate within lysosome like compartments, and it is believed that these compartments are not nutrient restrictive [61,62]. The maxicircle of the *Leishmania spp.* parasites has the same genes as the trypanosomes, but their editing patterns significantly vary. The pan-edited genes in *Leishmania spp.* are ND3, ND8, ND9, RPS12, CR3, and CR4, and the partially edited genes are A6, COII, COIII, MurfII, and ND7 [63].

### ***Phytomonas spp.***

*Phytomonas spp.* parasitize plants, utilizing their sucrose and polysaccharides as energy sources. Their insect vectors maintain a highly sap rich diet that allows the parasites to continually metabolize carbohydrates, unlike the *Trypanosoma spp.* Possibly as a result of this,

several metabolic pathways are incomplete. The pathways for beta oxidation of fatty acids or oxidation of amino acids are missing key enzymes, but the pathways for the synthesis of these metabolites are more complete [64]. The ETC of the mitochondria is also affected. Genes for all cytochromes are missing from nuclear and kDNA, and the cytochrome oxidase genes normally present on the maxicircle are also missing [64,65]. The other maxicircle genes are present, and ND3, ND8, ND9, RPS12, CR3 and CR4 are pan-edited, while ND7, A6 and MurfII are only partially edited, as with *Leishmania spp.*

### **Procyclic gRNA Transcriptome**

The gRNA transcriptome of insect stage (procyclic) *T. brucei* was previously sequenced [50]. This library was generated from the EATRO 164 cell line, grown in SDM79 medium, the most commonly used medium when culturing procyclic trypanosomes. As no reference genome exists for the minicircles, gRNAs could only be identified based on their function. Using a longest common substring algorithm, gRNAs were identified based on their complementarity to previously determined fully edited mRNA sequences. The RNA editing system tolerates G:U base pairs, so this program allows these base pairs in alignments, but these base pairs do not contribute to the overall alignment score as much as canonical Watson-Crick base pairs (1 point for G:U and 2 points Watson-Cricks). Guide RNAs with scores higher than 45 points were identified as editing gRNAs. Due to the fact that trypanosomes post transcriptionally add a poly-uridine (poly-U) tail to gRNAs, the sequences generated in this transcriptome possess the poly-U tail as well [66]. This program ignores the transcript's poly-U tail in the alignment, and it does not contribute to the score. Using this program, full complements of gRNAs were found for A6, COIII, CR4, CYb, and RPS12, and near full complements were identified for the

other edited genes. This study found that multiple different sequence classes of gRNAs (major classes) edited the same region of an mRNA (this group of gRNAs is called a population). Major classes within a population had many transition mutations, primarily A-G mutations, which appeared to be due to the editing system's toleration of G-U base pairs. Interestingly, populations of gRNAs varied extremely in transcript abundance, with abundance varying from <10 to >350,000 reads.

The gRNAs identified in this study possessed common characteristics. 64% of transcripts had 38-48 nucleotides (nt) of complementarity to their target mRNA. 84% of transcripts had 6 or fewer non-complementary nts at 5' end, and most transcripts had 0 nts non-complementary at 3' end prior to the poly-U tail. Conservation was observed in the gRNA transcription start site, with 74% of transcripts starting with 5'-ATATA-3'. Interestingly, a large proportion of transcripts had 5'-AAAAA-3' transcription start sites as well.

Beyond the identification of gRNAs directing the conventional mRNA edits, this study identified a number of gRNAs that could generate alternative edits. Most of these edits caused minor changes to the predicted mRNA and protein sequences, by either changing a single amino acid (ND8) or changing no amino acids at all (A6). However, some gRNAs were identified dramatically altered the mRNA or protein sequence. One edit in the essential A6 gene caused a frameshift that would alter and shorten the C-terminus of the protein [17]. Another generated a dramatically different sequence at the 3' end of CR3, and no gRNA was identified capable of editing that sequence. Interestingly, another study identified an alternative edit in COIII, that linked an open reading frame in the unedited 5' end of the transcript with the reading frame in

the edited 3' end of the transcript [67,68]. The gRNA required to generate this alternative edit was not identified in the procyclic gRNA transcriptome of *T. brucei*.

## **Evolution and retention of RNA editing in kinetoplastids**

The kinetoplastid RNA editing system is energetically expensive and, due to the system's sequential dependence, should be highly fragile. This means that with even high accuracy rates for each gRNA, the overall fidelity of the process is astonishingly low. Even a single point mutation could drastically change the editing pattern, and stop the editing process, aborting expression of the protein. A major question in the field has been why this fragile and metabolically expensive system of RNA editing would evolve and persist.

Initially, it was proposed that U-insertion/deletion editing (kRNA editing) was one of many RNA editing processes that were in fact relics of the RNA world. However, the very different mechanisms of the RNA editing systems in existence and their very limited distribution within specific groups of organisms indicate that they are more likely derived traits that evolved later in evolution [69,70]. The sheer complexity of the kRNA editing process, with no obvious selective advantage, led to the proposal that insertion/deletion editing arose via a constructive neutral evolution (CNE) pathway [71]. RNA editing in trypanosomes is always mentioned in support of CNE as an example of how seemingly non-advantageous, complex processes can arise [72,73]. More recently however, it has been hypothesized that RNA editing co-evolved with G-quadruplex structures found in the pre-edited mRNAs [74]. These structures can help regulate transcription in order to promote DNA replication and prevent kDNA loss, and thus provide an advantage to the organism. However, they must be removed by the RNA editing system prior to translation [74]. Another prominent hypothesis is that RNA editing is

advantageous because it is a mechanism by which an organism can fragment and scatter essential genetic information throughout a genome [75,76]. Kinetoplast DNA is far less stable than chromosomal DNA, and loss of minicircles due to asymmetric division of the kDNA network have been frequently observed, particularly in laboratory cultures of *Leishmania tarentolae* [77,78]. Buhrman et al. [76] suggest that the scattering of essential gRNA genes throughout the DNA network would prevent fast growing deletion mutants from outcompeting more metabolically versatile parasites during growth in the mammalian host. Using a mathematical model of gene fragmentation in changing environments (absence of functional selection), they showed a distinct advantage for gene fragmentation. In their model, the number of tolerable generations under periods of relaxed selective pressure was increased by more than 40% before loss of the ability to move to the next life cycle stage.

One mechanism for protecting small asexual populations is by increasing the severity of the mutations that can occur. If mutations severely impact fitness, deleterious mutations are selected out, preventing their fixation [79]. This phenomenon increases the ‘drift robustness’ of a population. One study modeled the acquisition of drift robustness mathematically and computationally [80]. This study showed that in a simulated environment, small populations evolved a lower fitness than large populations, but when the most common genotypes from these populations were placed in a scenario with extremely high genetic drift, the genotypes evolved from small populations experienced a smaller decline in fitness than the genotypes evolved in a large population. Furthermore, they examined the types of mutations in the fitness landscape nearest to the peaks that each simulated population had fixed on and found that in smaller populations, there was an excess of mutations possible that were neutral,

beneficial or strongly deleterious, whereas in larger populations, there were more small-effect deleterious mutations possible. As the RNA editing process may be operating as a proof-reading system to weed out mutations by making them lethal, these findings strongly support the hypothesis that RNA editing is beneficial to the trypanosomes by providing a level of drift robustness to the population as a whole.

While these hypotheses do address the evolution and retention of the RNA editing system itself, they do not address another key issue with the system as a whole: maintaining genetic material used in RNA editing while that material is not under selection. During the life cycle of *T. brucei*, trypanosomes undergo a severe bottleneck as they transition through the tsetse fly and into the mammalian host, and then within the bloodstream, they undergo multiple bottlenecks at each antigenic switch, as they evade the host immune system [22]. Such bottlenecks create additional forces of genetic drift, where genes can be lost even if their deleterious fitness effect is considerable. This life cycle should make *T. brucei* particularly sensitive to genetic drift, especially for those genes which are not under selection (Krebs cycle and ETC) and should make them extremely vulnerable to Muller's ratchet (the gradual increase of mutational load that eventually leads to extinction) [81–84].

During a reexamination of the EATRO 164 procyclic gRNA transcriptome, a number of gRNAs were identified capable of shifting the open reading frame of their respective transcripts. These gRNAs acted at the 5' end of edited mRNAs and either shifted the position of an existing start codon or generated a new start codon that would allow that transcript to be translated in an alternative reading frame. Surprisingly, the alternative reading frames spanned the full or nearly full length of their transcripts, suggesting that these transcripts were capable

of generating two distinctly different protein products. Based on these observations, we hypothesize that trypanosomes use dual-coding genes to protect genetic information by essentially hiding a gene not under selection (i.e. ETC genes) within one that remains under selection. Thus, the ability to access overlapping reading frames may be one explanation for how genetic material that is unused in one life cycle stage may be preserved while it is not under selection.

### **Dual-coding and dual-function genes**

Dual-coding genes are defined as a stretch of DNA containing overlapping open reading frames (ORFs) [85,86]. Overlapping reading frames are common in viruses and are thought to persist due to strong genome size constraints [87,88]. More recently however, overlapping genes have been identified in mammalian and bacterial genomes [89–92]. In these organisms, size is not an issue and the potential advantage of overlapping genes is less clear. Maintaining dual-coding genes is costly, as it constrains the flexibility of the amino acid composition of both proteins, constraining the ability of each protein to become optimally adapted [93]. As this constraint can be alleviated by gene duplication, it is thought that dual-coding regions can survive long evolutionary spans only if the overlap provides a selective advantage. In mammals, many of the identified dual-coding genes produce two proteins that bind and regulate each other [94,95]. For these proteins, dual-coding may be advantageous for the tight co-expression needed. An alternative model suggests that under high mutation rates, the overlapping of critical nucleotide residues is advantageous because it may reduce the target size for lethal mutations [96].

The use of genetic information with more than one function is not a new idea in *T. brucei*. The nuclear encoded α-ketoglutarate dehydrogenase E2 (α-KDE2) is known to be a dual-function protein, in that it plays important roles in both the Krebs cycle and in mitochondrial DNA inheritance [97]. RNAi knockdowns of this gene in bloodstream form (BF) trypanosomes also show a pronounced reduction in cell growth. Similarly, the Krebs cycle enzyme α-ketoglutarate decarboxylase (α-KDE1) is a dual-function protein with overlapping targeting signals that allow it to be localized to both the mitochondrion and glycosomes [98]. RNAi knockdowns of α-KDE1 in BF trypanosomes is lethal, suggesting that, in addition to its enzymatic role in the Krebs cycle, it plays an essential role in glycosomal function in *T. brucei* [98].

Another example of this was identified in the RNA editing system. Alternative editing of COIII is reported to generate a novel DNA-binding protein, Alternatively Edited Protein-1 (AEP-1), that functions in mitochondrial DNA maintenance [67,68]. In this transcript, one alternative gRNA generates sequence changes at two sites that links an open reading frame (ORF) found in the pre-edited 5' end, to the 3' transmembrane domains found in the COIII edited ORF. This was the first indication that one cryptogene could contain information for more than one protein. It has been previously suggested that both alternative editing and dual-function proteins are important mechanisms for expanding the functional diversity of proteins found in trypanosomes [67,97–99]. We hypothesize that because trypanosomes live exclusively extracellularly in their mammalian host, they are more sensitive to genetic drift, and an equally important role for these dual-coding/function genes may be the protection of genetic information.

## Project Summary

The goal of this work is to examine the impact of alternative editing on the protein diversity, editing efficiency, developmental regulation, and genetic integrity of *Trypanosoma brucei*. This investigation began with the generation of the gRNA transcriptome of bloodstream form EATRO 164 *T. brucei*. This analysis identified near full complements of gRNAs for the edited genes, as was discovered in the procyclic transcriptome. A detailed comparison of the gRNAs identified in both datasets revealed conserved characteristic, such as anchor length, length of complementarity, and transcription start site sequences, even though very few identical sequences existed between the two transcriptomes. Additionally, an interesting correlation was found that suggests a relationship between the relative abundance of initiating gRNAs between stages and the developmental pattern of mRNA editing.

During this comparison of the two transcriptomes, a number of alternative editing gRNAs were identified. Notably, three of these gRNAs were capable of shifting translation of ND7, RPS12, and CR3 into alternative reading frames. This discovery prompted the analysis of the mitochondrially encoded transcripts to determine which of the genes had the capacity to be dual-coding. Using mutational bias analysis, we show that as many as six cryptogenes in addition to the previously discovered COIII/AEP-1, encode more than one protein, and that RNA editing allows access to both reading frames.

In order to determine if mRNA transcripts with access to multiple open reading frames exist within the mitochondrial transcriptome, we deep sequenced the transcript populations of three putative dual coding genes: RPS12, the 5' editing domain of ND7 (ND7 5'), and CR3. Using the previously generated gRNA transcriptomes, we constructed detailed editing pathways for

each of these genes. We found evidence that CR3 and ND7 5' are dual-coding genes, based on the identification of transcripts that would translate into different reading frames. This study indicates that RNA editing can be used to access multiple open reading frames using two different methods: in ND7 5', different gRNAs bring alternate start codons into frame and in CR3, different gRNAs can shift the reading frame of the existing start codon. In addition, CR3 showed incredible editing diversity. In two different cell lines, highly divergent editing patterns were characterized, with the two cell lines using different sets of gRNAs to edit the CR3 cryptogene. This suggests that the use of a gRNA-guided editing system can also dramatically increase protein diversity in spite of an incredibly rigid and mutationally fragile system.

With a more complete understanding of the existing edits found in the mRNA transcriptomes of RPS12, ND7 5' and CR3, we used this knowledge to analyze the RNA editing system's ability to tolerate noise. Reexamining the procyclic gRNA transcriptome revealed many previously unidentified gRNAs, potentially capable of generating alternative edits or disrupting the editing system. Using a new program called ACORNS (Assemble Clusters Of Related Nucleotide Sequence), the gRNAs were grouped into clusters based on sequence homology. This allowed us to determine which unidentified gRNAs were related to previously identified gRNAs. This analysis showed that more than half of the unidentified gRNAs were not related to any gRNA of known function, suggesting that many more alternative edits are waiting to be discovered. In order to analyze the impact of the gRNAs that were related to previously identified gRNAs, another new program, GUIDE (gRNA Uridine Insertion/Deletion Editor), was created. This program is able to analyze the functionality of gRNA clusters generated by ACORNS, by simulating the editing process. Combining this data with the mRNA transcriptomes

previously generated, we conducted a detailed analysis of each population of gRNAs capable of editing these three genes. This analysis revealed a surprisingly high tolerance for mismatches and gaps in mRNA/gRNA alignments in the editing system, most notably in the editing of the essential RPS12 [100].

This project found that not only is alternative editing present in *T. brucei*, but that it is pervasive, and the system, as a whole, is surprisingly robust. We propose the hypothesis that the RNA editing system does in fact promote the genetic robustness of *T. brucei* through the facilitation of dual-coding genes, as well as the introduction of alternative edits that increase protein diversity and allow the editing system to continue to evolve.

## **CHAPTER 2: ANALYSIS OF THE *TRYPANOSOMA BRUCEI* EATRO**

### **164 BLOODSTREAM GUIDE RNA TRANSCRIPTOME**

#### **Abstract**

The mitochondrial genome of *Trypanosoma brucei* contains many cryptogenes that must be extensively edited following transcription. The RNA editing process is directed by guide RNAs (gRNAs) that encode the information for the specific insertion and deletion of uridylates required to generate translatable mRNAs. We have deep sequenced the gRNA transcriptome from the bloodstream form of the EATRO 164 cell line. Using conventionally accepted fully edited mRNA sequences, ~1 million gRNAs were identified. In contrast, over 3 million reads were identified in our insect stage gRNA transcriptome. A comparison of the two life cycle transcriptomes show an overall ratio of procyclic to bloodstream gRNA reads of 3.5:1. This ratio varies significantly by gene and by gRNA populations within genes. The variation in the abundance of the initiating gRNAs for each gene, however, displays a trend that correlates with the developmental pattern of edited gene expression. A comparison of related major classes from each transcriptome revealed a median value of ten single nucleotide variations per gRNA. Nucleotide variations were much less likely to occur in the consecutive Watson-Crick anchor region, indicating a very strong bias against G:U base pairs in this region. This work indicates that gRNAs are expressed during both life cycle stages, and that differential editing patterns observed for the different mitochondrial mRNA transcripts are not due to the presence or absence of gRNAs. However, the abundance of certain gRNAs may be important in the developmental regulation of RNA editing.

## **Author Summary**

*Trypanosoma brucei* is the causative agent of African sleeping sickness, a disease that threatens millions of people in sub-Saharan Africa. During its life cycle, *Trypanosoma brucei* lives in either its mammalian host or its insect vector. These environments are very different, and the transition between these environments is accompanied by changes in parasite energy metabolism, including distinct changes in mitochondrial gene expression. In trypanosomes, mitochondrial gene expression involves a unique RNA editing process, where U-residues are inserted or deleted to generate the mRNA's protein code. The editing process is directed by a set of small RNAs called guide RNAs. Our lab has previously deep sequenced the gRNA transcriptome of the insect stage of *T. brucei*. In this paper, we present the gRNA transcriptome of the bloodstream stage. Our comparison of these two transcriptomes indicates that most gRNAs are present in both life cycle stages, even though utilization of the gRNAs differs greatly during the two life-cycle stages. These data provide unique insight into how RNA systems may allow for rapid adaptation to different environments and energy utilization requirements.

## **Introduction**

The life cycle of *Trypanosoma brucei* involves two distinct environments, the animal host and the insect vector. These environments are distinct in temperature and nutrient composition, providing a unique challenge to *T. brucei* as it cycles between hosts. In the bloodstream, trypanosomes exist in two forms, the actively dividing slender form and the non-dividing stumpy form. The slender form is optimized to utilize its glucose rich environment, using glycolysis to generate energy [14]. The stumpy form appears to be transitional, activating

mitochondrial genes in preparation for uptake in a blood meal by its tsetse fly vector and subsequent transfer to a harsher environment [14]. Once inside the tsetse fly, the parasite utilizes proline to drive oxidative phosphorylation and ATP production in the mitochondrion [19]. While the activity of the mitochondrion is relatively low during the bloodstream stage (BS), expression of the mitochondrial genome is still essential [17,18]. In *T. brucei*, the mitochondrial genome consists of two types of DNA molecules, maxicircles and minicircles. Maxicircles are 22kb circular DNA that contain the genes for two ribosomal RNAs, 12S and 9S, and eighteen mRNA genes [5]. While some of the protein-coding genes do not require RNA editing prior to translation, most require extensive editing before they can be translated [6,7]. This process involves the insertion of hundreds of uridylates (U)s and less frequently deletion of Us, often doubling the size of the transcript. The sequence changes are guided by small complementary RNA molecules (the guide RNAs) that are encoded on the minicircles [8]. Minicircles make up the bulk of the kinetoplastid network (anywhere from 5,000–10,000 present in each network) with each minicircle encoding 3–5 gRNAs. In *T. brucei*, there are more than 200 different minicircle sequence classes (~1200 gRNAs) [8].

Distinct differences in mitochondrial transcript abundance, polyadenylation and the extent of RNA editing are observed during the complex life cycle (Table 1). The pattern of differential RNA editing observed is especially interesting. For example, the cytochrome b (CYb) and cytochrome oxidase II (COII) mRNAs are edited during the insect stage, but are primarily unedited in bloodstream forms [24,25]. In contrast, editing of the NADH dehydrogenase subunit transcripts (ND3, ND7, ND8 and ND9) and editing of the ribosomal protein subunit 12 transcript (RPS12) appears to occur preferentially in bloodstream forms [5,26–30]. Other

transcripts, cytochrome oxidase III (COIII) and ATPase subunit 6 (A6) are edited in both life cycle stages [31,32]. Early studies using both Northern blot and primer extension analyses on a limited number of gRNAs indicate that gRNAs are present in both insect and bloodstream forms, suggesting that the regulation of RNA editing is not at the level of gRNA availability [28,50,51]. Our lab has previously published deep sequencing results of the gRNA transcriptome of the *T. brucei* EATRO 164 procyclic form [9]. Here we present the deep sequencing data for the gRNA transcriptome of a bloodstream form of EATRO 164. A total of 211 populations of gRNAs were identified. We define a population as a group of gRNAs that may vary in sequence, but direct the editing of the same or near same region of the mRNA. Because kinetoplastid RNA editing allows G:U base pairing, most populations contain multiple sequence classes that can guide the generation of the same mRNA sequence. While the number of populations identified was similar to the number identified in the procyclic gRNA transcriptome (214 populations), the total number of gRNAs identified was much reduced and the coverage was less complete; full complements of gRNAs were only identified for COIII and CYb. In spite of the reduced number of gRNAs, an interesting correlation was found that suggests a relationship between the relative abundance of initiating gRNAs between stages and the developmental pattern of mRNA editing.

## **Materials and Methods**

### **Parasites, isolation of mitochondria and RNA extraction**

*T. brucei brucei* clone IsTar from stock EATRO 164 were grown in rats and isolated as previously described [101]. Bloodstream forms were virtually all long-slender forms isolated after 4 days of infection. Parasites were used immediately for isolation of mitochondria using

differential centrifugation as previously described or stored frozen at -80°C until RNA extraction [9]. Both total RNA from whole parasites and mitochondrial RNA (mtRNA) from purified mitochondria were isolated by the acid guanidinium-phenol-chloroform method [102].

### **Ethics statement**

Rats were raised according to the animal husbandry guidelines established by Michigan State University. All vertebrate animal use procedures were approved by MSU's Institutional Animal Care and Use Committee (Application 03/11-051-00). MSU has filed with the Office of Laboratory Animal Welfare (OLAW) an assurance document that commits the university to compliance with NIH policy and the Guide for the Care and Use of laboratory Animals.

### **Library preparation and Illumina sequencing**

Samples of mtRNA and total RNA were both treated with DNase RQI and size fractioned on a polyacrylamide gel as previously described [20]. Guide RNAs were extracted from the gel and prepped for sequencing using the Illumina 'Small RNA' protocol as previously described [9]. Libraries from both mtRNA and total RNA samples were deep sequenced on Illumina GAIix. Reads were then processed and trimmed as previously described [9]. Data with two or more Ns, shorter than 20nts after trimming or with an overall mean Q-score < 25 were discarded. Redundant reads were then removed, while maintaining the number of redundant reads and reads containing fewer than 4 consecutive Ts were removed.

### **Identification of gRNAs**

To identify gRNAs, each transcript read was aligned to the conventionally edited mRNAs based on known base pairing rules (canonical Watson-Crick base pairs and the G-U base pair). In the initial screen, no gaps were allowed in the alignment, allowing the formulation of the

gRNA-mRNA alignment as an extended longest common substring (LCS) problem as previously described [25]. Matched gRNAs were then scored (two points for G:C and A:U base pairs and one point for G:U base pairs). gRNAs with scores >45 were identified as guiding a specific region based on the identified mRNA fully edited sequence. Additional searches with reduced stringency (scores >30) were performed on regions with low gRNA coverage. The matched gRNAs were sorted into populations based on their guiding positions, and the populations analyzed and sorted into major sequence classes.

## **Results**

Much of the initial characterization of RNA editing in *T. brucei* was done using the EATRO 164 strain. These experiments suggested that RNA editing was developmentally regulated in that certain genes were shown to be more fully edited in some stages than others (Table 1) [24–32,46]. It was also reported that the developmental regulation was not controlled by gRNA availability, as gRNAs were found in both life cycle stages [28,50,51]. In these early studies, however, only a small number of gRNAs were investigated. In this study, we used deep sequencing to compare the gRNA transcriptomes of a bloodstream form to a procyclic form of *T. brucei* EATRO 164. The EATRO 164 strain was isolated in 1960 from *Alcephalus lichtensteini* and maintained in the lab of Dr. K. Vickerman until being obtained by Dr. Stuart in 1966 [103]. Dr. Stuart derived the procyclic form from the Bloodstream culture in 1979 [103]. Both cell lines have been maintained in separate culture since that time.

Trypanosomes from the EATRO 164 strain were grown in Wistar rats to a parasitemia of 1–2 × 10<sup>9</sup> trypanosomes per mL and isolated using DEAE cellulose columns. Mitochondria and gRNAs were purified as previously described [9]. Libraries were generated using gRNAs isolated

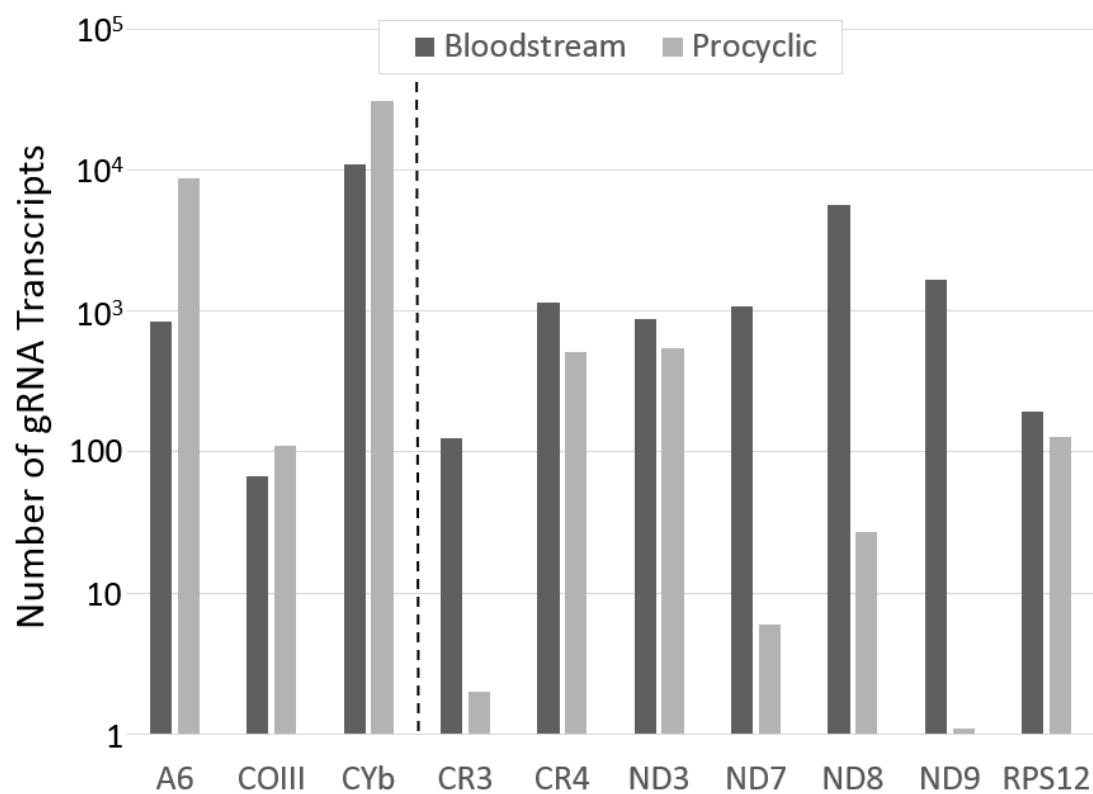
from whole cell RNA and gRNAs isolated from mitochondrial RNA. Both bloodstream gRNA libraries were searched using conventionally accepted fully edited mRNA sequences, and a total of 1,024,604 gRNA reads were identified. Surprisingly, the library generated using gRNAs isolated from whole cell RNA had more than twice as many identified gRNA reads as the data generated using gRNAs isolated from mitochondrial RNA. To insure sufficient abundance and gRNA coverage, the two data sets were combined for the analyses presented here. In contrast, over 3 million gRNA reads were identified in our procyclic gRNA transcriptome generated from gRNAs isolated from mitochondrial RNA. Of the 1,024,604 reads identified from the bloodstream transcriptomes, 982,450 reads were sorted into major sequence classes.

The overall ratio of identified procyclic gRNA reads to BS gRNA reads was 3.5:1. This ratio varies significantly by gene (Table 2), and by populations within genes (APPENDIX A) and, except for the initiating gRNA, no apparent trend relating gRNA abundance and developmental editing pattern was observed. Interestingly, for the initiating gRNA, mRNAs that are fully edited in the procyclic stage only, or are fully edited in both life cycle stages had initiating gRNAs with more reads in the procyclic data set (Figure 1) [31,32,40,46]. In contrast, mRNAs that are only fully edited or are more abundant in the BS, had more initiating gRNAs reads in the BS data set (Figure 1) [26–30].

Because the identified gRNAs from the BS cells were less abundant, the rule used to identify major gRNA sequence classes was relaxed. Instead of using a strict cut off for the minimum number of reads required, the cut off was assessed on a case-by-case basis. For example, if the total population only had 100 reads, a sequence class with only 10 reads would still be identified as a major sequence class. Once all major classes were identified, 657

**Table 2. Number of gRNA transcripts in procyclic and bloodstream major classes and ratio of procyclic transcripts to bloodstream transcripts for each gene.**

Gene	Bloodstream gRNA Reads	Procyclic gRNA Reads	Ratio of PC to BS Reads
A6	41,628	266,532	6.40
COIII	371,139	948,845	2.56
CR3	13,316	236,808	17.78
CR4	25,753	51,979	2.02
CYb	11,022	31,622	2.87
MurfII	157	2,605	16.59
ND3	13,567	75,739	5.58
ND7	291,927	702,061	2.40
ND8	112,868	584,639	5.18
ND9	83,924	72,027	0.86
RPS12	17,191	403,131	23.45
Total	982,492	3,375,988	3.44



**Figure 1. The abundance of the initiating gRNA of all edited mRNAs in each stage.** mRNAs to the left of the dashed line are constitutively edited or are edited only in the procyclic stage [31,32,40,46]. mRNAs to the right of the dashed line are only fully edited or more abundant fully edited in the bloodstream stage [26–30].

sequence classes were identified that could be sorted into 211 populations (Table 3).

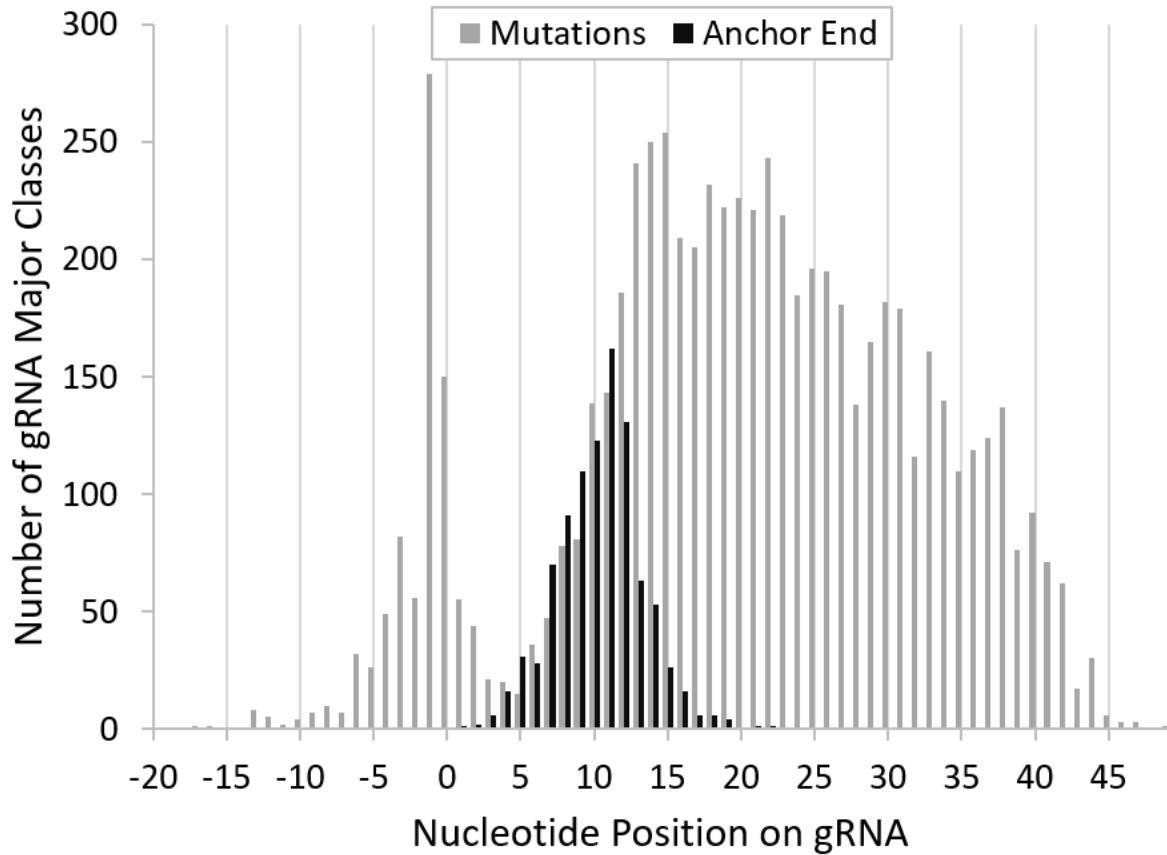
Although the overall gRNA numbers were down in comparison to the procyclic data set, most of the populations found in that stage (214 gRNA populations) were also identified in the BS transcriptome. However, there were a number of populations that were unique to either the procyclic or BS stage.

**Table 3. Summary of the gRNA data coverage for each gene.**

Gene	Populations		Unique Populations		Average* gRNA Overlap (nts)		Gaps		Weak Overlaps	
	BS	PC	BS	PC	BS	PC	BS	PC	BS	PC
A6	29	28	1	0	18	20	1	0	0	0
COIII	42	39	4	1	19	22	0	0	0	0
CR3	9	9	0	0	19	14	1	0	0	2
CR4	16	18	0	2	17	18	2	0	0	0
CYb	2	2	0	0	12	14	0	0	0	0
MurfII	1	1	0	0	N.A.	N.A.	1	1	0	0
ND3	12	12	0	0	15	15	1	1	0	0
ND7	45	48	2	5	17	21	7	2	4	1
ND8	20	21	2	3	17	21	2	1	1	0
ND9	24	23	1	0	16	16	1	0	0	2
RPS12	11	13	0	2	17	21	2	0	1	0
Total	211	213	10	13	17	19	18	5	6	5

<sup>a</sup>The average gRNA overlaps were determined excluding any regions where neighboring gRNAs shared no overlap.

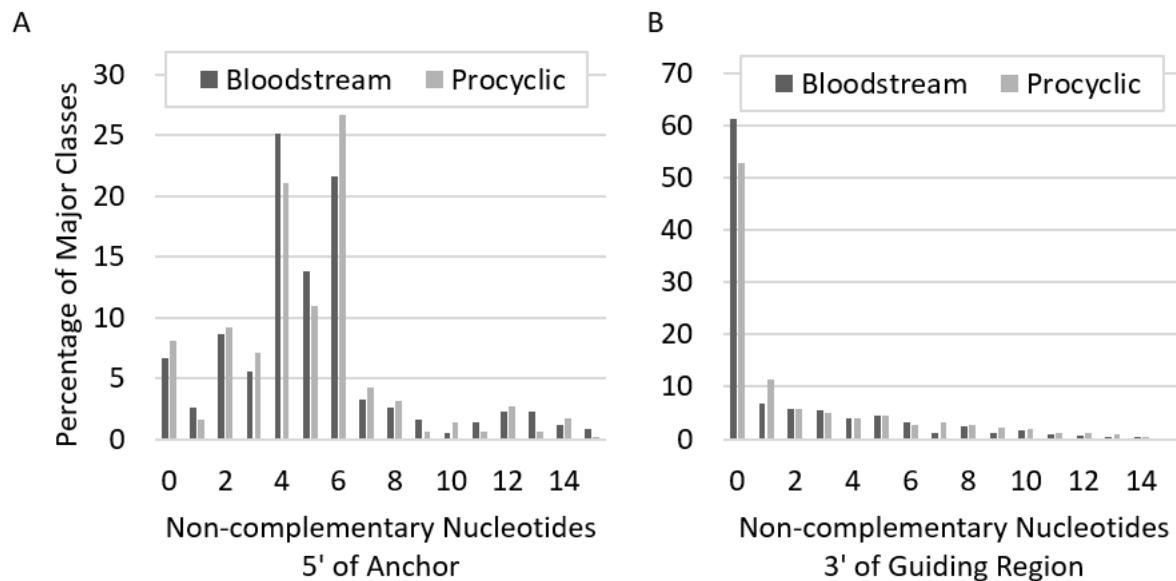
Surprisingly, when the bloodstream and procyclic data sets were compared, only 37 identical major sequence classes were found in both. However, distinctly related sequence classes could be identified when comparing the BS and procyclic populations. Comparing the related major classes from each transcriptome (BS vs procyclic) revealed a median value of ten single nucleotide variations per gRNA. Interestingly, nt variations were much less likely to occur in the consecutive Watson-Crick anchor region of the gRNA than in the rest of the gRNA indicating a very strong bias against G:U base pairs in this region (Figure 2).



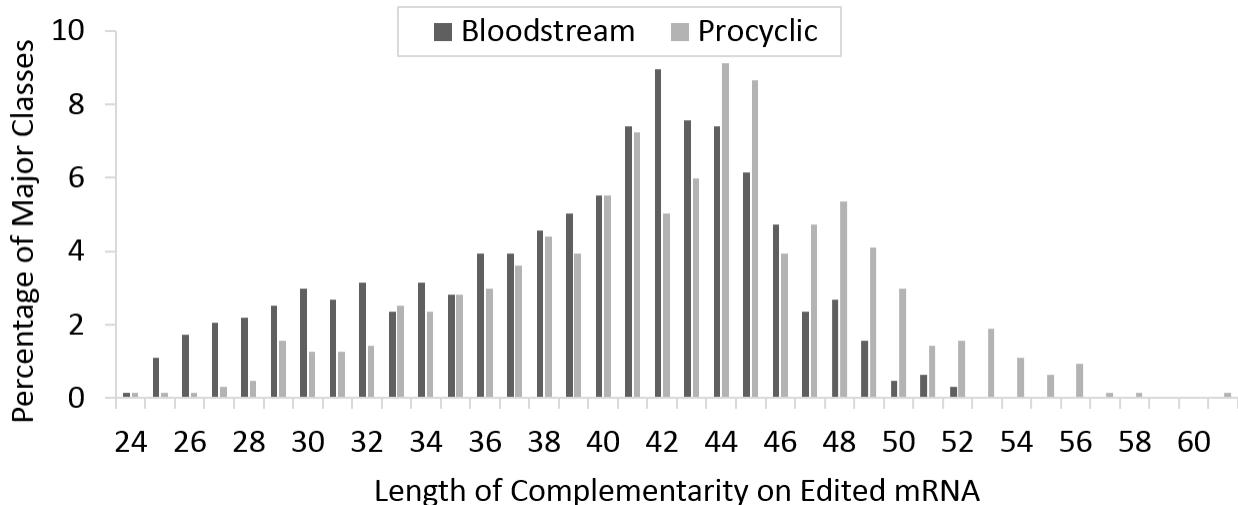
**Figure 2. The frequency of nt variations versus nucleotide position in the gRNA.** Gray bars indicate the number of gRNA sequence classes with an identified nt difference between related procyclic and bloodstream gRNAs. Nucleotide numbering for each gRNA was normalized by setting the start of the Watson-Crick anchor region to zero. Black bars indicate the number of gRNA sequence classes whose contiguous Watson-Crick anchors end at that position (start of Watson-Crick = zero, so this is an indication of the length of the contiguous Watson-Crick region).

The Watson-Crick anchors (defined as the number of consecutive nts in the 5' region with only G:C and A:U base pairs) had a median length of eleven nucleotides and anchor length did not vary between the two forms. The vast majority of major classes of gRNAs had consecutive Watson-Crick anchors greater than seven nts long (92.5%). In addition, most gRNAs with Watson-Crick anchors shorter than eight nts were not an abundant major class for their respective populations. Consistent with observations made from the procyclic data set, most gRNAs had zero non-base pairing nucleotides 5' to the poly-uridine tail and 4 to 6 non-base pairing nucleotides 5' to the anchor region (Figure 3). Also consistent with procyclic data, most

of the gRNAs (59%) had 38 to 48 nts of complementarity (including anchor regions) with their respective mRNAs (Figure 4). Transcription start sites also did not vary, as preference for an RYAYA start site was observed (Table 4).



**Figure 3. Comparing the number of non-complementary nucleotides 5' of the anchoring region (A) or 3' of the guiding region (excluding the U-tail) (B) in procyclic and bloodstream gRNAs.**



**Figure 4. Length of gRNA complementarity (including anchors) to fully edited mRNAs for both bloodstream and procyclic gRNAs.**

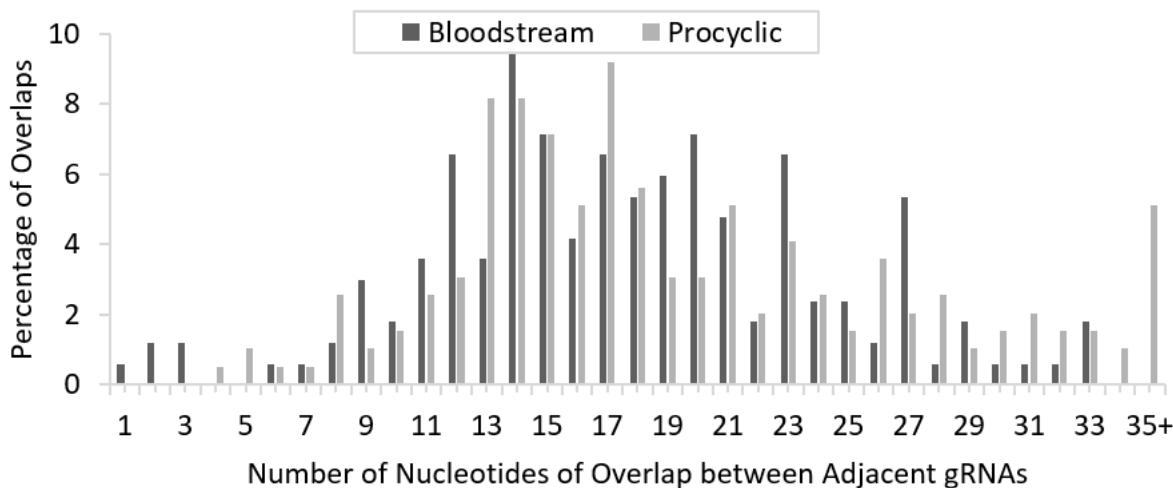
**Table 4. Most common gRNA transcription start sites in procyclic and bloodstream data.**

Initiating Sequence	Stage	Percentage of Sequence Classes	Percentage of Transcripts
ATATAT	Bloodstream	32.20%	33.60%
	Procyclic	35.20%	37.40%
ATATAA	Bloodstream	20.00%	17.90%
	Procyclic	21.10%	24.70%
AAAAAA	Bloodstream	3.60%	1.60%
	Procyclic	4.60%	1.30%
ATATAC	Bloodstream	3.80%	1.50%
	Procyclic	4.30%	3.80%
ATACAA	Bloodstream	2.40%	1.30%
	Procyclic	2.80%	1.70%
ATATTA	Bloodstream	4.90%	13.40%
	Procyclic	2.60%	0.90%
ATATAG	Bloodstream	2.20%	0.70%
	Procyclic	2.60%	7.60%
ATAAAT	Bloodstream	2.50%	1.00%
	Procyclic	2.60%	0.70%
ATACAT	Bloodstream	2.70%	3.20%
	Procyclic	2.20%	2.20%
ATAAAA	Bloodstream	1.70%	0.20%
	Procyclic	2.10%	1.10%

## Coverage and gaps

In order to determine if the BS gRNA transcriptome contained a full complement of guide RNAs, the gRNA populations were aligned to the fully edited mRNAs (APPENDIX B). We note, that for an mRNA to be fully edited, not only must all editing sites on the mRNA be covered by a gRNA, the downstream gRNA must generate the anchor binding site for the subsequent gRNA. Therefore, adjacent gRNAs must overlap. Overall, there was an average of 17 nts of overlap between adjacent gRNAs, with the average overlap varying slightly by gene (Table 3). As the median Watson-Crick Anchor is 11 nts, in most cases, the overlap extends beyond the Watson-Crick anchor of the subsequent gRNA. However, we did observe a number

of regions where the overlap is minimal. Currently, there is no data that stipulates the minimum anchor needed for efficient editing. However, we postulate that similar to microRNAs, for an anchoring sequence to be sufficiently specific, it should be at least six nucleotides [104]. Indeed, when examining the overlaps between most gRNAs, there are only ten (four procyclic and six BS) that are less than six nucleotides (Figure 5).



**Figure 5. The percentage of different nucleotide overlaps found between adjacent gRNAs.** gRNAs were aligned to their fully edited mRNA sequence and the number of mRNA nts with complementarity to both adjacent gRNAs determined.

We therefore used six nucleotides as a cut off to identify regions with potential missing guide RNAs for both life cycle stage transcriptomes. In contrast to the procyclic data, where full complements of gRNAs were identified for five of the mRNA transcripts (A6, COIII, CR4, CYb, and RPS12), in the BS transcriptome, a full complement of gRNAs was only identified for COIII and CYb. Overall, there are 12 edited regions where no gRNAs were identified, and five regions with weak gRNA overlaps in the BS data (Table 5). Of these 17 regions, seven belong to ND7 alone. Interestingly, nine of the 17 missing populations are in very low abundance in the procyclic data, having 100 or fewer reads. Because the number of reads in the BS data is ~3.5 fold less abundant, this could account for some of these regions of poor coverage. There are six

regions that lack gRNA coverage in both data sets. These are found in CR3, MurfII, ND3 and ND7 (Table 5). Interestingly, three of these regions are close to the 3' end of their respective genes. Regions of weak overlap (ND9(238–242), ND9(609–612)) and regions without gRNA coverage (CR3(278–292), ND8(541–553)) that are unique to the procyclic transcriptome were also observed. Interestingly, the regions of poor procyclic coverage are found in CR3, ND8 and ND9, all transcripts that are preferentially edited in the BS form [5,28,29,47].

**Table 5. Identified gaps or weak overlaps (less than 6 nucleotides) between populations of gRNAs observed in both data sets.**

Gene	Stage Missing Coverage	Range	Gap or Overlap	Abundance of Equivalent gRNA
A6	Bloodstream	669-670	2 nt Gap	39,063
CR3	BS/P <sup>a</sup>	233/226-230	1 nt G/5 nt O	Missing in Both Stages
CR3	Procyclic	278-292	15 nt Gap	125
CR4	Bloodstream	143-165	23 nt Gap	7,175
CR4	Bloodstream	302-306	5 nt Gap	643
MurfII	BS/P	80-85	6 nt Gap	Missing in Both Stages
ND3	BS/P	389-401	13 nt Gap	Missing in Both Stages
ND7	BS/P	92-94	3 nt Gap	Missing in Both Stages
ND7	Bloodstream	95-120	26 nt Gap	1
ND7	Bloodstream	292-293	2 nt Overlap	3,259
ND7	Bloodstream	325-326	2 nt Gap	888
ND7	BS/P	485-486	0 nt Overlap	Missing in Both Stages
ND7	Bloodstream	1000-1000	1 nt Overlap	101
ND7	Bloodstream	1079-1085	7 nt Gap	44
ND7	BS/P	1086/1086-1088	1 nt G/3 nt G	Missing in Both Stages
ND7	Bloodstream	1225-1232	8 nt Gap	123
ND7	Bloodstream	1269-1270	2 nt Overlap	63
ND8	Bloodstream	54-56	3 nt Overlap	1
ND8	Bloodstream	153-159	7 nt Gap	4
ND8	Bloodstream	386-389	4 nt Gap	2
ND8	Procyclic	541-553	13 nt Gap	413
ND9	Procyclic	238-242	5 nt Overlap	652
ND9	Bloodstream	340-342	3 nt Gap	36
ND9	Procyclic	609-612	4 nt Overlap	7
RPS12	Bloodstream	122-132	11 nt Gap	3
RPS12	Bloodstream	156-158	3 nt Overlap	62
RPS12	Bloodstream	337-349	13 nt Gap	128

**Table 6. Summary of populations found in both data sets that have more reads in the bloodstream data set than in the procyclic data set.**

Gene	PC and BS Shared Populations	Populations more abundant in BS	Percentage of populations more abundant in BS
A6	28	6	21%
COIII	38	12	32%
CR3	9	5	56%
CR4	16	10	63%
CYb	2	0	0%
MurfII	1	0	0%
ND3	12	6	50%
ND7	43	20	47%
ND8	18	8	44%
ND9	23	16	70%
RPS12	11	4	36%
Total	201	87	43%

### Gene specific gRNA characteristics

**ATPase 6.** In the BS gRNA transcriptome, a total of 29 gRNA populations containing 86 different major sequence classes were identified that could guide the editing of A6 (Table 3; APPENDIX C Part A). One population was identified that was unique to the BS transcriptome (gA6(281-329)). The gRNAs bordering this population share extensive overlap, so its absence in the procyclic transcriptome would not impact the editing process (APPENDIX C Part A). We note that two of the gRNAs identified have single nucleotide mismatches. The bloodstream gA6(640–668) has an identified mismatch (C:U) that disrupts the complementarity of the gA6(640–668) population (APPENDIX B Part A). The second mismatched gRNA (gA6(520–533)) would introduce a frameshift. Excluding these two mismatched regions, there is complete coverage of ATPase 6. In contrast to the procyclic data, where the conventional initiating gRNA and the gRNA immediately following it were extremely rare, both of these gRNAs, gA6(773–822), previously identified as gA6-14 and gA6(745–789), were fairly abundant, each having hundreds

of reads. The alternative initiating gRNA identified in the procyclic data set was not found. This finding is similar to that found in the *T. brucei* Lister strain 427 where authors identified alternative initiating gRNAs not found in the EATRO 164 procyclic gRNA transcriptome [105].

Another disparity between the two life cycle data sets was found when comparing the abundance of gRNAs implicated in a potential alternative edit. In the procyclic gRNA transcriptome, a gRNA was identified that would guide the insertion of 11 U-residues instead of the needed 12 between G555 and A568 [9]. This gRNA (pA6(557–593)) was 25-fold more abundant than the conventional gRNA (pA6(549–593)). In the BS data set however, more than 400 reads of the 12U gRNA were identified and only one read was found that would encode the alternative 11U edit. Surprisingly, while G555-A568 would be correctly edited (insertion of 12Us), the next editing site (A549-G555) is edited by bsA6(520–553), the gRNA that introduces the 1 nt frameshift. This frameshift would generate a predicted protein with nearly the same amino acid sequence as the procyclic 11U frameshift edit (two amino acid changes) (Figure 6).

160	170	180	190	200	210	220
Conv	-FFDFYFIEVFFFYGVFCYWFILFIFVFCFCLLFYVFLYLLDLFAAILQLFIFCNMILQLIMDFLLFLLFV					
11U	-FFDFYFIEVFFF <u>MVFFVIDLFYLFLCFV</u> VYYFMCFYICWIYLPPYYSYLFFMWFCSW					
4U	-FFDFYF <u>I</u> ILFFF <u>MVFFVIDLFYLFLCFV</u> VYYFMCFYICWIYLPPYYSYLFFMWFCSW					

**Figure 6. Alignment of conventional ATPase 6 protein sequence to hypothetical proteins generated by the 11U alternative edited mRNA and the 4U alternatively edited mRNA.** Double underlined residues show where the alternative sequences differ from the conventional sequence. The shaded residues in the 4U sequence show where it differs from the 11U sequence.

**Cytochrome oxidase subunit III.** Forty-two gRNA populations, guiding the editing of COIII were identified in the BS transcriptome; three more than in the procyclic data set (APPENDIX C Part B). This disparity is caused by the presence of several unique populations. While the procyclic data set contained one unique population, the BS data contained four gRNA

populations not previously identified. Of these four unique populations, three of them are required for full overlapping coverage in the bloodstream. They are not however, required for full coverage in the procyclic stage. These three unique gRNA populations all span relatively small regions of weak overlap (Figure 7, APPENDIX B Part B).

420            430            440            450            460            470            480

GuuuGCUUUCGUuuuuuuGuuuACCuuAuAuGuuuuGuuuAuuAuGuGAuuAuGGGuuuuGuuuuuuAu  
 |||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
 pCO3 (461-497) 14 TACTTATAGTGTACTGATGTTAAGACAGAAGATA  
 |||:||||:||||:||||:||||:||||:||||:||||:  
 bscO3 (457-499) 14 TATAAGTAGTATATTAGTGTCAAGAGTAAAAGATA  
 :||||:||||:||||:||||:||||:||||:||||:  
 TAAATGGAAGTGAAGAATAGATGGAATGTGTAAACAAACAATAATATCATA 5' pCO3 (418-467)  
 |||:||||:||||:||||:||||:  
 TAAATGTAGGTAAAAAGTGAATGGAATATGCAAAACAACAATTATA 5' bscO3 (427-460)  
 |||:||||:||||:||||:  
 10 TAAAAGTAGATAGAATATGCAAGATAGTAAATAACTAATA 5' bscO3 (443-474)

**Figure 7. Editing sites 420–489 of COIII aligned with the gRNAs identified for that region in the procyclic (grey) and bloodstream (black) data sets.** The gRNA covering 443–474 was only found in the bloodstream data set.

An alternative edit of COIII has been described, involving distinct edits at two adjacent sites that links the open reading frame of the edited 3' end to an ORF found in the 5' pre-edited sequence [67]. The previously identified alternative gRNA that can generate the needed editing events was not found in either the BS or procyclic transcriptomes.

**C-rich regions 3 and 4.** In the BS data set, nine populations and 34 major sequence classes were identified that direct the editing of the CR3 transcript. The coverage of edited CR3 is nearly complete in the bloodstream data set with only a one nt gap in coverage (editing site 233) (APPENDIX B Part C). This is in contrast to the procyclic transcriptome, where gRNAs that matched the published sequence downstream of nt 196 were very rare (<10 copies) and no gRNAs were identified that could direct editing near the 3' end (nucleotides 275–292).

A full consensus sequence for edited CR4 has only been found in BS *T. brucei* [48]. Using this sequence, 16 gRNA populations, containing 62 major sequence classes were identified in

the BS transcriptome (APPENDIX C Part D). In contrast to the procyclic data, where a full complement of gRNAs were identified, there are two gaps in the BS coverage (Table 5).

**Cytochrome b and maxicircle unidentified reading frame II.** RNA editing in the Cytochrome b (CYb) transcript is limited to the 5' end and two gRNA populations are sufficient to guide the small number of edits needed to render the CYb transcript functional. Both populations were observed in both data sets, with a total of 6 major classes. Interestingly, in both data sets, the initiating gRNA is significantly more abundant than the second gRNA, being approximately 30 fold more abundant in the procyclic data set and approximately 200 fold more abundant in the bloodstream data set (APPENDIX C Part E). This is in contrast to most of the other transcripts where the initiating gRNAs are not very abundant. In addition, almost all of the CYb gRNA major classes have an A-run transcription start site, deviating from the common RYAYA initiation site pattern.

Editing in MurfII is also limited to the 5' end and requires only two gRNAs. One of these gRNAs (gMurfII(30–79)) is encoded on the maxicircle [106]. While this gRNA was observed in both data sets, the gRNAs identified were not identical. A purine-purine transition near the 3' end of the gRNA differentiates the procyclic and BS forms (APPENDICES B and C Parts F). An initiating gRNA is needed to generate the 3' most edits that create the anchor sequence for gMurfII (30–79). This gRNA was not found in either data set, despite additional searches with reduced search stringency.

**NADH dehydrogenase subunits 3, 7, 8, and 9.** In the initial characterization of RNA editing in *T. brucei* EATRO 164, fully edited ND subunit transcripts were only found in RNA isolated from the BS stage. We were therefore surprised to find that fewer ND gRNA

populations were identified in the BS transcriptome and a full complement of gRNAs was not identified for any of the ND subunits. The most complete coverage was found for ND3 and ND9. For ND3, the BS data set contained twelve populations and 41 major classes of gRNAs. One gap in coverage was observed, from 389–401. This region overlaps a region that has no clear consensus sequence, 375–395 [26]. ND9 is the only gene in this study whose bloodstream gRNA reads outnumber the procyclic gRNA reads identified (Table 2). Twenty-four bloodstream gRNA populations were identified with all edited nucleotides covered if gRNAs with a single base pair mismatch are taken into account (APPENDIX B Part J).

While 45 gRNA populations were identified for ND7 in the BS data set, the gRNA coverage was significantly worse when compared to the identified procyclic gRNAs (Table 5). Despite the poor coverage, two unique gRNA populations (bs gRNA (772–816) and (1128–1182)) were identified (APPENDICES B and C Parts H). ND8 also had poor gRNA coverage (Table 5). Interestingly, there are several populations in ND8 that contain highly abundant gRNA sequence classes with mismatches that shorten the complementarity of the gRNA. These usually have a single mismatch in the gRNA that would otherwise guide conventional editing (APPENDIX C Part I).

**Ribosomal protein S12.** The BS data set contained 11 populations and 26 major sequence classes that direct editing of RPS12 (Table 3). While the procyclic transcriptome contained a full complement of gRNAs, the BS RPS12 data contains one gap in coverage and one region of poor overlap, (Table 5). This was surprising, as RPS12 has been shown to be essential in both life cycle stages [100,107]. The region of the mRNA with poor coverage has a high percentage of C residues and gRNAs covering this region may utilize C:A base pairs. If this

is the case, some classes of gRNAs may not have been detected, as the program used to search for gRNAs does not allow for C:A base pairs (APPENDIX B Part K).

## **Discussion**

This is the first comprehensive characterization of the mitochondrial gRNA transcriptome from the bloodstream stage of *Trypanosoma brucei brucei*. As we have previously characterized the insect stage gRNA transcriptome, these data allow the comparison of gRNA characteristics across the two main life cycle stages [9]. In the EATRO 164 BS gRNA transcriptome, gRNAs for every edited gene were identified. Interestingly, while the number of populations identified in this data set was only slightly lower than that reported in the procyclic data set, the total number of gRNA transcript reads identified was considerably lower despite the fact that multiple transcriptome libraries were combined. While this may be a reflection of the down regulation of mitochondrial transcription in the bloodstream stage (see Table 1), it is impossible to rule out technical problems in the generation and sequencing of the libraries. It has been previously reported that gRNA presence did not correlate with developmental RNA editing patterns in *T. brucei* and our data does not challenge this [50,51]. The data did however, show an interesting trend in the abundance of the initiating gRNAs as relates to their developmental editing patterns. It may be that the abundance of the initiating gRNAs is regulated in order to control editing of their target mRNAs. However, we cannot rule out the possibility that not all of the populations of initiating gRNAs were identified. For the pan-edited mRNAs, the initiating gRNAs direct sequence changes that are often downstream of the stop codon. Sequence changes in this region would be tolerated, as long as the anchor sequence for the next gRNA is maintained. This type of mutation was observed in the 3' end of ATPase 6 [19].

In addition, characterization of the initiating gRNAs in the Lister 427 *T. brucei* cell line identified several gRNAs that would direct an alternative editing pattern, suggesting a high tolerance for sequence changes near the mRNA 3' ends. [105].

As expected, general gRNA characteristics are conserved across the two life-cycle stages. Populations retain the general location of their anchors, there is relatively little shift in the location of populations, and the lengths of complementarity are very similar. We did observe that considerable nucleotide variations were found in the guiding regions of the gRNAs from the different life cycle strains of the EATRO 164 cells. This particular cell line dates back to 1960 when the BS form was originally acquired [103]. Procyclic cells were derived from the BS stock in 1979 and the two cell lines maintained separately since that date [103]. Mixed trypanosome genotypes are detected frequently in field isolates from both tsetse flies and mammals and it may be that separation into different culture conditions allowed different genotypes to predominate in each life cycle strain [22,108,109]. Because gRNAs utilize both canonical (Watson-Crick) as well as G:U base-pairing to direct the change in sequence, most transition mutations in the gRNA, would not lead to changes in the mRNA sequence and would not be selected against [33]. We do note however, that a very strong bias against A to G transitions is observed in the anchor regions of the gRNAs. This suggests that transition mutations in this region are not tolerated. This suggests that the editing machinery recognizes and selects for a conventional base-paired double helix in the initial gRNA/mRNA pairing. The ability to discriminate against G:U base-pairs in the initial interaction would greatly increase the accuracy of the gRNA targeting event. Considering the sequential nature of the overall editing process, this would be very advantageous.

## Coverage

Surprisingly, complete gRNA coverage was observed only for the pan-edited COIII and for CYb, where editing is limited to the 5' end. The identification of the CYb gRNAs was expected, as it has been previously reported that the gRNAs are present in both life cycle stages even though editing of CYb is limited to the procyclic stage [8,24]. The full coverage of COIII was also not surprising, as COIII was shown to be fully edited and equally abundant in both stages [32]. However, we expected to see complete coverage of ATPase 6 and RPS12 as both of these transcripts have been shown to be essential in both life cycle stages [17,100,107,110]. For ATPase 6, we did identify a total of 29 gRNA populations that do cover all of the editing sites. However, one of the gRNAs (bsA6(643–667)) has a single nucleotide mismatch (C:U) and one would introduce a frameshift (bsA6(520–553)). The C:U mismatch occurs near the middle of the gRNA, placing the C:U mismatch in a region that is unusually high in Gs and Cs (APPENDIX B Part A). It may be that the G:C base pairs immediately upstream of the mismatch stabilize the gRNA/mRNA interaction, allowing it to be tolerated. The frameshift gRNA is also interesting, as it occurs just upstream (1 editing site) of another site where we had previously observed a frameshift sequence anomaly. Both frameshifts (the BS 4U and the Procyclic 11U) generate a predicted protein with nearly the same amino acid sequence. As the frameshifts occur downstream of the highly conserved amino acid region involved in proton translocation [31], it may be that this different carboxyl terminus is tolerated.

Near full coverage is also observed for RPS12. For this transcript, one BS identified gRNA (bsRPS12(96–121)) has an A-nt insertion that disrupts the gRNA complementarity. Surprisingly, the other mRNA transcript found with near complete coverage was ND9 (one gRNA has a single

nt mismatch). All of the other mitochondrially encoded Complex I members did have substantial gaps in coverage. Currently, there is considerable debate on the necessity of Complex I subunits for either stage of the trypanosome life cycle. Studies using RNAi and knockout cell lines of nuclear-encoded members of Complex I have shown that the complex is unnecessary for survival in either life cycle stage [111,112]. However, the nuclear encoded Complex I member genes are maintained [42], and while we did not identify full coverage for the ND transcripts, a vast majority of the gRNAs were found in both life cycle stages.

This study used high-throughput sequencing to characterize the gRNA transcriptome during the bloodstream stage of the trypanosome life cycle. This work suggests that gRNAs are expressed during both life cycle stages, and that differential editing patterns observed for the different mitochondrial mRNA transcripts are not due to the presence or absence of gRNAs.

### **Accession Numbers**

SAMN04302078, SAMN04302079, SAMN04302080, and SAMN04302081 NCBI's Sequence Read Archive.

### **Acknowledgments**

The authors dedicate this work in memory of David Judah, MS, DVM. He was a wonderful colleague.

We also acknowledge the work of Joshua Foster, Mark Johnson, James Rauschendorfer, Heather Tyler, Callie Vivian, and Alexis Weber who were involved in sorting and identifying gRNAs, the MSU RTSF for their contribution in deep sequencing and Ken Stuart and Jason

Carnes at the Center for Infectious Disease Research for supplying the *T. brucei* strains used in this study.

# **CHAPTER 3: MITOCHONDRIAL DUAL-CODING GENES IN**

## ***TRYPANOSOMA BRUCEI***

### **Abstract**

*Trypanosoma brucei* is transmitted between mammalian hosts by the tsetse fly. In the mammal, they are exclusively extracellular, continuously replicating within the bloodstream. During this stage, the mitochondrion lacks a functional electron transport chain (ETC). Successful transition to the fly, requires activation of the ETC and ATP synthesis via oxidative phosphorylation. This life cycle leads to a major problem: in the bloodstream, the mitochondrial genes are not under selection and are subject to genetic drift that endangers their integrity. Exacerbating this, *T. brucei* undergoes repeated population bottlenecks as they evade the host immune system that would create additional forces of genetic drift. These parasites possess several unique genetic features, including RNA editing of mitochondrial transcripts. RNA editing creates open reading frames by the guided insertion and deletion of U-residues within the mRNA. A major question in the field has been why this metabolically expensive system of RNA editing would evolve and persist. Here, we show that many of the edited mRNAs can alter the choice of start codon and the open reading frame by alternative editing of the 5' end. Analyses of mutational bias indicate that six of the mitochondrial genes may be dual-coding and that RNA editing allows access to both reading frames. We hypothesize that dual-coding genes can protect genetic information by essentially hiding a non-selected gene within one that remains under selection. Thus, the complex RNA editing system found in the mitochondria of

trypanosomes provides a unique molecular strategy to combat genetic drift in non-selective conditions.

## **Author Summary**

In African trypanosomes, many of the mitochondrial mRNAs require extensive RNA editing before they can be translated. During this process, each edited transcript can undergo hundreds of cleavage/ligation events as U-residues are inserted or deleted to generate a translatable open reading frame. A major paradox has been why this incredibly metabolically expensive process would evolve and persist. In this work, we show that many of the mitochondrial genes in trypanosomes are dual-coding, utilizing different reading frames to potentially produce two very different proteins. Access to both reading frames is made possible by alternative editing of the 5' end of the transcript. We hypothesize that dual-coding genes may work to protect the mitochondrial genes from mutations during growth in the mammalian host, when many of the mitochondrial genes are not being used. Thus, the complex RNA editing system may be maintained because it provides a unique molecular strategy to combat genetic drift.

## **Introduction**

Trypanosomes are one of the most successful parasites in existence, inhabiting an incredibly wide range of hosts [2,3]. The dixenous members cycle between two distinct hosts and can encounter different environments with distinct metabolic constraints. These parasites are unique in that they all possess glycosomes (where glycolysis occurs) as well as mitochondria [16]. The salivarian trypanosomes (e.g. *T. brucei*, *T. vivax*) are especially interesting, because

they are exclusively extracellular in their mammalian hosts, continuously replicating within the bloodstream over periods of months. During this stage of the life cycle, the mitochondrion is down-regulated, lacking both Krebs cycle enzymes and a functional electron transport chain (ETC) [21]. Successful transition to the fly vector, requires activation of the ETC and ATP synthesis via oxidative phosphorylation. This unique lifecycle leads to a major problem: when the mitochondrial genes are unused, they are not under selection, hence the integrity of these genes are threatened by genetic drift [75,113]. Exacerbating this, salivarian trypanosomes undergo a severe bottleneck as they transition through the tsetse fly and into the mammalian host, and then within the bloodstream, they undergo multiple bottlenecks at each antigenic switch, as they evade the host immune system [22]. Such bottlenecks create additional forces of genetic drift, where genes can be lost even if their deleterious fitness effect is considerable. These parasites possess several unique genetic features, including RNA editing of the mitochondrial transcripts. RNA editing creates open reading frames in “cryptogenes” by insertion and deletion of uridylate residues at specific sites within the mRNA. The U- insertions/deletions are directed by small guide RNAs (gRNA) and can repair frameshifts, generate start and stop codons and more than double the size of the transcript (for review see [4]). While the mRNA cryptogenes are encoded on maxicircles ( $25\pm50$  copies per DNA network), the guide RNAs are encoded on thousands of 1 kb minicircles, encoding  $3\pm5$  gRNA genes each [8]. This effectively means that the genetic information for the edited mitochondrial mRNAs is dispersed between the mRNA cryptogenes on the maxicircles and the thousands of gRNA coding minicircles. The extensive editing of a single transcript can require more than 40 gRNAs and hundreds of editing events [9]. While the initial gRNA can interact with the 3' end of the

pre-edited transcript, all subsequent gRNAs anchor to edited sequence created by the preceding gRNA. Hence, editing proceeds from the 3' end to the 5' end of the mRNA transcript with the terminal gRNA (last one in the cascade) often creating the start codon needed for translation. This sequential dependence means that with even high accuracy rates for each gRNA, the overall fidelity of the process is astonishingly low. A major question in the field has been why this fragile and metabolically expensive system of RNA editing would evolve and persist.

Another level of complexity in the kinetoplastids RNA editing process was the detection of an alternative editing event that leads to the production of a functionally discrete protein isoform. Alternative editing of Cytochrome Oxidase III (COIII) is reported to generate a novel DNA-binding protein, AEP-1, that functions in mitochondrial DNA maintenance [67,68]. In this transcript, one alternative gRNA generates sequence changes at two sites that links an open reading frame (ORF) found in the pre-edited 5' end, to the 3' transmembrane domains found in the COIII edited ORF. This was the first indication, that one cryptogene could contain information for more than one protein. Here, we show that as many as six additional cryptogenes also encode for more than one protein. Analyses of the terminal gRNA populations indicate that gRNA sequence variants exist that can alter the choice of the start codon and the open reading frame by alternative editing of the 5' end of the mRNA. Mutational bias analyses indicate that six of the mitochondrial genes may be dual-coding, with RNA editing allowing access to both reading frames. Dual-coding genes are defined as a stretch of DNA containing overlapping open reading frames (ORFs) [85,86]. Of particular interest are dual-coding genes that contain two ORFs read in the same direction: a canonical protein (normally annotated as

protein coding in the literature) and an alternative ORF. Maintaining dual-coding genes is costly, as it constrains the flexibility of the amino acid composition of both proteins. Hence, it is thought that dual-coding genes can survive long evolutionary spans only if the overlap is advantageous to the organism [93]. We hypothesize that trypanosomes use dual-coding genes to protect genetic information by essentially hiding a non-selected (ETC) gene within one that remains under selection. Thus, the ability to access overlapping reading frames may be added to a growing list of gene protective strategies made possible by the complex RNA editing process [74,75,113].

## **Materials and Methods**

### **Trypanosome growth**

*T. brucei* procyclic clones from IsTAR (EATRO 164), TREU 667 and TREU 927 cell lines were grown in SDM79 at 27°C and harvested at a cell density of 1-3x10<sup>7</sup>. The TREU 667 cell line was originally isolated from a bovine host in 1966 in Uganda [114]. The TREU 927 cell line was originally isolated from *Glossina pallidipes* in 1970 in Kenya [115]. The EATRO 164 strain was isolated in 1960 from *Alcephalus lichensteini* and maintained in the lab of Dr. K. Vickerman until being obtained by Dr. Ken Stuart in 1966 [103]. Dr. Stuart derived the procyclic form from the bloodstream form culture in 1979.

### **Guide RNA isolation, preparation, and sequencing**

Mitochondrial mRNAs and gRNAs were isolated as previously described [9]. All RNAs were treated with Promega DNase RQI. In order to isolate gRNAs from TREU 667 and TREU 927 cells, RNAs were size fractionated on a polyacrylamide gel as previously described [9]. Guide

RNAs were then extracted and prepped for sequencing using the Illumina Small RNA protocol [9]. Libraries from TREU 667 and TREU 927 were deep sequenced on the Illumina GAIIx; reads were processed and trimmed as previously described [9].

### **Messenger RNA preparation and sequencing**

In order to isolate target mRNAs, isolated TREU 667 mitochondrial RNAs were reverse transcribed using the Applied Biosystems High Capacity cDNA Reverse Transcription Kit. CR3 cDNAs were amplified via PCR using the following primers (underlined portions are gene specific and non-underlined portions are tag regions used in deep sequencing reaction):

CR3DS5'NEV:ACACTGACGACATGGTTCTACAAGAAATATAATATGTGTATG

CR3DS3'170:TACGGTAGCAGAGACTTGGTCTCAATAACCCATTAAATAAAAAACAAAAATCC

After amplification, the products were purified using the QIAquick PCR Purification Kit, and paired end Illumina deep sequencing was performed on the Illumina Miseq (2x 250 bp paired end run). Low quality results were removed using FaQCs, adapters were removed using Trimmomatic and PEAR was used to merge paired end reads. Finally, Fastx was used to compile identical reads while maintaining the number of redundant reads. CR3 edited transcripts were identified by comparing sequence downstream of the 5' never edited region to the edited CR3 sequence. Guide RNAs were identified by using the mRNA sequences as queries against our existing gRNA databases, as previously described [9].

### **Mutational frequency and editing conservation analysis**

Mitochondrial pan-edited genes were categorized as potentially dual-coding based on identification of extended alternative reading frames and/or presence of identified gRNAs that generate alternative 5' end sequences. These genes include CR3, CR4, ND3, the 5' editing

domain of ND7, ND9 and RPS12. Nondual-coding pan-edited genes include ATPase 6, COIII, ND8 and the 3' editing domain of ND7. Partially edited genes include CYb, Murf II and COII. Never edited genes include COI, ND1, ND2, ND4 and ND5. For all analyses, ND7 was considered as two separate coding regions: the 5' editing domain (ND7N) and the 3' editing domain (ND7C) [27]. As we hypothesize that only the 5' editing domain of ND7 is dual-coding, mutation calculations for ND7N was pooled with the dual-coding genes and ND7C was pooled with nondual-coding pan-edited genes. *T. brucei* and *T. vivax* mRNA sequences of mitochondrial encoded genes were aligned based on protein sequence using Clustal Omega [116]. Nucleotide sequence mutations were identified and their effects on the amino acid sequence were classified as silent, missense or nonsense mutations. Missense mutations were further divided into three groups based on the PAM 250 matrix where conversions with a value <0 were considered not conserved, conversions with a value 0x0.5 were considered modestly conserved, and conversions with a value >0.5 were considered strongly conserved [117]. Mutation frequencies were normalized for each gene using nucleotide sequence length. Frequencies were compared using unpaired t-tests.

The extent of editing conservation between *T. vivax* and *T. brucei* was calculated by aligning the pan-edited genes based on ACG sequence. For each alignment, each location between an A, C or G nucleotide where a U-residue was inserted or deleted in either sequence was considered an editing site. Editing sites were classified as identical in both sequences, altered in insertion or deletion length, having switched from an insertion site to a deletion site, or only occurring in one of the sequences. Percent editing conservation was based on total number of editing sites within each mRNA. Percentages were compared using unpaired t-tests.

A principal component analysis (PCA) was performed on all three reading frames of the pan-edited genes using the scikit-learn principal component analysis tool [118]. For this analysis, the predicted protein sequences for all three reading frames were aligned using Clustal Omega [116]. Missense, nonsense, and indel mutations were quantified. Missense mutations were further divided into three groups as described above. Each mutation type was quantified and the relative frequency of each mutation calculated based on protein amino acid length. The variables used in the PCA include the protein mutation frequencies and the percentage of identical editing sites in each mRNA. The first reading frame of each gene is defined as the ORF published in the literature.

## **Data availability**

CR3 sequence accession number: SAMN06318039. TREU 927 gRNA sequence accession number: SAMN06318154. TREU 667 gRNA sequence accession number: SAMN06318153. NCBI's Sequence Read Archive.

## **Results**

In *T. brucei*, analyses of the gRNA transcriptome for the pan-edited transcripts indicate that full editing involves a large number of gRNA populations [9,119]. In addition, most of the gRNA populations (population defined as guiding the same or near same region of the mRNA) contain multiple sequence classes. The sequence classes most often differ in R to R or Y to Y mutations, hence guide the generation of the same mRNA sequence (A:U and G:U base pairs allowed). During these analyses, we noted that the terminal gRNA population for Cytosine-rich Region 3 (CR3) (putative NADH dehydrogenase subunit 4L [120]), had 3' sequences that would extend editing beyond the previously identified translation start codon. In addition, this

population had several sequence variants that would generate different edited sequences in this region. The most abundant terminal gRNA would introduce a stop codon in-frame with two alternative AUG start codons found near the 5' end (Figure 8A). Other sequence classes however, would either bring the upstream AUGs into frame, or shift the reading frame. Intriguingly, the alternative +1 reading frame (ARF) did not contain any premature termination codons. In order to determine if these gRNAs were utilized, we used Illumina deep sequencing to identify the most abundant forms of fully edited CR3 transcripts. Surprisingly, we identified multiple forms of the mRNA (Figure 8A, 8B and 8C and APPENDIX D). The first was the fully edited sequence predicted by the most abundant gRNA identified (Figure 8A). The other transcripts however, had unique editing patterns at the 5' end (Figure 8B and 8C and APPENDIX D). Use of these 5' CR3 sequences allowed us to identify novel gRNAs. Predicted translation of these mRNA sequences indicate that they use the +1 reading frame, and that the protein generated would be the same length as the ORF previously identified. This suggests that CR3 is dual-coding, and that selection of the terminal gRNA determines which reading frame will be used. A re-examination of the terminal gRNAs for the pan-edited genes indicated that at least two other transcripts, NADH dehydrogenase subunit 7 (ND7) and ribosomal protein subunit 12 (RPS12), have identified gRNA sequence variants within the terminal gRNA population that allow access to alternative reading frames (Figure 8D and 8E and APPENDICES E and F). Interestingly, the alternative gRNA for ND7 generates a +2 frameshift with a 65 amino acid open reading frame. The ND7 transcript is differentially edited in two distinct domains separated by 59 nts that are not edited in the mature transcript (the HR3 region) [27]. Only the 5' domain is edited in both life cycle stages; full editing of the 3' domain was only found in

bloodstream form (BF) parasites. The stop codon for the +2 frameshift is found within the HR3 region, therefore this alternative protein would be generated by full editing of only the 5' domain. While the most abundant gRNA in the Eatro BF transcriptome (~50,000 reads) would generate a sequence utilizing the identified ND7 ORF (Figure 8D ORF), the most abundant gRNA (>100,000) in the Eatro 164 procyclic library is the +2 ARF gRNA (Figure 8D ARF and APPENDIX E). In RPS12, the alternative gRNA deletes an additional U-residue downstream of the existing start codon, shifting the reading frame into the +1 ARF (Figure 8E). Interestingly, in *Leishmania tarentolae*, a gRNA, gRPS12VIIIa, has been identified that would also shift the frame of the existing start codon into the +1 ARF [121].

The identification of gRNAs that could alter the reading frame led us to re-analyze the ORFs of the edited transcripts. In addition to CR3, we found extended ORFs in two different frames for Cytosine-rich Region 4 (CR4) and NADH Dehydrogenase (ND) subunit 9, while several others had shorter, but still significant ORFs in alternative frames (Figure 9). We do note that the original sequence publications for both CR4 and RPS12 (CR6) had indicated that the fully edited sequence contained extended ORFs in two different frames [30,48]. Additionally, NADH Dehydrogenase subunit 3 (ND3) was also considered to be potentially dual coding, based on mutational analysis described below.

As we did find potential ARFs in the edited transcripts, we analyzed the predicted ORFs for biases in their mutational pattern. Dual-coding genes often display an atypical codon mutation bias due to constraints imposed by the need to maintain protein function in both genes. In single-coding genes, changes in the third nucleotide of a codon give rise to synonymous amino acids, so this position (N3) is much less constrained. In contrast, in dual-

### A. CR3 ORF

M C M I Y N STOP M F D C L V L L F F Y C L F V H F F C  
 AUGUAGUAUGAUAAUAAAAAuuuAuuuuCAuuuuAuGuuuGA\*\*\*\*UUGuuuGGuuuGuuGuuuUUUAuuGuuuuGuuuuGuACAUuuuuuuuG  
 |||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
 nUAU---AGUAGUAUAAAGUGACAAGAGAGUACAGGCAAACAUAAAUAUA5'  
 nUUAAAUGUGAAAGUGAGAUAGACU---AACGAGCCAAAACAAGUUA5'

### B. CR3 ARF(+1)

M C M I Y I I I Y L F S L C L I V W F C C C F F I V C L Y I F F V  
AUGUGUAUGAUUAUAuuAuuuuuAuuuuuAuuuuuCuuuAuGuuuGA\*\*\*\*UUGuuuGGuuuuGuuGuuuuUUUAuuGuuuGuuuGuACAuuuuuuuuuG  
|||||:|||||:|||||:|||||:|||||:|||||:  
„UAUGAUAGUAAGUGAGUGGAGGUAAUAGGCCU---AACAAACCAAUUAUA5'

### C. CR3 ARF (+1)

D. *ND7*

ORF  
 M I S I I L C Y F W ST  
 M T T W ST M L F L V V F L H L Y R F T F G P Q  
 AUGACUACAUAGUAAGUAAuCAuuuuuAuGuuAuuuuGGuAGuuuuuuuuACAuuuGuAuCGuuuuACAuuuG\*GUCCACAU  
 :|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:  
 „UAUAGUAAGAUGCAUAGAAAGCCGUCAAGAGAAUGUAACAUAAA5’

**ARF (+2)**  
M T T W Y S I I Y V I F G S F F T F V S F Y I W S T A S R  
AUGACUACAUAGUAuAGUAuCAuuuAuGuuuuuuuGGuAGuuuuuuuACuuuGuAuCGuuuuACuuuG\*GUCCACAGCAuCCC  
|||:||||:|||:|||:|||:|||:|||:|||:|||:  
„UUUUAGUGUUACGGUGAGUGUUACUGAGAAGUUAAUUAUAU5'

### E. *RPS12*

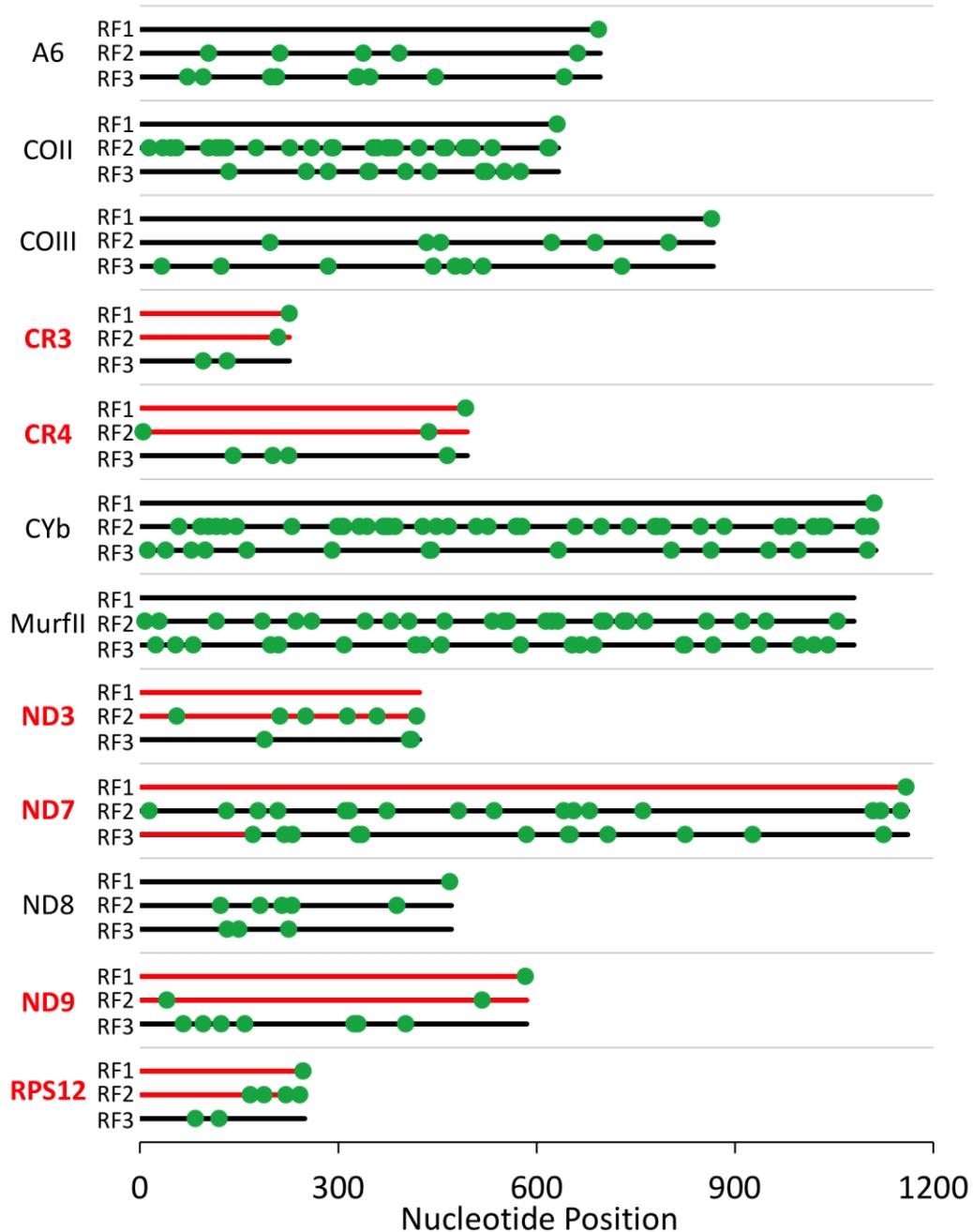
ORF M W F L Y G C C L R F V L F V  
 CAAACUAAAAGUAAuAuAu**uuAGuuuuuu**UCGUAuGuGA\*UUUUU**GUAUG\*****GuuGuuGuuu**AC\***GuuuuGuuuu**AuuuGu  
 „**UAUAUAGUUAAGAAGAUGCAUGUACU**-**AGAAGUUAAC**-**CAACAACAUUA5'**

**ARF (+1)** M W F C M V V V Y V L F Y L F  
 CAAACUAAGUAAuuuAAuuuuGuuuuuuuGCGuuGuGA\***\*UUUU****GUAUG\*****GuuGuu**GuuuAC\*GuuuuGuuuuAuuuuGu  
 |||||:||||:||||:||||:||||:||||:||||:||||:  
 „UUUUUUUAGAGUAGAGAAGUGCAUAUACU--AAGACAUAC-CAAUAUAUA' 5'

**Figure 8. Alternative editing of the 5' end of pan-edited genes results in access to different reading frames.** CR3 (A, B, C): Sequenced mRNA variants are aligned with gRNAs and predicted protein sequences. Inserted U-residues are lowercase while deleted U-residues are shown as asterisks. Canonical Watson-Crick base pairs (||); G:U base pairs (:). Previously identified start codons are doubled underlined. Potential upstream AUG start codons are indicated by wave underlines. Alternatively edited nucleotides are shown in red. Common anchor regions are shown in blue. ND7 (D), and RPS12 (E): Predicted mRNA and protein sequences, based on identified gRNAs.

coding genes, the N3 position in one frame is the N1 or N2 position in the alternative frame.

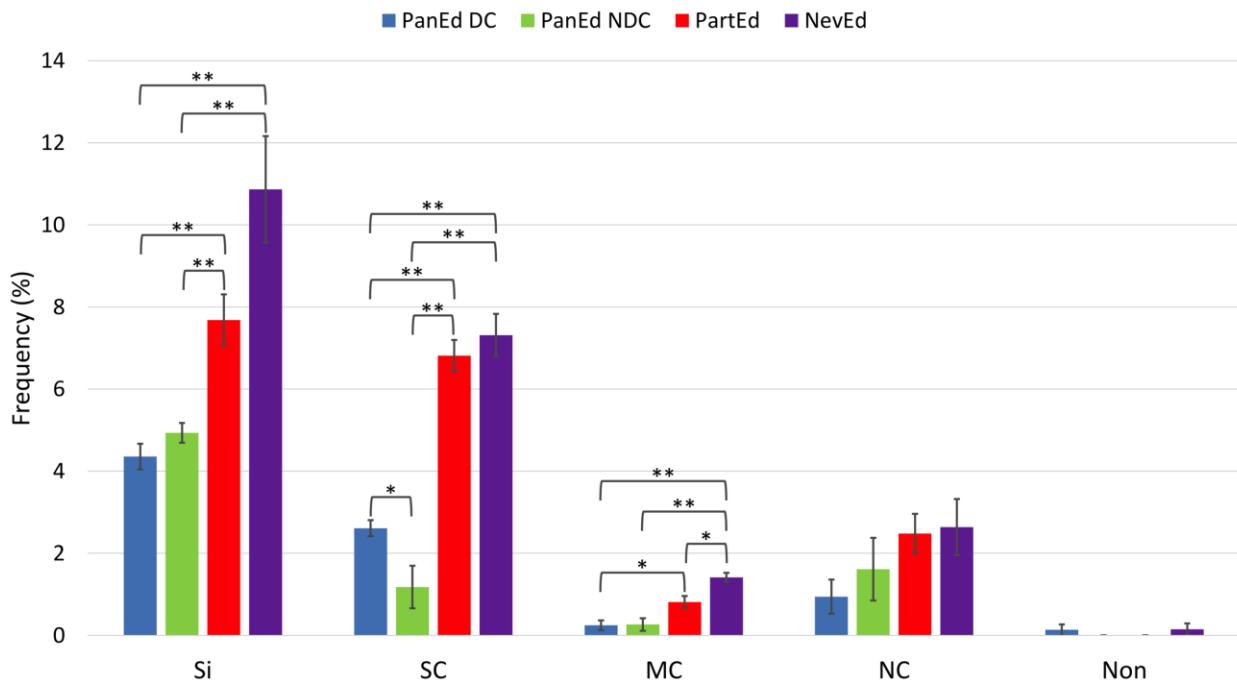
Therefore, they have low rates of synonymous mutations [122]. This codon bias has been used to develop algorithms to detect novel overlapping genes [123,124]. These algorithms however, cannot be used in the analysis of our edited transcripts as the two-component genetic system (mRNAs created by gRNA editing) introduces another layer of mutational constraint [125]. In addition, the edited sequence of the transcripts is known for only a limited number of



**Figure 9. Positions of stop codons on all RFs of the edited genes in *T. brucei*.** For each gene, reading frame 1 (RF1) is designated as the protein ORF previously identified. Hypothetical dual-coding reading frames are shown in red. A6 = ATPase 6; CO = Cytochrome Oxidase; CYb = Cytochrome b; Murf = Maxicircle unidentified reading frame; ND3 ±ND9 = NADH Dehydrogenase subunits.

kinetoplastids, and only the salivarian trypanosomes have the same general life cycle; other kinetoplastids, like *Leishmania* and *T. cruzi*, have evolved different infective cycles and are

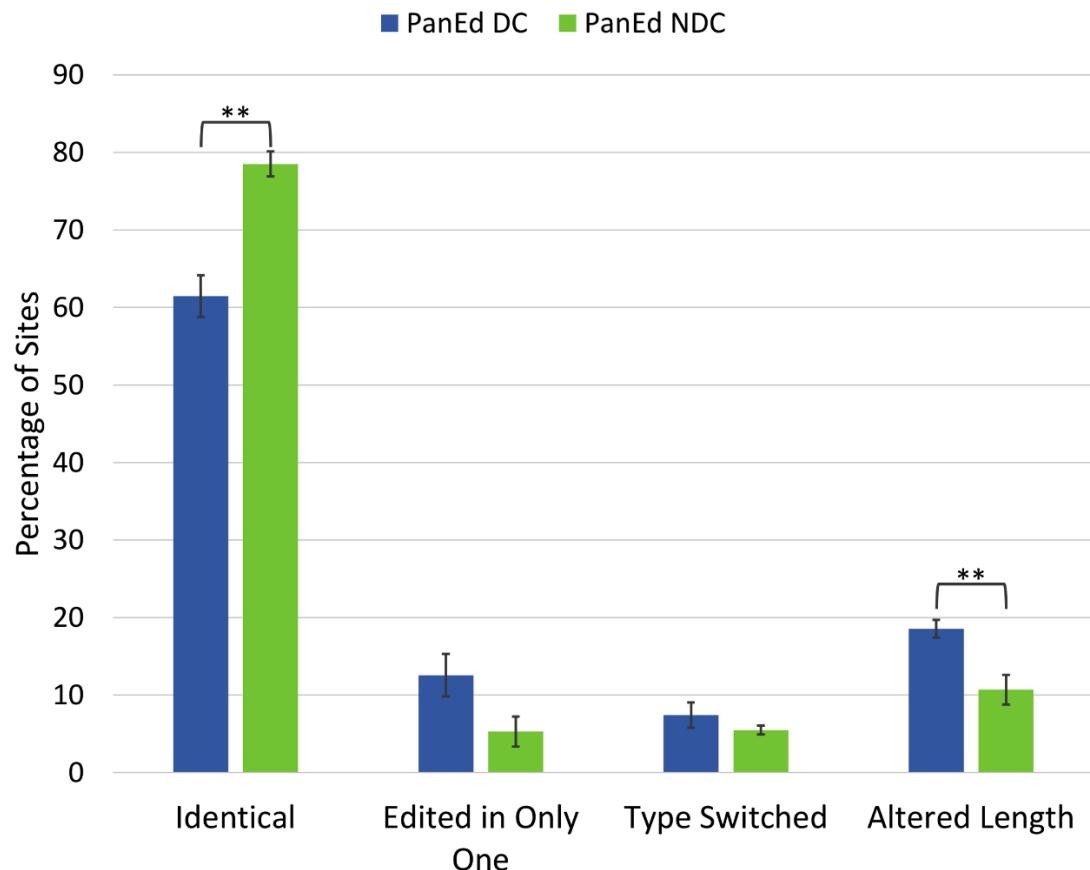
under very different selective pressures [113,126,127]. Fully edited sequences are known for *T. vivax*, the earliest branching salivarian trypanosome [52,128]. *T. vivax* differs from *T. brucei* in that they complete the insect phase of their life cycle entirely within the proboscis of the fly. This parasite has been described as an intermediate stage in the evolutionary pathway from mechanical transmission (ancestral) to full adaptation to the midgut and salivary glands of the tsetse fly [129]. Using the *T. vivax* sequence, we analyzed mutation patterns in all of the mitochondrial-encoded mRNAs (Figure 10). mRNA sequences were aligned by codons based on their protein alignments (Clustal Omega [116]). Mutated codons were identified and classified as silent, missense and nonsense mutations. Missense mutations were further divided into three groups based on the PAM 250 matrix [117]. These data clearly show that the RNA editing process significantly constrains the types of mutations tolerated within the mitochondrial genome. In comparison to the genes that are not edited (ND1, ND2, ND4, ND5, COI) or have limited editing (CYb, Murf II and COII), a distinct suppression of silent mutations and strongly conserved missense mutations were observed for all of the pan-edited genes, consistent with previous observations (Figure 10) [125]. A suppression of mutations that lead to moderately conserved amino acid replacements was also observed, but these were not as striking due to the low frequency of this type of mutation. No significant difference was observed in the frequency of not conserved missense mutations, though a trend towards a lower frequency of these mutations in the putative dual-coding genes (CR3, CR4, ND3, ND9, 5'ND7 and RPS12) was noted. This was complemented by a significant increase in the frequency of strongly conserved missense mutations in the putative dual-coding genes in comparison to the other pan-edited genes (3'ND7, ND8, A6 and COIII).



**Figure 10. Mutational frequencies in mitochondrially encoded genes categorized by effect on amino acid sequence.** *T. brucei* and *T. vivax* pan-edited dual-coding (PanEd DC), pan-edited nondual-coding (PanEd NDC), partially edited (PartEd), and never edited (NevEd) mRNA sequences were aligned based on their amino acid alignment (reading frame 1, defined as the reading frame encoding the gene product previously annotated in the literature) [116]. Mutations were categorized as silent (Si), strongly conserved (SC), modestly conserved (MC), not conserved (NC) or nonsense (Non). The amount of conservation was determined using the PAM 250 matrix, where conversions with a value  $0 < x \leq 0.5$  were considered modestly conserved, and conversions with a value  $> 0.5$  were considered strongly conserved, and conversions with a value  $\leq 0$  were considered not conserved. Error bars depict standard error. \*  $p < 0.05$ , \*\*  $p < 0.01$  (unpaired t-test).

Surprisingly, while the overall mutational frequency of the fully edited pan-edited genes was similar, a comparison of the conservation of editing patterns did show a significant difference between the putative dual-coding and the other pan-edited genes (Figure 11). The dual-coding genes consistently had a lower conservation of their editing pattern. Upon further examination, we found that most changes in the editing pattern resulted from thymidine insertions and deletions within the maxicircle DNA sequence, which was then corrected by the editing machinery. These types of mutations do not result in a change to the final mRNA

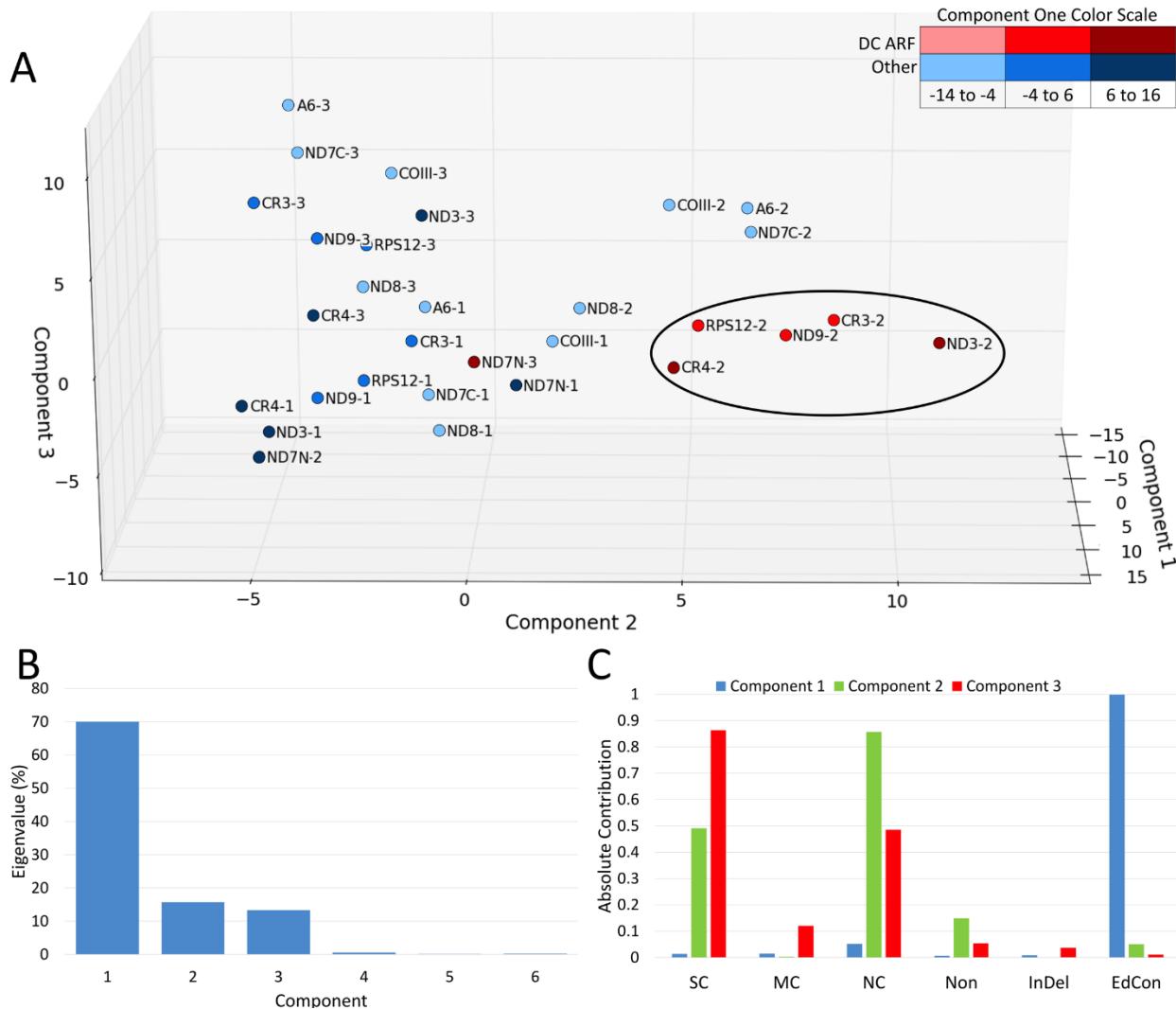
sequence once edited. The *T. brucei* (Tb) dual-coding genes appeared to consistently insert more U-residues, while *T. vivax* (Tv) had more U-residues encoded within the DNA sequence. Indeed, comparisons of the length of the coding regions of Tb and Tv cryptogenes (unedited sequence) show that the putative dual-coding genes are almost 10% shorter in Tb. In contrast, the nondual coding cryptogenes are not significantly shorter (~2.5%). Some of the other changes in editing patterns did generate small internal frameshifts as previously described by Landweber and Gilbert [125]. However, the high prevalence of internal frameshifts reported for COIII by Landweber and Gilbert is reflected in our analysis only for COIII and A6.



**Figure 11. Percent conservation of editing patterns between *T. brucei* and *T. vivax*.** Alignment of the fully edited mRNAs was based on ACG sequence (see APPENDIX G). Each editing site was defined as a site on at least one of the two aligned mRNAs where an editing event occurred. Sites were then classified as identical, only identified in one of the two sequences, type switched (one site is an insertion and the other is a deletion), or altered in length. Error bars depict standard error. \*  $p<0.05$ , \*\*  $p<0.01$  (unpaired t-test).

Since differences in the types of amino acid mutations were observed, we performed a principal component analysis on the frequency of mutation types for all three reading frames of the pan-edited genes (Figure 12). In addition, we included the percentage of editing site conservation as a variable. This analysis clearly clustered the putative +1 dual-coding transcripts (reading frame 2). The first component (z-axis,) is strongly based on editing conservation, and separates the dual-coding genes from the other pan-edited genes as expected. While component 2 (x-axis) separated ORF1 and ORF3 from ORF2 of each gene, component 3 clearly separated the dual-coding ORF2s from nondual-coding ORF2s. The ND7N ORF3 was the only exception, and the gRNA data suggests that it is a dual-coding gene using the +2 (ORF3) reading frame. This suggests that an additional layer of mutational constraint beyond that imposed by the RNA editing process can be detected for six of the extensively edited transcripts.

Because dual-coding genes are often conserved in multiple species, we analyzed the available sequences of other kinetoplastids (*Leishmania tarentolae* (Lt), *Leishmania mexicana amazonensis* (Lma), *Phytomonas serpens* (Ps), *Perkinsela CCAP1560/4* (Pk)) to determine if they also contain multiple overlapping reading frames with homology to those found in *T. brucei*. Interestingly, many of the alternative reading frames did show some homology to the ARFs found in Tb. However, most of these ARFs are punctuated with stop codons (APPENDIX H). Extended alternative reading frames are found in CR3, 5'ND7 and RPS12 in Ps. However, the extended ARF in the Ps CR3 is in the +2 reading frame and the ND7 and RPS12 ARFs shows very little homology with the Tb/Tv ARF (APPENDIX H Parts A, D and F) [130]. Interestingly, while *Perkinsela* has lost many of the genes in the mitochondria, RPS12 was retained [131]. The Pk RPS12 ARF possesses a near full open reading frame with one stop codon three codons after an



**Figure 12. Principal component analysis of frequency of amino acid mutation types and editing conservation between *T. brucei* and *T. vivax* pan-edited transcripts.** A. First factorial plan (z-axis: first component, x-axis: second component, y-axis: third component). ND7N = ND7 5' editing domain, ND7C = ND7 3' editing domain. ORF2 and ORF3 are defined as the +1 and +2 reading frames, respectively. B. Histogram of eigenvalues for first six components. Eigenvalues represent the amount of the variance accounted for by each component. C. Absolute contribution of each analyzed mutation frequency to components 1, 2, and 3. Amount of conservation was determined using the PAM 250 matrix as described in Figure 10. Mutation type: SC = strongly conserved, MC = moderately conserved, NC = Not conserved, Non = Nonsense, InDel = insertion or deletion. Editing conservation (EdCon) was determined using alignments of edited mRNAs (APPENDIX G). Aligned editing sites were characterized as identical or altered.

in frame start near the 5' end of the edited transcript. This pattern is reminiscent of the conventionally edited sequences of CR3 and ND7, and could suggest that an alternative edit

may remove the stop codon, allowing access to the ARF. The *L. tarentolae* CR4 orthologue also has two extended ORFs. Interestingly, the published sequence for Lt CR4 appears to switch between the two ORFs (switch appears to occur in a stretch of 13 inserted Us) [77]. This may explain why only the carboxyl half of the published Lt CR4 showed good homology with Tb and Lma [132]. Translation of the Lt ARF does generate a protein with the N-terminus showing high homology to the conventional Tb and Lma CR4, while translation of the published ORF shows some homology to the Tb CR4 ARF (APPENDIX H Part B). These data are intriguing enough that these sequences should be re-examined. While most of the other pan-edited transcripts had multiple stop codons in the +1 and +2 reading frames, many did show good homology to the Tb ARF sequences. Particularly intriguing are the ND3, ND8 and ND9 alignments. While internal stop codons are found in Tv, Lt and Lma ND9 ARFs, they show strong homology to the Tb ND9 ARF throughout the protein (APPENDIX H Part E). In ND3, the amino ends of the ARFs show strong homology between all four of the *Trypanosoma* and *Leishmania* species (APPENDIX H Part C). This homology decreases after an internal stop codon found in the same position in 3 of the 4 species. As ND8 is the only other pan-edited gene in Lt, Lam and Ps, we also examined the conservation of the ND8 ORF and ARF, even though our mutational analyses did not tag the ND8 gene as dual-coding. While the ND8 ARFs were punctuated by multiple stop codons, they surprisingly also showed areas of strong homology between all 5 species, especially downstream of an internal methionine (APPENDIX H Part G). We do note that we cannot rule out the possibility that alternative editing can remove stop codons observed in the ARFs.

Analyses of the ARF predicted proteins suggest that they are all short transmembrane proteins with two or more predicted transmembrane alpha helices (Figure 13) [133,134]. While

functional homologues are often difficult to detect in trypanosomes, searches using the predicted protein sequence of each ARF did identify small molecule transport proteins with limited confidence. Using Phyre2, the ND7 ARF was identified as a homolog of the bacterial sugar transporter SemiSWEET (61.5% confidence) [135,136]. SemiSWEET, which forms homodimeric structures, is also a distant homolog of the yeast mitochondrial pyruvate carrier 1 (MPC1). This protein has two transmembrane alpha helices and forms a heterodimer with either of the other two pyruvate carrier proteins [137,138]. While still very speculative, it is intriguing that the small ARF proteins might oligomerize to form small mitochondrial membrane transporters.

**CR3ARF:** MCM1YKNNVYYVVLFWFWLWYIFFVFYLFVICFYVCYLVVFYWIFVFYLIWVYCCVYFFFLFYHLICCYHFYCI (75 A.A.)

**CR4ARF:** FWLCIFFFFFIWCVCVLCTVYGYIFYCCFVFCFCLFVWWVCFICFVIVVCFLLFWVVIFYWCDFCIVYFFCDVIILFIFFFYFVLCFLYCCFYLVCLVFFCIFCCVLCYFLIYFLCCFLFWVFLFFFVYVCYFLWLLLFC (143 A.A.)\*

**ND3ARF:** SKNPRFLFTLVCYHYFYICFWYCGLLFYL!VFFYVFYFYCIFLIVFVVVCGFRVVCMIWIHVWCFIHWIYVLLVVCFLYC (80 A.A.)\*

**ND7ARF:** MTTWYSIIYVIFGSFFTFSFYIWSTASRSTWCFMLFIVFLWWIYCLYWLYYRLFASWYRKVMWI (65 A.A.)

**ND9ARF:** FYFIVCVVDGVLFVLLIVVFCFFIVLFLVFFCFIVCFYFLICDFCFYIIVVICYWLIFVVVFVVLCCCIFYFVCFCVFCVLFCVVCLYFLDCVLVLVVFVMRFFCCWNANVLICLFLVILLVMIIFYIVYLLIGFLVFFCWSVIHYLVCYFVCCWRR (158 A.A.)\*

**RPS12:** MWFCMVVYVLFYIWFYIIWVRDCPVPVTDVYCMPPYFIYIILFGCCVVFFVVLLV (55 A.A.)

**Figure 13. Amino acid sequences of ARFs of dual-coding genes.** Predicted transmembrane regions are shaded in gray [133,134]. \*No start codon was identified and the amino acid sequence shown begins at the 5' end or after the first stop codon at the 5' end. Exclamation point indicates premature termination codon.

## Discussion

The work presented here, suggests that as many as six of the extensively edited mRNAs in *T. brucei* are dual-coding and that it is alternative editing using different terminal gRNAs that allows access to the two different reading frames. Deep sequencing of the 5' end of CR3 indicates fully edited transcripts that have access to both reading frames are present in the mitochondrial transcriptome and gRNA analyses indicate that three different cell lines contain

gRNAs that can alternatively edit the 5' ends of CR3, RPS12 and ND7. In addition, analyses of the mutational bias in pan-edited genes suggest that an additional layer of mutational constraint is observed in the putative dual-coding genes. While the overall mutational frequency observed for the fully edited mRNAs is similar for all pan-edited genes, the types of amino acid changes that appear to be tolerated are significantly different. This is consistent with these genes having to maintain functional proteins in two different reading frames. Analyses of other trypanosomes, do show that some of the ARFs have intriguing homology to the ARFs identified in *T. brucei* and *T. vivax*. However, most of the ARFs are punctuated with stop codons. These data are difficult to interpret because we cannot rule out the possibility that the stop codons are removed by alternative editing events. In addition, the other trypanosome species have evolved very different infective life cycles and are under different selective pressures. For example, *P. serpens* is a pathogen that infects important crops and is transmitted by sap-feeding bugs. These parasites have glucose readily available in both life cycle stages and are unique in that they lack a fully functional respiratory electron transport chain [64,65]. For *Leishmania*, all life cycle stages possess an active Krebs cycle and ETC linked to the generation of ATP [61,62,139]. These unique adaptations to different hosts suggest that they may not be under the same evolutionary pressure to maintain dual-coding genes.

Overlapping reading frames are common in viruses, and are thought to persist due to strong genome size constraints [87,88]. More recently however, over-lapping genes have been identified in mammalian and bacterial genomes [89–92]. In these organisms, size is not an issue and the potential advantage of overlapping genes is less clear. For dual-coding genes, the need to maintain both ORFs constrains the ability of each protein to become optimally adapted [93].

As this constraint can be alleviated by gene duplication, it is thought that dual-coding regions can survive long evolutionary spans only if the overlap provides a selective advantage. In mammals, many of the identified dual-coding genes like Gnas1 and XBP1, produce two proteins that bind and regulate each other [94,95]. For these proteins, dual-coding may be advantageous for the tight co-expression needed. An alternative model, suggests that under high mutation rates, the overlapping of critical nucleotide residues is advantageous because it may reduce the target size for lethal mutations [96]. This may be particularly important for organisms that have evolved to exist in dual-metabolic environments (two hosts). We hypothesize that the trypanosome mitochondrial ARFs encode small metabolite transporters that provide a distinct growth advantage to bloodstream form parasites. The complete overlap of these small transporter genes with electron transport chain (ETC) genes would protect the integrity of the ETC genes that are required only in the insect host. Thus, in trypanosomes, dual-coding genes may be a mechanism to combat genetic drift during extended periods of growth in non-selective environments. In *T. brucei*, it is known that a number of bloodstream form essential proteins are functionally linked to Krebs cycle or ETC genes. While not a classic dual-coding gene in that production of the alternative protein does not involve overlapping reading frames, the pan-edited COIII gene does contain the information for two distinct proteins, COIII and AEP-1. AEP-1 is important for kinetoplastid DNA maintenance and overexpression of the DNA-binding domain results in a dominant negative phenotype including decreased cell growth and aberrant mitochondrial DNA structure [68]. The nuclear encoded  $\alpha$ -ketoglutarate dehydrogenase E2 ( $\alpha$ -KDE2) is known to be a dual-function protein, in that it plays important roles in both the Krebs cycle and in mitochondrial DNA inheritance [97]. RNAi

knockdowns of this gene in bloodstream form (BF) trypanosomes also show a pronounced reduction in cell growth. Similarly, the Krebs cycle enzyme  $\alpha$ -ketoglutarate decarboxylase ( $\alpha$ -KDE1) is also a dual-function protein with overlapping targeting signals that allow it to be localized to both the mitochondrion and glycosomes [98]. RNAi knockdowns of  $\alpha$ -KDE1 in BF trypanosomes is lethal, suggesting that in addition to its enzymatic role in the Krebs cycle, it plays an essential role in glycosomal function in *T. brucei* [98]. It has been previously suggested that both alternative editing and dual-function proteins are important mechanisms for expanding the functional diversity of proteins found in trypanosomes [67,97–99]. We hypothesize, that in salivarian trypanosomes, an equally important role for these dual-coding/function genes may be the protection of genetic information.

The ‘why’ of the unique RNA editing process in kinetoplastids has been a long-standing paradox. The complex machinery and the sheer number of gRNAs required to direct the thousands of U-insertion/deletions indicate that this process is metabolically very costly. Initially, it was proposed that U-insertion/deletion editing (kRNA editing) was one of many RNA editing processes that were in fact relics of the RNA world. However, the very different mechanism of the RNA editing systems in existence, and their very limited distribution within specific groups of organisms indicate that they are more likely derived traits that evolved later in evolution [69,70]. The sheer complexity of the kRNA editing process, with no obvious selective advantage, led to the proposal that insertion/deletion editing arose via a constructive neutral evolution (CNE) pathway [71]. Indeed, RNA editing in trypanosomes is always mentioned in support of CNE as an example of how seemingly non-advantageous, complex processes can arise [72,73]. More recently however, it has been hypothesized that RNA editing

co-evolved with G-quadruplex structures found in the pre-edited mRNAs [74]. These structures are thought to be advantageous in that they can help regulate transcription in order to promote DNA replication and prevent kinetoplast DNA loss. However, they must be removed by the RNA editing system prior to translation [74]. Another prominent hypothesis is that RNA editing is advantageous because it is a mechanism by which an organism can fragment and scatter essential genetic information throughout a genome [75,76]. Kinetoplast DNA is far less stable than chromosomal DNA, and loss of minicircles due to asymmetric division of the kDNA network have been frequently observed, especially in laboratory cultures of *Leishmania* [76,77]. Buhrman et al. [76] suggest that the scattering of essential guide RNA genes throughout the DNA network, would prevent fast growing deletion mutants from outcompeting more metabolically versatile parasites during growth in the mammalian host. Using a mathematical model of gene fragmentation in changing environments (absence of functional selection), they showed a distinct advantage for gene fragmentation. In their model, the number of tolerable generations under periods of relaxed selective pressure was increased by more than 40% before loss of the ability to move to the next life cycle stage. If the dual-coding ARFs give BF trypanosomes a selective growth advantage similar to that observed by the COIII alternative protein AEP1, then the number of ‘essential’ gRNA genes would increase greatly. Currently, only AEP1, A6 and RPS12 mitochondrial genes have been experimentally shown to be essential [68,100,110]. In addition, the presence of alternative editing and dual-coding genes would complement the protection provided by gene fragmentation by also shielding the genes from deleterious point mutations within critical ETC genes. This suggests that the complex RNA editing system found in the mitochondria may therefore provide multiple molecular strategies

to increase genetic robustness. Protection of the mitochondrial genome during growth in the mammal would increase the capacity for successful transfer to an insect vector and maximize the parasites long-term survival and spread.

### **Acknowledgments**

We thank the Ken Stuart Lab for trypanosome cell lines and Chris Adami for helpful discussions. We would also like to acknowledge the Dr. Marvis A. Richardson Endowed Fellowship Fund for their recognition of LEK.

# **CHAPTER 4: ANALYSIS OF THREE PAN-EDITED MRNAS REVEALS**

## **DUAL-CODING GENES AND COMPLEX MULTIPATH EDITING**

### **Abstract**

*Trypanosoma brucei* is a single celled eukaryote that possesses a highly complex RNA editing system. In this system, a large set of small RNAs, called guide RNAs direct the insertion and deletion of uridines in mitochondrial mRNAs. These changes extensively alter the target mRNAs, up to doubling them in length. Recently, mutational analysis showed that several of the edited genes possessed capacity to encode two different protein products. These overlapped reading frames could be accessed through alternative RNA editing, that shifts the translated reading frame. In this study, we analyzed the editing patterns of three putative dual-coding genes, ribosomal protein S12 (RPS12), the 5' editing domain of NADH dehydrogenase subunit 7 (ND7 5'), and C-rich region 3 (CR3). We found evidence that fully edited ND7 5' and CR3 are can translate in more than one reading frame. Moreover, we found that CR3 has a complex set of editing pathways that vary substantially between cell lines, and that changing available energy sources also alters the editing preferences of CR3 and ND7 5'. These findings suggest that editing patterns can be influenced by the current environment, and that alternative editing may be utilized by the trypanosomes to introduce variation within this fragile editing system.

## Introduction

*Trypanosoma brucei* is a member of the Kinetoplastea, a group of protozoans characterized by a large network of DNA in their mitochondria known as the kinetoplast [1]. The kinetoplast is composed of two types of concatenated circular DNA molecules: maxicircles and minicircles. The maxicircles all encode mitochondrial ribosomal RNAs as well as 18 protein coding genes, most of which are components of the electron transport chain. The approximately 30-50 identical copies of the maxicircle make up a relatively small proportion of the kinetoplast [5]. Most of the DNA network is composed of 5,000 and 10,000 1 kb minicircles, each of which encodes 2-5 small non-coding guide RNAs (gRNAs) [8,33]. These gRNAs are used in the process of RNA editing. In *T. brucei*, RNA editing consists of specific uridine insertion and deletion events that render 12 of the 18 mitochondrially encoded mRNAs translatable [4]. The gRNAs act as templates for the large editosome complex which cleaves the mRNA, inserts or deletes the correct number of uridines and then re-ligates the mRNA in an energy intensive process. This is repeated until the mRNA is complementary to the small gRNA. Each gRNA directs edits that generate the anchor region for the next gRNA, thus the RNA editing process is sequentially dependent on correct editing by each gRNA. As editing of some of the extensively edited mRNAs can involve upwards of 40 gRNAs, this renders the process incredibly fragile. [140]. We hypothesize, that such an expensive and fragile process evolved in response to the unique life cycle of *T. brucei*.

*T. brucei* is a dixenous parasite, invading the bloodstream of a mammalian host and being transmitted between hosts by bite of a tsetse fly. Once taken up in a blood meal by the tsetse fly, *T. brucei* transitions into the replicating procyclic state in the midgut, and the energy

*T. brucei* requires for this replication is gained through metabolism of amino acids [19,20]. This is accomplished through use of a portion of the Krebs cycle and the electron transport chain (ETC), thus most of the ATP required is produced by the mitochondria [19–21]. This stage of the life cycle is followed by a dramatic bottleneck when the trypanosomes transition from the midgut to the salivary glands of the tsetse fly [22,23]. From the salivary glands, trypanosomes are then refluxed into their next mammalian host during a bloodmeal. Once the parasite is deposited into its mammalian host, it quickly transitions to utilizing glycolysis for its energy generation, removing the requirement for ATP production in the mitochondria [15]. While in the mammalian host, *T. brucei* lives entirely extracellularly. It is frequently subject to attacks by the host's adaptive immune system, and the population evades these attacks through antigenic variation [12]. This part of the life cycle can be quite long, with the longest known infection lasting 29 years [13]. This life cycle should make *T. brucei* particularly sensitive to genetic drift, especially for those genes which are not under selection (Krebs's cycle and ETC) and should make them extremely vulnerable to Muller's ratchet (the gradual increase of mutational load that eventually leads to extinction) [81–84]. One mechanism for protecting small asexual populations is by increasing the severity of the mutations that can occur. If mutations severely impact fitness, mutated individuals are selected out, preventing their fixation [79]. Recently, computer modeling studies suggest that small asexual populations can evolve this type of mechanism (termed “drift robustness”) in order to maintain fitness [80]. The sequential dependence of the kRNA editing process implies that the system is inherently fragile to mutations. Even a single point mutation can drastically change the editing pattern, and stop the editing process, aborting expression of the protein. Hence, the RNA editing process may

operate as a proof-reading system to weed out mutations by making them lethal. This is effective however, only if the mitochondrial genes are under selection. Previously, we showed that many of the mitochondrially pan-edited genes have a distinct mutational bias that is suggestive of dual-coding genes (coding two proteins by overlapping reading frames) [141]. The overlapping of ETC genes not under selection in the bloodstream stage with genes that are under selection during this stage of the life cycle, would prevent the accumulation of mutations. As the extensively overlapped genes share most gRNAs, this strategy would ensure that almost all of the genetic material is protected.

Our analyses suggested that out of the twelve pan-edited genes in *T. brucei*, six are potentially dual coding, and that the RNA editing system is used to determine which reading frame is accessed. In order to determine if mRNA transcripts with access to multiple open Reading frames (ORFs) exist within the mitochondrial transcriptome, we deep sequenced the mRNA transcript populations of three putative dual coding genes: ribosomal protein S12 (RPS12), the 5' editing domain of NADH dehydrogenase subunit 7 (ND7 5'), and C-rich region 3 (CR3). Using the previously generated gRNA transcriptomes, we constructed detailed editing pathways for each of these genes. The editing pathway of RPS12 was primarily linear, reflecting the high degree of conservation required for a gene that is essential [100,107]. We found no evidence of utilization of the gRNA that provides access to the alternative reading frame [141]. In contrast, we did identify transcripts using different reading frames for both CR3 and ND7 5'. This study indicates that RNA editing can be used to access multiple open reading frames using two different methods: in ND7 5', different gRNAs bring alternate start codons into frame and in CR3, different gRNAs can shift the reading frame of the existing start codon. In addition, CR3

showed incredible editing diversity, in that two different cell lines showed very different editing patterns, using different sets of gRNAs to edit the CR3 cryptogene. This suggests that the use of a gRNA-guided editing system can also dramatically increase protein diversity in spite of an incredibly rigid and mutationally fragile system.

## **Materials and Methods**

### **T. brucei culture and RNA Isolation**

*T. brucei* clones from strains EATRO 164 and TREU 667 were grown in SDM79 and harvested as previously described [9]. EATRO 164 cells grown in SDM79 were then gradually transitioned to SDM80 using serial 1:3 dilutions when cells reached a density of at least  $5 \times 10^6$  cells/mL. SDM80 was prepared as described by Lamour et al. with the exception of using undialyzed FBS, and reducing the amount of FBS added by half [142]. This results in the final concentration of glucose being 0.5 mM instead of 0.15 mM. This concentration is still well below that of SDM79, which has a glucose concentration of 6 mM. Once cells had been acclimated to SDM80, cells were harvested as previously described [9]. Mitochondrial vesicles were isolated using differential spins and mitochondrial RNA was then isolated from vesicles as previously described [9].

### **Preparation, Sequencing, and Analysis of mRNAs**

cDNAs were generated from isolated RNAs using the Applied Biosystems High Capacity cDNA Reverse Transcription Kit. CR3, RPS12, and ND7 5' editing domain cDNAs were amplified via PCR using the following primers (underlined sequences are gene specific and non-underlined sequences are tag regions used in deep sequencing reaction):

CR3 5': AACTGACGACATGGTTCTACAAGAAATATAAATATGTG

CR3 3' Short: TACGGTAGCAGAGACTTGGTCTACAAAAATTATTCATACTT

CR3 3' Extended: TACGGTAGCAGAGACTTGGTCTACAAAAATTATTCATACTTTT

RPS12 5': AACTGACGACATGGTTCTACACTAACACTTTG

RPS12 3': TACGGTAGCAGAGACTTGGTCTAAAAACATATCTTAT

ND7 5': AACTGACGACATGGTTCTACAGATAACAAAAAACATGAC

ND7 3': TACGGTAGCAGAGACTTGGTCTTTATATTCACATAACTTTCTGTAC

Amplified cDNAs from EATRO 164 cells grown in SDM79, EATRO 164 cells grown in SDM80, and TREU 667 cells grown in SDM79 were individually barcoded and combined in equal molar amounts. Samples were sequenced in a 2x250bp paired end format (PE250) using an Illumina MiSeq Standard flow cell and 500 cycle reagent cartridge, version 2. Sequence data was preprocessed as previously described [141].

Sequence data was then separated by cell line, growth media and gene. Sequence data was then analyzed using a new pipeline and program called SKETCH (Segmentation of Kinetoplast Edited Transcripts to Characterize editing Heterogeneity). This program allowed us to classify mRNAs at the block editing level and determine which editing patterns were most prominent.

For each set of sequences, SKETCH would remove low quality sequences whose sequences containing more than 5 mismatches to the unedited template, disregarding uridines. In order to classify the editing patterns observed in the mRNA transcripts, SKETCH requires a set of template sequences. Initially, the templates supplied to SKETCH were the conventional fully edited and unedited sequences for each of the three genes examined. These sequences were

then segmented based on the editing blocks previously defined by the locations of gRNA populations [9]. Each transcript was then classified by editing block, with each block being classified as matching the unedited sequence, matching the fully edited sequence or being unknown. After the initial characterization of the transcripts, the most abundant unknown sequences for each editing block were then added to the reference pool. Sequences were then reclassified by SKETCH based on the newly added reference sequences. This process was repeated until the most abundant forms of editing were identified. SKETCH code is available upon request. To validate the newly identified editing patterns as true alternatives, the new sequences were screened against the gRNA transcriptome as previously described [9]. Sequences with a gRNA match were then considered valid alternative edits.

### **Uridine deletion analysis**

For RPS12 and ND7 5', once full editing pathways were characterized, editing sites with DNA encoded Us were identified. Each encoded U site was then characterized based on the proximity of the preceding gRNAs' 3' poly-U tail as well as whether the all uridines at the site in question were deleted in the final fully edited sequence. For each site, a window was defined consisting of the 6 sites upstream and 6 sites downstream of the site in question. Using these parameters, the mRNA transcripts of RPS12 and ND7 were analyzed at each deletion window. The window of each transcript for each encoded U was examined and classified as unedited, fully edited or partially edited. Partially edited sequences were then classified based on the editing state of the encoded U site. For sequences with total deletions, each editing sequence was then classified based on the states in the 3' end of the window as either matching the fully edited sequence or not. Code available upon request.

## Results

In order to confirm that transcripts with access to two reading frames exist *in vivo*, we analyzed the mRNA transcriptomes for three of the putative dual-coding genes, RPS12, ND7 5' and CR3. This mRNA deep sequencing data was then used in combination with the sequenced gRNA transcriptomes, to generate precise editing pathway maps. In order to determine how robust the observed editing pathways were, we characterized editing in two different cells lines, TREU 667 and EATRO 164. In addition, we examined the effect of energy source on these editing pathways by using two different media, SDM79 and SDM80. SDM79 is the standard medium used to grow the procyclic stage parasite. However, it contains 6mM glucose, and experiments have shown that under these levels of glucose, the procyclic stage can grow in the absence of electron transport chain (ETC) activity [112,142–147]. The SDM80 medium was developed to more closely resemble insect gut conditions and has very low glucose concentrations [142]. Trypanosome growth in this medium requires ATP production using the ETC [142].

RPS12 is an essential component of the mitochondrial ribosome [30,100,107]. RPS12 is extensively edited (pan-edited) with 132 Us inserted and 28 Us deleted. Full editing is directed by 12 populations of gRNAs (defined as a group of gRNAs that edit the same region of an mRNA) [9,30]. In this analysis, we identify 10 populations, with three of the previously identified populations being combined with other populations that shared a very high amount of overlap. One new population (F) was identified through a search of the gRNA transcriptome under reduced stringency. Analyses of the canonical editing pattern indicate that there are two long ORFs, and mutational bias analyses indicate that both ORFs may be selected for [141]. The

longest ORF encodes the RPS12 protein and encompasses a second shorter ORF of unknown function [30]. Northern blots revealed that edited RPS12 mRNAs were found in both life cycle stages, however, edited mRNAs were more abundant in bloodstream form than procyclic form trypanosomes [30].

Because RPS12 is essential, we expected it to have a very robust editing pattern in both cell lines, as well as under both energy conditions. In contrast, neither ND7 or CR3 appear to be essential in the insect stage of the parasite [112]. The canonical ND7 has two separate editing domains that are edited independently [27]. Interestingly, while the 3' editing domain is fully edited only in the bloodstream life cycle stage, the 5' editing domain is edited in both life cycle stages [5,27]. In addition, the mutational bias analyses indicate that only the 5' editing domain has characteristics indicative of a dual coding gene. The canonically edited CR3 is also a putative Complex I member (ND4L) and is preferentially edited in the BS stage [47,120]. Complex I has been shown to be non-essential in both life cycle stages, and other mitochondrially encoded complex I subunits, ND3, ND8, and ND9, have been shown to be preferentially edited in the bloodstream stage [5,26,28,29,111,112].

RNA seq data was generated by reverse transcribing all mtRNAs using random primers. For both RPS12 and CR3, transcripts were then selectively amplified using sequence specific primers targeted to the terminal 5' and 3' never edited regions as to not bias against any possible editing pattern. For ND7, the 5' editing domain was selectively amplified using sequence specific primers targeted to the 5' never edited region and the homology region 3 (HR3) that separates the 5' and 3' editing domains [27]. The HR3 is a span of 59 nts that is also never edited, hence should not bias the analysis. The targeted transcriptome libraries were

generated from TREU 667 cells grown in SDM79 and EATRO 164 cells grown in SDM79 and SDM80. Additionally, for CR3, we generated another library using TREU 667 cell line mRNA by selecting for transcripts of a larger size, instead of taking transcripts of all sizes (SDM79). This allowed us to enrich the library for transcripts that had initiated the editing process. Amplified cDNAs were then gel purified, barcoded and combined in equal molar amounts for sequencing. While the number of total reads obtained did vary by cell line and media used, surprisingly few transcript were fully edited (canonical AUG + ORF). For both RPS12 and CR3, the majority of reads (>80%) were completely unedited (Table 7). CR3, which has previously been shown to be preferentially edited in the BS stage, had the lowest percentage of fully edited transcripts, with only 0.1% – 0.2% translatable transcripts detected in both cell lines and under both growth conditions. In contrast, while RPS12 had similar levels of unedited transcripts, a larger percentage of translatable transcripts were found. For this essential transcript, the number of fully edited transcripts differed between the two different cell lines; 2.3% in TREU 667 and 0.9% in EATRO 164. Growth of the EATRO cells in low glucose media (SDM-80) did result in a substantial jump in the both the number of transcripts that initiated the editing process, and the number of fully edited transcripts (4.16%). This suggests that energy source may influence editing efficiency. While the predominance of completely unedited transcripts found for both CR3 and RPS12 was surprising, these numbers are in line with those found in other studies [131,148,149].

The ND7 5' transcriptome analyses differed substantially from both RPS12 and CR3 in that the majority of these transcripts had initiated the editing process. The TREU cell line showed the highest editing efficiency with ~80% of transcripts having initiated editing and 9.7%

**Table 7. Editing efficiencies of RPS12 (A), ND7 (B), and CR3 (C).**

Transcript	Cell line and Media	Total # Reads	% unedited	% partial edited	% fully edited
RPS12	Treu 667, SDM79	787,584	89.8%	7.9%	2.3%
RPS12	Eatro 164, SDM79	846,549	92.6%	6.5%	0.9%
RPS12	Eatro 164, SDM80	1,381,092	81.3%	14.5%	4.2%
ND7 5'	Treu 667, SDM79	1,141,322	20.3%	70.0%	9.7%
ND7 5'	Eatro 164, SDM79	915,610	47.0%	52.8%	0.2%
ND7 5'	Eatro 164, SDM80	313,657	27.4%	72.1%	0.5%
CR3	Treu 667, SDM79	18,832	84.9%	15.0%	0.1%
CR3	Treu 667 enriched.	50,589	18.1%	73.1%	8.8%
CR3	Eatro 164, SDM79	348,210	93.2%	6.6%	0.2%
CR3	Eatro 164, SDM80	53,000	90.6%	9.3%	0.1%

of the transcripts fully edited and translatable. In contrast, in EATRO cells, only 53% of the transcripts had initiated editing, and a scant 0.2% had completed the editing process. As with RPS12, we did see an increase of efficiency in the cells grown in SDM80, with over 70% initiating editing. However, even with the large increase in initiation of the editing process, only a scant 0.5% of transcripts were fully edited (Table 7). The sharp drop in the ability to complete the editing process appears to be due to loss of an optimal gRNA for one region of this transcript (described below).

## Editing Cascade and Reading Frame Analyses

In order to determine if the low editing efficiencies were due to any one step in the editing cascades, a full analysis of each editing step was done. For these analyses, we developed a pipeline that used our gRNA database to distinguish true alternative edits from both mis-edited and partially edited transcripts. This pipeline uses two programs, Segmentation of Kinetoplast Edited Transcripts to Characterize Editing Heterogeneity (SKETCH), and the gRNA database search program previously described [9]. The SKETCH program analyzes segments of transcripts that are defined by the relative range of coverage of each gRNA population used in

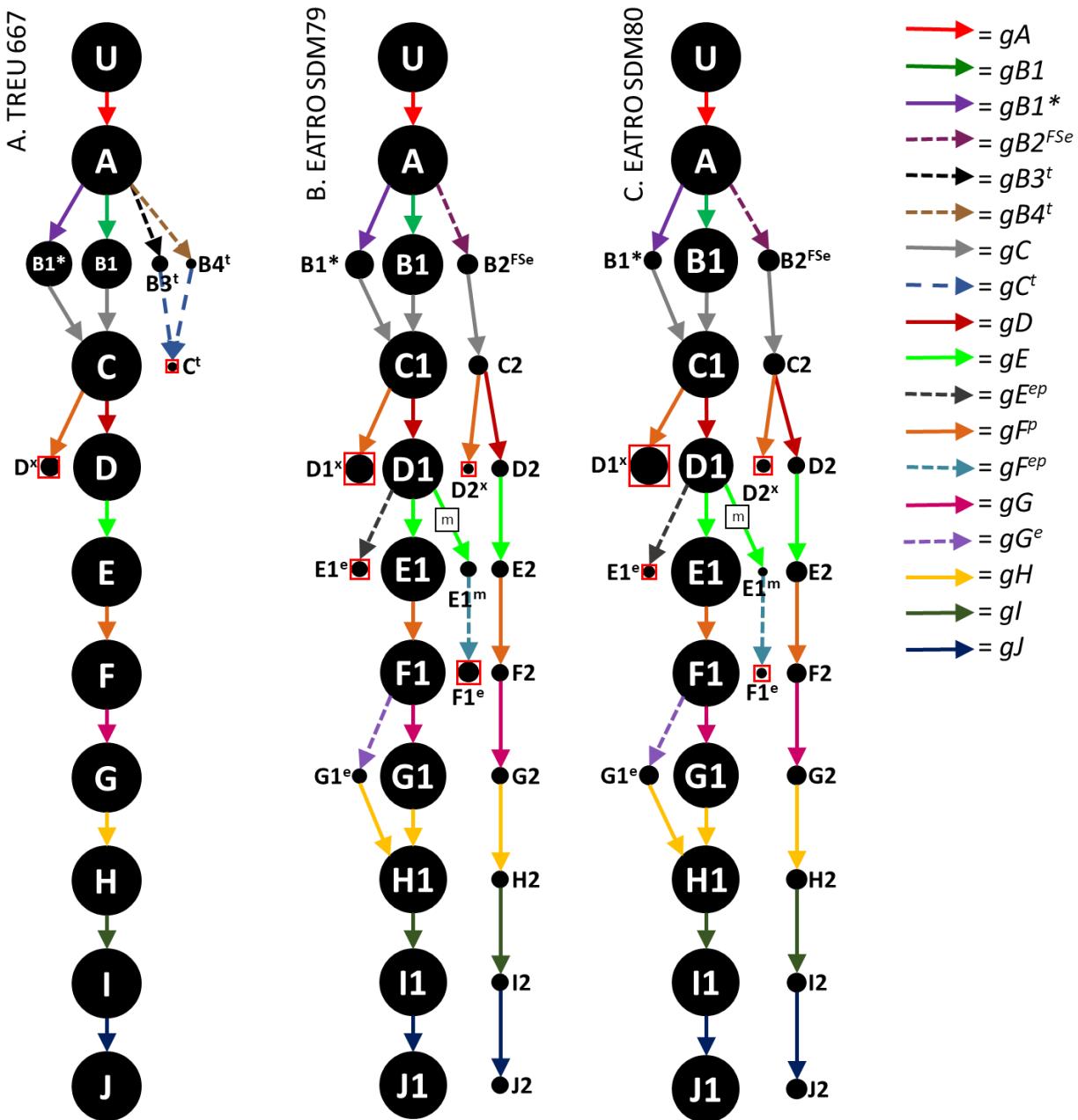
conventional editing patterns. Block sequences are compared to both the unedited sequence and the fully edited conventional sequence and then classified into unedited, fully edited and “unknown” blocks. Once the most abundant sequences of all segments are identified, transcripts containing each most abundant “unknown” sequences are used as queries against the gRNA database. If a gRNA is identified that can generate the edit, the sequence is considered a true alternative edit. If no plausible gRNA is identified the edit is considered a misedit or a junction, depending on the sequence and the status of other segments on that transcript with this sequence. By examining segments of transcripts independently, we were able to identify both branching and converging editing pathways.

## RPS12 Analysis

As expected, the essential RPS12 showed the most robust editing path. In all three analyses, the majority of transcripts used the same series of 10 gRNA populations (A – J) (For full gRNA sequences and alignments see APPENDICES I and J). Use of the final gRNA population (gJ) in the cascade lead to only the RPS12 ORF, and we found no evidence of an alternative AUG or frameshift leading to utilization of the second ORF. We do note that there is a downstream start codon, that if translated, would be read in the alternative reading frame (APPENDIX J). While the editing cascades were relatively straight forward, we did see some minor deviations (Figure 14). Editing of block B could utilize a number of different gRNAs, including several that were used in one cell line only (dashed arrows). gRNAs B1 and B1\* are variants of the same gRNA, with gB1\* introducing a single amino acid (aa) change (V/Y) (Figure 15). While editing using the TREU specific gB3<sup>t</sup> and gB4<sup>t</sup> gRNAs lead to a distinct editing “dead end” (dead end = disruption of the next canonical anchor sequence, and no detection of any further editing), the

EATRO specific gB2<sup>e</sup> did not disrupt the editing cascade. Use of this variant, however, did introduce a frameshift seven amino acids (aa) from the C-terminus (Figure 15). Because gB2<sup>e</sup> did not disrupt editing, a significant percentage (5.3% in SDM79 and 7.6% in SDM80) of translatable RPS12 transcripts did contain the alternative C-terminus (J2 transcripts). This alternative C-terminus was previously reported in the 29-13 strain [148], however, it appears to be absent in the TREU cell line.

A drop in editing efficiency was seen at the D to E block transition due to the incorrect utilization of the gF<sup>p</sup> guide RNA (D<sup>x</sup>) that disrupted the editing cascade (Table 8). While mis-editing by gF<sup>p</sup> was limited in the TREU 667 cell line (7.1% of D-block edited transcripts), its use was much more prominent in the EATRO cell line (17.5%), leading to a significant drop in transcripts that could continue past D-block editing. Interestingly, growth in SDM80 lead to a significant increase in mis-editing by gF<sup>p</sup>, with over 32.5% of transcripts using gF<sup>p</sup> incorrectly, resulting in a significant portion of dead-end transcripts. The EATRO cell line had additional minor dead-end pathways at the D to E transition. Misediting by a ND7 gRNA (gE<sup>ep</sup>) again disrupted any further editing, and mis-anchoring by the gE guide RNA (marked with box m) also led to the generation of an anchor sequence that could be used by a ND8 gRNA (gF<sup>ep</sup>) disrupting any further editing. Interestingly, the editing efficiency did not drop as transcripts transitioned to the next block of editing (Table 8). In EATRO-SDM80 cells, the editing efficiency at level F is ~5.6%, and at level G it actually increases to 5.9%. Editing efficiency at the block level is



**Figure 14.** Observed RPS12 editing pathways in the TREU 667 cell line (A) and the EATRO 164 cell line grown in SDM79 (B) and SDM80 (C). U = unedited transcripts. Dot sizes are proportional to the percent of block level edited transcripts using the gRNA indicated. Colored arrows indicate the gRNA population used. Dashed arrows with closed heads represent gRNA populations used in only one cell line (superscript 'e' or 't'). gRNA names with superscript 'p' represent promiscuous gRNAs. Dots enclosed by a red box represent end point mRNAs with no AUG start codon.  $gF^p$  is a promiscuous gRNA that edits both in the D and F editing block of RPS12. Arrows with a boxed 'm' represent a gRNA that has mis-anchored.

calculated based on the number of transcripts that match any of the fully edited sequences in that block, regardless of the condition of earlier blocks. Analysis of editing intermediates suggest that this increase occurs due to the ability of the downstream gRNA (the gF<sup>P</sup> population) to overwrite transcripts that have been previously edited through the G-level. Because of the overwriting gF<sup>P</sup> population, mRNAs exist that are fully edited at the G editing block but are in a transition state in block F.

```

T. brucei B1      --MWFLYGCCLRFVLFVLCYYMSPRLPSSGNRRVLYAVFYLYNFVWMLRCFFCC-FIGLVMSLFIIEGGGFVDLPGKYYTRIVS-----
T. brucei B2FSE --MWFLYGCCLRFVLFVLCYYMSPRLPSSGNRRVLYAVFYLYNFVWMLRCFFCC-FIGLVMSLFIIEGGGFVDLPGYKILFTYCKLDLDIYVF
T. vivax          --MWFLYGCCLRFVLFVLCYYMSPRLPSSGNRRVLYAVFYLYNFVWMLRCFFCC-FIGLVMSLFIIEGGGFVDLPGYKILFTYCKLDLDIYVF
L. tarentolae    MRVLFLYGLCVRFLYFCLVLYLSPRLPSSGNRRCLYAICYMFNILWFLC-VFCCVCFL-NHLLFIVEGGGFIDLPGVKYFSRFFLNA-----
L. donovani        VRVLYLYGLCVRFLFFSLVLYLSPQLPSSGNRRCLYAISIMFNILWFL-VFCCVFV-VHLLFIVEGGGFIDLPGVKYFSRFFCKS-----
L. amazonensis    VRVLYLYGLCVRFLFLCLVLYLSPRLPSSGNRRCLYAISIMFNILWYFL-VFCCVFV-IFQLFIVEGGGFIDLPGVKYFSRFCNVS-----
: :***: *;*: : * :*;:***: ***: ***: : :*;*: : .***: : ***: : ***: ;***: ***: ***: *
```

**Figure 15. Alignment of RPS12 proteins from *T. brucei*, *T. vivax*, *Leishmania tarentolae*, *Leishmania donovani*, and *Leishmania amazonensis*.** Asterisks indicate identical residues, colons indicate conserved residues. Highlighted amino acids are changes introduced by alternative editing (S>P, gG<sup>e</sup>, V>Y gB1\*) Alignment was generated by Clustal Omega [116]. RPS12 signature sequence is shown in bold [150].

**Table 8. Editing efficiency for each RPS12 gRNA population.** Percentages were calculated based on the number of transcripts that had completed each editing level out of the total number of RPS12 transcripts.

Block	Percent complete editing of block		
	TREU 667 (SDM 79)	EATRO 164 (SDM 79)	EATRO 164 (SDM 80)
Initiated Editing	10.2	7.4	18.7
A	8.9	4.6	15.9
B	6.9	4.4	14.6
C	6.0	4.2	13.5
D	4.9	3.1	11.7
E	4.5	2.2	6.9
F	3.8	1.4	5.6
G	3.7	1.4	5.9
H	3.7	1.3	5.7
I	3.4	1.2	5.4
J	2.3	0.9	4.2

The only other minor variation was the use of the EATRO specific gG<sup>e</sup> guide that occurs in a highly cytosine-rich region (Figure 16). Previous examinations of the gRNA coverage in this

region identified only rare gRNAs with multiple C:A basepairs, alignment mismatches and with gaps between adjacent gRNAs [9,119]. While this analysis did extend the identified gRNA population and eliminated the gap region, we did not identify either mRNA sequences or gRNAs that improved the alignment mismatches (Figure 16). The use of alternative base pairs is not unheard of. A study of *in vitro* deletions found that alternative base pairs such as C:A, C:U, and C:C were tolerated to varying extents [151]. Interestingly, this portion of RPS12 encodes the signature sequence, which is nearly universal [150]. Use of the gG<sup>e</sup> variant gRNA results in a single point mutation, substituting a proline in place of a serine within this important sequence.

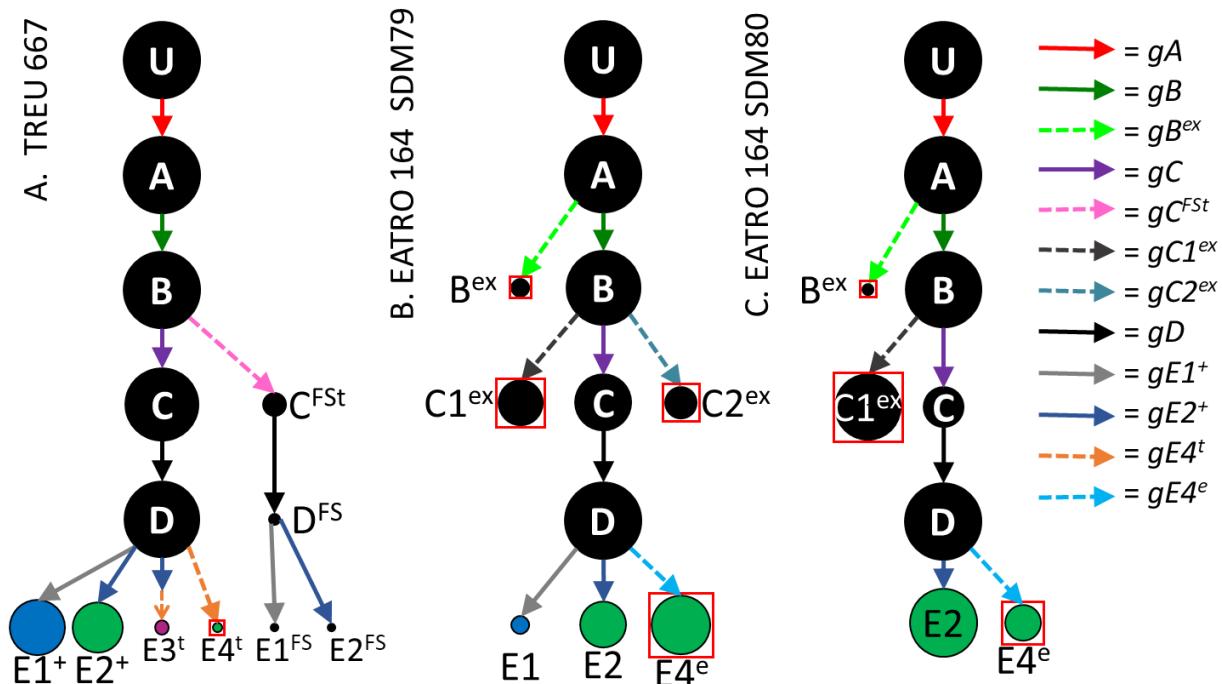


**Figure 16. Regions with poor gRNA coverage and functionally conserved residues in RPS12.** Functionally important aa residues are underlined [150]. Pipes (|) indicate Watson/Crick base pairs and colons (:) indicate G/U base pairs. Red highlighted hashtags (#) indicated gaps or mismatches, green highlighted # indicate C:A basepairs. The introduced substitution mutation introduced by use of the gG<sup>e</sup> gRNA is highlighted in yellow (S>P).

## ND7 5' analysis

Analyses of the ND7 5' targeted transcriptomes, indicate that full editing of the 5' domain requires five gRNA populations for both cell lines (Figure 17, APPENDIX K). Two variants of the terminal population (gE1 and gE2) were identified that resulted in different 5' terminal editing patterns (APPENDIX L). Translation of these editing patterns yields two different protein products in two different reading frames, one (RF1) encoding the canonical ND7 protein (E1) and the other (RF3) encoding a putative metabolite transporter (E2, see below) [141]. While transcripts for both open reading frames were found in both cell lines, there were notable

differences in the populations. The TREU 667 cell line had the highest editing efficiency with over 80% of the transcripts initiating the RNA editing process and ~9.7% of the transcripts fully edited through Block E. Use of the gE1 or gE2 gRNAs appeared to be equally efficient, resulting in nearly equal amounts of RF1 and RF3 fully edited transcripts. A small percentage of transcripts (4.9% of transcripts that completed block E editing) were observed that appeared to be mis-edited by a TREU specific gRNA (gE4<sup>t</sup>), leading to a dead-end product (no ORF). In addition, gE4<sup>t</sup> also appeared to be able to overwrite editing directed by gE2, to generate a small number of transcripts that could be translated in RF2 (pink E3<sup>t</sup>).



**Figure 17. Observed ND7 5' editing pathways in the TREU 667 cell line (A) and the EATRO 164 cell line grown in SDM79 (B) and SDM80 (C).** U = unedited transcripts. For arrow and gRNA naming descriptors see Figure 14. Dots enclosed by a red box represent end point mRNAs with no AUG start codon. + indicates that more than one mRNA form was condensed into this circle to simplify the figure (See APPENDIX M). Condensed forms encode largely the same amino acid sequence with only small variants. Terminal dots are colored blue for reading frame 1, magenta for reading frame 2, or green for reading frame 3. Boxed green dots have no functional start codon, but are translatable into reading frame 3 with the use of an alternative start codon (UUG).

While the number of “dead-end” pathways were very limited in the TREU cell line, use of the gC guide RNA population appeared to be very inefficient, resulting in a large drop in the percent of Block C-edited transcripts (25.8% drop, Table 9). A mutant gC gRNA ( $gC^{FSt}$ ), did result in a small percentage of transcripts with a frameshift C-terminus. Interestingly, while 9.1% of C block transcripts used the  $gC^{FSt}$  gRNA, only 2.4% of the transcripts that have completed D-block editing come from this minor branch. This suggests that this alternative edit decreases the efficiency of use of the subsequent gRNAs. In contrast, full editing of the ND7 5’ domain in the EATRO 164 cell line was very inefficient. While transcripts were able to initiate the editing process relatively efficiently (~50 – 70%, dependent on growth medium used), less than 1% of ND7 transcripts were fully edited at level E (Table 7). This appears to be due to the use of several EATRO specific gRNAs that disrupt further editing (Table 9). Again, the largest drop in editing efficiency occurred at the B to C-block transition. In addition, the EATRO specific use of  $gB^{ex}$ ,  $gC1^{ex}$  and  $gC2^{ex}$  all disrupted the editing cascade (Table 9). This compounded the editing efficiency problem, with a majority of C-block edited transcripts (47% in SDM79 and 73.4% for SDM80), no longer editing competent. The 5’ end of ND7 has multiple AUG sequences not created by the editing process. Translation predictions of these editing blocked transcripts ( $C1^{ex}$ ,  $C2^{ex}$ ) indicate that they do have ORFs that extend through the HR3 region. The protein product of  $B^{ex}$  transcripts is in the ARF, but is ten amino acids shorter, while the  $C1^{ex}$  and  $C2^{ex}$  products, which translate in the canonical ND7 reading frame, produce proteins that are both three amino acids shorter. Further drops in efficiency occurred due to an anchor mis-match (A:A) found in the gD guide RNA population (Figure 18, APPENDIX L). While the gD mutation is also observed in TREU, this cell line contains a sizable population of non-mutated gD guide

RNAs. Editing by the gE4<sup>e</sup> guide, results in a transcript with no in-frame AUG. However, translation of this transcript (E4<sup>e</sup>) in RF3 has no stop codons and we cannot rule out the possibility of a non-canonical START codon.

**Table 9. Editing efficiency for each ND7 5' domain gRNA population.** Percentages were calculated based on the number of transcripts that had completed each editing level out of the total number of ND7 transcripts.

Block	Percent complete editing of block		
	TREU 667 (SDM 79)	EATRO 164 (SDM 79)	EATRO 164 (SDM 80)
Initiated Editing	79.7	52.4	72.6
A	45.8	47.0	68.5
B	44.7	45.8	66.7
C	18.9	13.4	16.9
D	11.7	0.4	1.1
E	9.7	0.2	0.5



## CR3 Analysis

Previous work indicated that C-Rich region 3 is a putative Complex I member and that it is preferentially edited in the Bloodstream stage [47,120]. However, CR3 gRNAs are present in both life cycle gRNA transcriptomes, and PCR amplification of 5' edited transcripts were successfully cloned and sequenced in the TREU 667 procyclic cell line [9,119,141]. These studies indicated that multiple forms of the mRNA did exist that used different reading frames suggesting that CR3 is dual-coding and that it is selection of the terminal gRNA that determines which reading frame will be used [141]. In this study, we used primers flanking the editing domain in order to analyze the entire CR3 sequence. Interestingly, while 15% of the TREU CR3 transcripts had initiated the editing process, only 2.2% had completed editing by the initiating gA guide RNAs (Table 10). This suggests that the large drop in editing efficiency occurs due to incomplete editing by the block A guides. These gRNAs are fairly abundant, and we see no alignment issues, so it is unclear why editing of Block A is so inefficient (APPENDIX N).

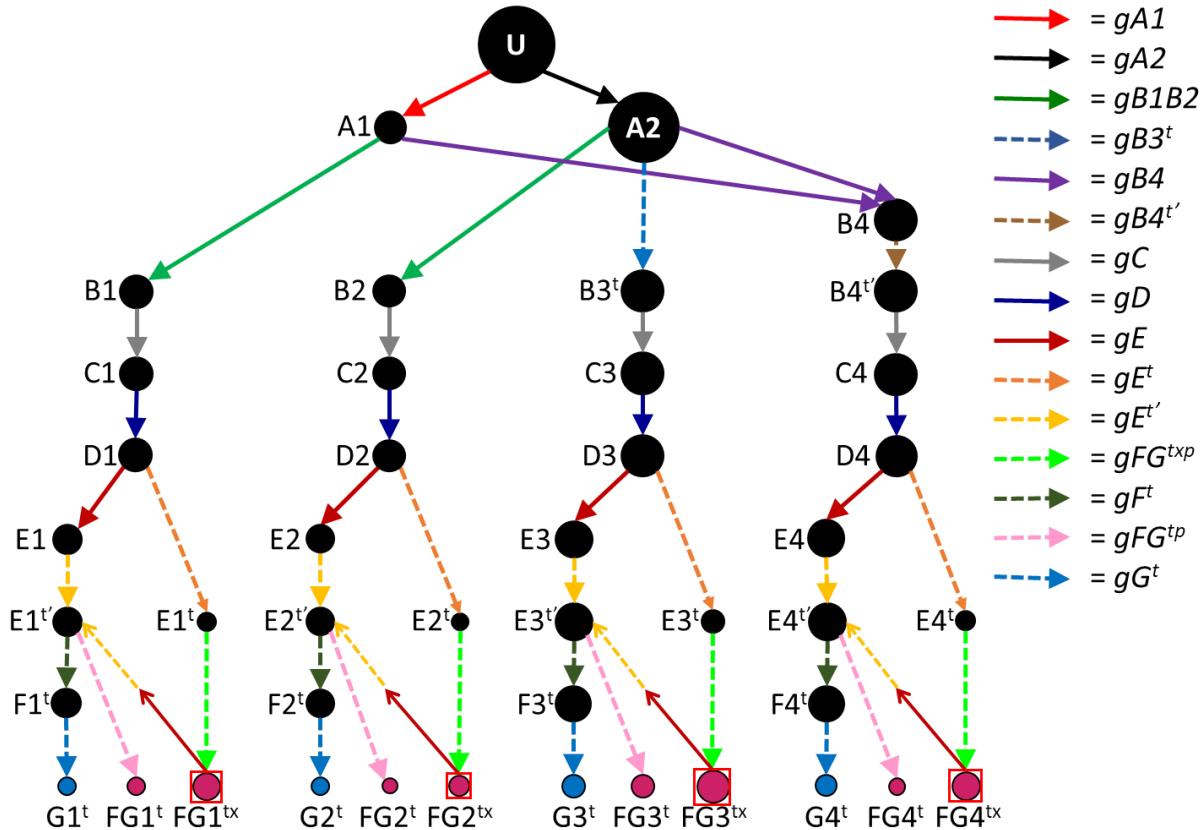
**Table 10. Editing efficiencies by block level of CR3.** Transcripts whose gRNAs covered two blocks (DEs and FGs) were included in both blocks for this calculation.

Level	Percentage Complete			
	TREU 667	TREU 667 (Enriched)	EATRO 164 (SDM 79)	EATRO 164 (SDM 80)
Initiated editing	15.1	81.9	6.8	9.4
A	2.2	68.3	1.9	2.3
B	1.8	56.6	1.2	1.8
C	1.7	54.6	0.9	1.5
D	1.6	52.0	0.6	1.1
E	1.2	39.1	0.6	1.0
F	0.8	22.3	0.2	0.4
G/FG	0.4	8.8	0.2	0.3

While the percentage of fully edited transcripts was very low percentage (0.2 – 0.4%, Table 10), we were able to again identify the major 5' alternative editing patterns that direct translation to either the ORF or to the +1 Alternative Reading Frame (RF2). To increase the robustness of the analyses, we also generated a biased CR3 transcriptome, by size selecting for longer transcripts during the amplification process. Analyses of the TREU transcriptome indicates that the full CR3 editing pathway has multiple branches, resulting in a total of 12 major forms of fully edited CR3 (Figure 19). These 12 forms are comprised of three major 5' editing patterns, paired with any of four different 3' editing patterns. The two initiating gRNAs identified (gA1 and gA2), direct identical editing patterns except gA2 inserts an additional three U- 1 phenylalanine). The gB guide RNAs all anchor in different areas (Figure 20A) and do introduce substantial AA changes near the 3' end (Figure 20B). However, all gB guide RNAs generate the anchor binding site (ABS) that is recognized by gC, hence all 4 nodes merge to a common sequence guided by gC and gD (APPENDIX N).

The 5' end editing patterns begin to diverge after Block D editing. FG<sup>tx</sup> transcripts are generated by the use of two subsequent gRNA populations, gE<sup>t</sup> and gFG<sup>txp</sup>. gFG<sup>txp</sup> is a promiscuous gRNA (previously identified as a ND7 gRNA) that spans both the F and G editing blocks. These transcripts were more abundant than both G<sup>t</sup> and FG<sup>t</sup>, however final editing using this gRNA does not generate a AUG start codon. It has been proposed that trypanosomes can use UUG as an alternative start codon, thus we cannot rule out the possibility that FG<sup>tx</sup> transcripts can be translated (Figure 20B). Analyses of intermediates suggest that the gE guide (red arrow) can in fact “overwrite” gE<sup>t</sup>, indicating that a proportion of these may still be re-edited into other forms. Editing via the gE population required an additional gRNA to generate

the anchor for either  $gF^t$  or  $gFG^{tp}$ . Generation of  $G^t$  transcripts (canonical CR3) requires 2 additional gRNAs, while  $FG^t$  (+1 ORF) transcripts are generated by a single gRNA population ( $gFG^{tp}$ ), another promiscuous gRNA (CR4).



**Figure 19. Observed CR3 editing pathways in the TREU 667 cell line.** U = unedited transcripts. For arrow and gRNA naming descriptors see Figure 14. Dots enclosed by a red box represent end point mRNAs with no AUG start codon. Terminal dots are colored blue for reading frame 1 or magenta for reading frame 2. Boxed magenta dots have no functional start codon, but are translatable into reading frame 2 with the use of an alternative start codon (UUG).

Surprisingly, when we examined the editing pathways of CR3 in the EATRO 164 libraries, we discovered that while three of the four initial 3' editing patterns were found in this library, editing beyond those patterns was completely divergent (Figure 21, APPENDIX O). A completely different set of gRNAs were used to generate fully edited CR3 transcripts (APPENDIX

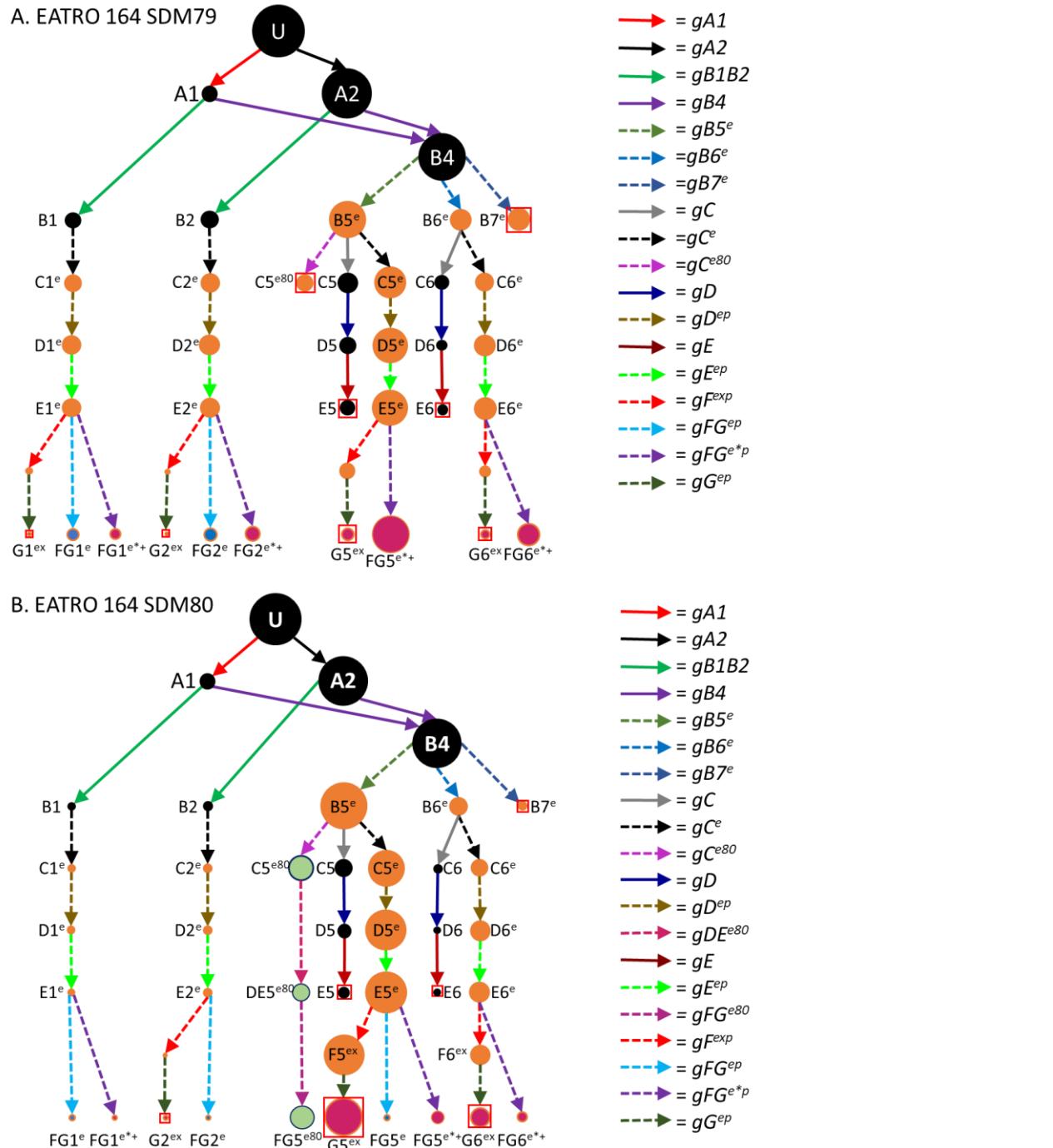
P). The divergent pathway did show some superficial similarities to the editing patterns observed in

A.	
B1	AuUGuuGuGuuuuAuAuuACAGAuuuu <u>AGuGuuAuCA</u> --- <u>uUAuuAuuGuAuAaAAGu<u>UUUCGUUAUAGAUUA</u></u>
B2	AuUGuuGuGuuuuAuAuuACAGAuuuu <u>AGuGuuAuCAuuu</u> <u>uUAuuAuuGuAuAaAAGu<u>UUUCGUUAUAGAUUA</u></u>
B3 <sup>t</sup>	AuUGuuGuGuuu <u>AuAuuAuuuCA</u> <u>GGuAuCAuuu</u> <u>uUAuuGuAuAaAAGu<u>UUUCGUUAUAGAUUA</u></u>
B4B4'	AuUGuuGuGuuu <u>AuAuuAuuuCA</u> <u>GuGuAuCAuuu</u> <u>uUAuuAuuGuAuAaAAGu<u>UUUCGUUAUAGAUUA</u></u>
B.	
ORF	
G1 <sup>t</sup>	MFDCLVLLFFYCLFVHFFCFLFVCDLFLCLLFSFCFLDFCFLFMGLLLCFILQIFS VII - IIVYKFSLLD
G2 <sup>t</sup>	MFDCLVLLFFYCLFVHFFCFLFVCDLFLCLLFSFCFLDFCFLFMGLLLCFILQIFS VII <b>I</b> IIVYKFSLLD
G3 <sup>t</sup>	MFDCLVLLFFYCLFVHFFCFLFVCDLFLCLLFSFCFLDFCFLFMGLLL <b>C</b> YYFRFY <b>G</b> I <b>I</b> IIVYKFSLLD
G4 <sup>t</sup>	MFDCLVLLFFYCLFVHFFCFLFVCDLFLCLLFSFCFLDFCFLFMGLLL <b>C</b> FFFFILSFDM <b>L</b> LSF <b>L</b> LYISFRY
ARF	
FG1 <sup>t</sup>	MCM <sup>b</sup> YKNNVYVVVLFWFWLYIFFV <sup>b</sup> FYL <sup>b</sup> VICFYVCYLV <sup>b</sup> V <sup>b</sup> FYW <sup>b</sup> IFV <sup>b</sup> FYLI <sup>b</sup> WVYCCVLYYRFLVLS- <b>LLY</b> ISFRY
FG2 <sup>t</sup>	MCM <sup>b</sup> YKNNVYVVVLFWFWLYIFFV <sup>b</sup> FYL <sup>b</sup> VICFYVCYLV <sup>b</sup> V <sup>b</sup> FYW <sup>b</sup> IFV <sup>b</sup> FYLI <sup>b</sup> WVYCCVLYYRFLVLS <b>F</b> LLYISFRY
FG3 <sup>t</sup>	MCM <sup>b</sup> YKNNVYVVVLFWFWLYIFFV <sup>b</sup> FYL <sup>b</sup> VICFYVCYLV <sup>b</sup> V <sup>b</sup> FYW <sup>b</sup> IFV <sup>b</sup> FYLI <sup>b</sup> WVYCCV <b>YIISDFMVSF</b> LLYISFRY
FG4 <sup>t</sup>	MCM <sup>b</sup> YKNNVYVVVLFWFWLYIFFV <sup>b</sup> FYL <sup>b</sup> VICFYVCYLV <sup>b</sup> V <sup>b</sup> FYW <sup>b</sup> IFV <sup>b</sup> FYLI <sup>b</sup> WVYCCV <b>YFFF</b> LY <b>H</b> LC <b>YY</b> <b>C</b> I
FG1 <sup>tx</sup>	LVVYCVYHC <sup>b</sup> IFLW <sup>b</sup> IFV <sup>b</sup> YVCYLV <sup>b</sup> V <sup>b</sup> FYW <sup>b</sup> IFV <sup>b</sup> FYLI <sup>b</sup> WVYCCVLYYRFLVLS- <b>LLY</b> ISFRY
FG2 <sup>tx</sup>	LVVYCVYHC <sup>b</sup> IFLW <sup>b</sup> IFV <sup>b</sup> YVCYLV <sup>b</sup> V <sup>b</sup> FYW <sup>b</sup> IFV <sup>b</sup> FYLI <sup>b</sup> WVYCCVLYYRFLVLS <b>F</b> LLYISFRY
FG3 <sup>tx</sup>	LVVYCVYHC <sup>b</sup> IFLW <sup>b</sup> IFV <sup>b</sup> YVCYLV <sup>b</sup> V <sup>b</sup> FYW <sup>b</sup> IFV <sup>b</sup> FYLI <sup>b</sup> WVYCCV <b>YIISDFMVSF</b> LLYISFRY
FG4 <sup>tx</sup>	LVVYCVYHC <sup>b</sup> IFLW <sup>b</sup> IFV <sup>b</sup> YVCYLV <sup>b</sup> V <sup>b</sup> FYW <sup>b</sup> IFV <sup>b</sup> FYLI <sup>b</sup> WVYCCV <b>YFFF</b> LY <b>H</b> LC <b>YY</b> <b>C</b> I

**Figure 20. Four different 3' end sequences found in the TREU 667 transcriptome for the CR3 transcript (A) and CR3 protein sequences (B).** U-residues inserted by editing are indicated by lowercase; different sequences created by the different gRNAs are highlighted in RED. Thick underline sequence indicates the anchor binding site (ABS) for the initiating gRNAs (*gA1* and *gA2*). **Green** = ABS for *gB1B2*; **Blue** = ABS for *gB3*; **Purple** = ABS for *gB4*. Bolded amino acids show sequence variants and shaded sequence shows position of predicted transmembrane domains [134]

TREU 667 cells. While both the B1 and B2 transcripts were directly edited by gC<sup>e</sup>, the B4 transcripts required an additional gRNA to generate the ABS recognized by gC<sup>e</sup>. In EATRO cells, B4 transcripts could be edited by 3 different gRNAs (gB5<sup>e</sup>, gB6<sup>e</sup> and gB7<sup>e</sup>). While gB7<sup>e</sup> disrupted editing, both gB5<sup>e</sup> and gB6<sup>e</sup> generated the anchor that could be used by either gC or gC<sup>e</sup>. Surprisingly, while the conventional CR3 gC guide RNA was clearly used by B4 transcripts, we saw no evidence of its use in the B1/B2 pathways. Transcripts using gC could be further extended by both gD and gE guide RNAs, however, no evidence of editing beyond the gE guides was observed. In contrast, use of the alternative gC<sup>e</sup> guide RNA population, could be extended by a series of additional guide RNAs, generating transcripts with functional AUG start codons.

However, many of the gRNAs used were promiscuous, in that they had been previously identified as gRNAs of other transcripts. As with the TREU editing pathway, we observe transcripts capable of being translated in two reading frames with the FG<sup>e</sup> mRNAs translating in RF1, and the FG<sup>e\*</sup> mRNAs translating in RF2 (Figure 21A, Figure 22). In addition, the G<sup>e</sup> mRNAs, while not having a functional “AUG” do translate into RF2 if the first “UUG” is used. As with ND7, we observed a shift in editing pattern preference when the EATRO 164 cells were changed from SDM79 medium to SDM80. Interestingly, a new fully edited form of CR3 appeared in the EATRO164 SDM80 library only. The gRNA gC<sup>e80</sup> is used in the EATRO SDM79 pathway, but editing appears to cease here. Cells grown in SDM80 continue this editing pathway with two additional gRNAs, gDE<sup>e80</sup> and gFG<sup>e80</sup> (Figure 21B). This mRNA is translatable, but produces a distinctly different and shorter protein product (Figure 21A). The protein products of the two different cell lines are highly dissimilar. Using bioinformatics tools to predict the secondary structure of these proteins, we find that the difference is most noticeable in the RF1s of the two cell lines (Figure 23). Interestingly, the RF2s have a very similar predicted secondary structure. This evidence suggests that the two different cell lines are able to use the CR3 transcript to create distinctly different protein products.



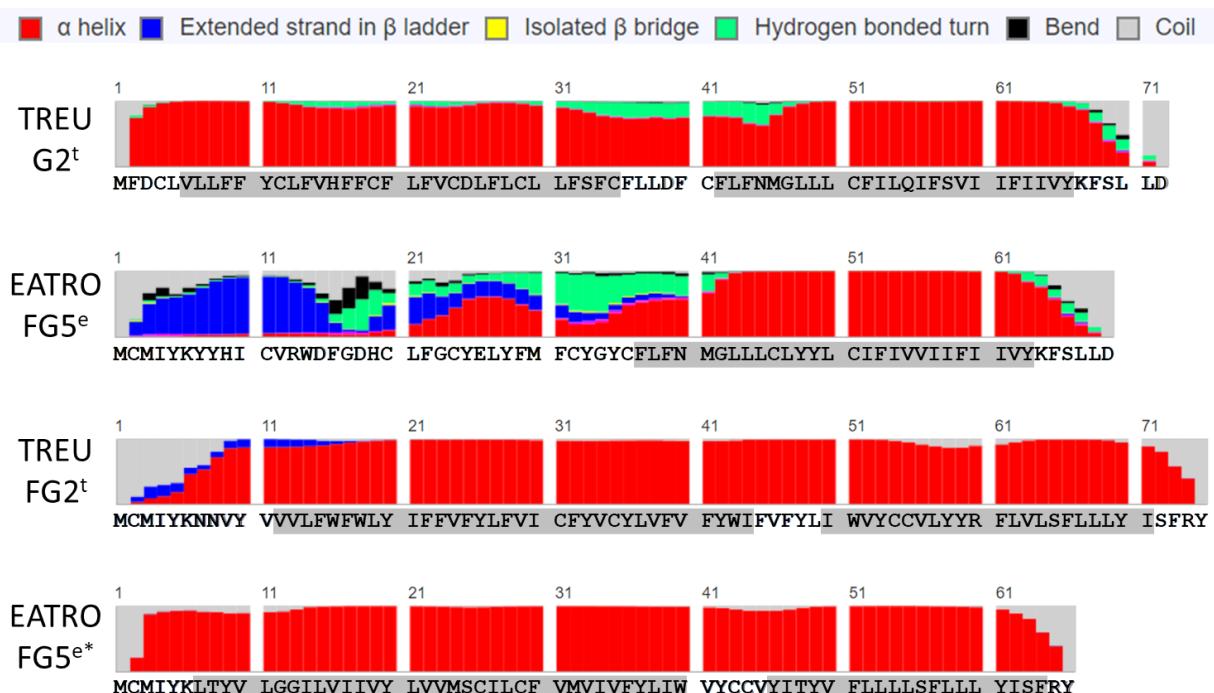
**Figure 21. Observed CR3 editing pathways in the EATRO 164 cell line grown in SDM79 (A) and SDM80 (B).** U = unedited transcripts. For arrow and gRNA naming descriptors see Figure 14. Dots enclosed by a red box represent end point mRNAs with no AUG start codon. Terminal dots are colored blue for reading frame 1 or magenta for reading frame 2. Boxed magenta dots have no functional start codon, but are translatable into reading frame 2 with the use of an alternative start codon (UUG). + indicates that more than one mRNA form was condensed into this circle to simplify the figure (See Figure 22). Condensed forms encode largely the same amino acid sequence with only small variants.

ORF	
FG1e	MCM <del>IYKYYHICVRWDFGDHCLFGCYELYFMFCYGYCFLFNMGLLL</del> CF--ILQIFS <del>VII</del> -IIVYKFSLLD
FG2e	MCM <del>IYKYYHICVRWDFGDHCLFGCYELYFMFCYGYCFLFNMGLLL</del> CF--ILQIFS <del>VII</del> <b>F</b> IIVYKFSLLD
FG5e	MCM <del>IYKYYHICVRWDFGDHCLFGCYELYFMFCYGYCFLFNMGLLL</del> CF <b>L</b> <b>Y</b> <b>L</b> <b>C</b> <b>I</b> <b>F</b> <b>I</b> <b>V</b> <b>V</b> <b>I</b> <b>I</b> <b>F</b> IIVYKFSLLD
ARF	
FG1e*v1	MCM <del>IYKLTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> -YYRFL-VLS-LLLYISFRY
FG2e*v1	MCM <del>IYKLTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> -YYRFL-VLS <b>F</b> LLL <del>Y</del> ISFRY
FG5e*v1	MCM <del>IYKLTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> <b>Y</b> <b>I</b> <b>T</b> <b>Y</b> <b>V</b> <b>F</b> <b>L</b> <b>L</b> <b>L</b> <b>S</b> <b>F</b> <del>Y</del> LLL <del>Y</del> ISFRY
FG6e*v1	MCM <del>IYKLTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> <b>Y</b> <b>I</b> <b>I</b> <b>L</b> <b>C</b> <b>I</b> <b>F</b> <b>I</b> <b>V</b> <b>V</b> <b>I</b> <b>I</b> <b>F</b> IIVYKFSLLD
FG1e*v2	MCM <del>IYKLTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> -YYRFL-VLS-LLLYISFRY
FG2e*v2	MCM <del>IYKLTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> -YYRFL-VLS <b>F</b> LLL <del>Y</del> ISFRY
FG5e*v2	MCM <del>IYKLTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> <b>Y</b> <b>I</b> <b>T</b> <b>Y</b> <b>V</b> <b>F</b> <b>L</b> <b>L</b> <b>L</b> <b>S</b> <b>F</b> <del>Y</del> LLL <del>Y</del> ISFRY
FG6e*v2	MCM <del>IYKLTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> <b>Y</b> <b>I</b> <b>I</b> <b>L</b> <b>C</b> <b>I</b> <b>F</b> <b>I</b> <b>V</b> <b>V</b> <b>I</b> <b>I</b> <b>F</b> IIVYKFSLLD
FG1e*v3	MCM <del>IYKNTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> -YYRFL-VLS-LLLYISFRY
FG2e*v3	MCM <del>IYKNTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> -YYRFL-VLS <b>F</b> LLL <del>Y</del> ISFRY
FG5e*v3	MCM <del>IYKNTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> <b>Y</b> <b>I</b> <b>T</b> <b>Y</b> <b>V</b> <b>F</b> <b>L</b> <b>L</b> <b>L</b> <b>S</b> <b>F</b> <del>Y</del> LLL <del>Y</del> ISFRY
FG6e*v3	MCM <del>IYKNTIVLGGILVIIIVYLVVMSCILCFVMIVFYLIWVYCCVL</del> <b>Y</b> <b>I</b> <b>I</b> <b>L</b> <b>C</b> <b>I</b> <b>F</b> <b>I</b> <b>V</b> <b>V</b> <b>I</b> <b>I</b> <b>F</b> IIVYKFSLLD
G1ex	LLFGVLFICFVYFIVYL <del>VVVMSCILCFVMIVFYLIWVYCCVL</del> -YYRFL-VLS-LLLYISFRY
G2ex	LLFGVLFICFVYFIVYL <del>VVVMSCILCFVMIVFYLIWVYCCVL</del> -YYRFL-VLS <b>F</b> LLL <del>Y</del> ISFRY
G5ex	LLFGVLFICFVYFIVYL <del>VVVMSCILCFVMIVFYLIWVYCCVL</del> <b>Y</b> <b>I</b> <b>T</b> <b>Y</b> <b>V</b> <b>F</b> <b>L</b> <b>L</b> <b>L</b> <b>S</b> <b>F</b> <del>Y</del> LLL <del>Y</del> ISFRY
G6ex	LLFGVLFICFVYFIVYL <del>VVVMSCILCFVMIVFYLIWVYCCVL</del> <b>Y</b> <b>I</b> <b>I</b> <b>L</b> <b>C</b> <b>I</b> <b>F</b> <b>I</b> <b>V</b> <b>V</b> <b>I</b> <b>I</b> <b>F</b> IIVYKFSLLD

EATRO SDM 80 Only FGe80

MCM~~IYKNNGSCFGVFWRLGYCYCECCSFCMIIL~~

**Figure 22. Alignment of CR3 predicted protein variants from the EATRO 164 cell line.** Bolded amino acids show sequence variants and shaded sequence shows position of predicted transmembrane domains [134].

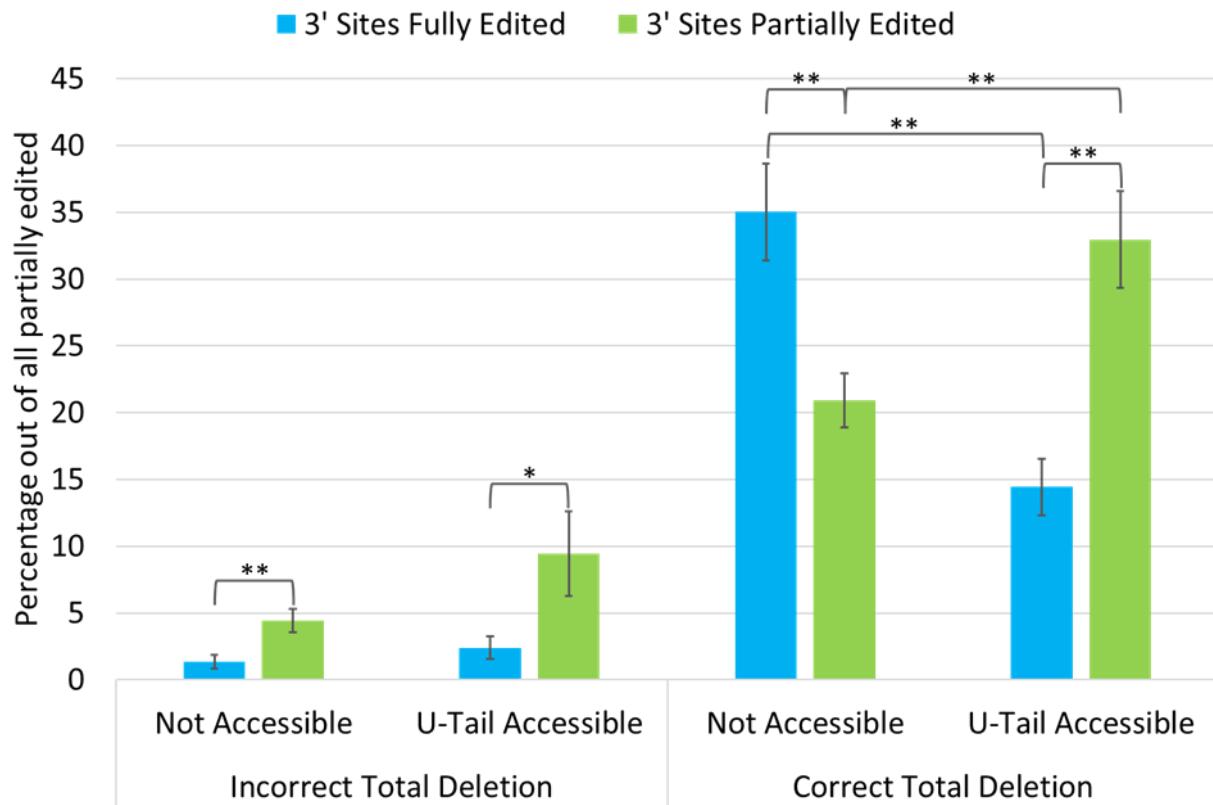


**Figure 23. Predicted secondary structures of most abundant CR3 predicted proteins.**

Secondary structure predictions were generated by RaptorX [153–155]. Shaded regions indicate predicted transmembrane alpha helices predicted by Phobius [134].

### **Deletion of Encoded Uridines Directed by the Poly-U Tail**

During the analyses of the mRNA deep sequencing data, we observed many partially edited transcripts where the deletion of encoded uridines appeared to occur early (prior to 3' insertion events). We hypothesize that the poly-U tail of preceding gRNAs may direct the removal of these encoded uridines. To examine this hypothesis, we examined partially edited RPS12 and ND7 5' transcripts with early deletion events, and determined their relative proximity to the preceding gRNA (Figure 24). (CR3 was excluded from this analysis due to the variability in gRNA location caused by the multiple branching editing pathways). These analyses revealed that for both RPS12 and ND7 5', encoded U sites that are near the poly-U tail of the preceding gRNA (U-Tail accessible sites) have a higher frequency of total deletions in the partially edited transcripts. This suggests that the proximity of the preceding gRNA's poly-U tail does impact deletion of encoded uridines and supports our hypothesis that poly-U tails can guide the deletion of encoded uridines.



**Figure 24. Frequencies of early total deletions of DNA encoded uridines in partially edited ND7 and RPS12 transcripts.** Editing sites with DNA encoded Us were identified, and each encoded U site was then characterized based on the editing access of the preceding gRNAs' poly-U tail (U-Tail Accessible or Not Accessible) as well as whether the site in question was totally deleted in the final fully edited sequence (Correct Total Deletion or Incorrect Total Deletion). For each site, a window was defined consisting of the 6 sites upstream and 6 sites downstream of the site in question. Each editing sequence was then classified based on the states in the 3' end of the window as either matching the fully edited sequence or not. Error bars depict standard error. (\*= $p<0.05$  \*\*= $p<0.01$  unpaired T-test)

## Discussion

In this work, we developed a new transcriptome analysis pipeline to fully characterize the editing pathways for three putative dual-coding genes, RPS12, ND75' and CR3. The pipeline uses a new program, SKETCH, in combination with our database search program [9]. Combining these two programs allowed us to separate true alternative edits from partially edited transcripts and allowed the precise mapping of the full progression of the editing process. This

characterization was done in two different cell lines (TREU 667 and EATRO 164) and under different energy conditions in order to determine the robustness of the editing process. Surprisingly, distinct differences in both editing progression as well as editing efficiency were observed in the two different cell lines. In addition, growth of parasites under different energy conditions also appeared to be able to influence the editing process. In both cell lines, the editing process appeared to be very inefficient, with most of the transcripts completely unedited. A comparison of the two cell lines grown in SDM79 did suggest that overall, the TREU 667 cells were more efficient in editing these three pan-edited transcripts. However, when the EATRO 164 cells were transferred from a high glucose medium (SDM79) to a glucose-restricted medium (SDM80), the number of transcripts that initiated the editing process more than doubled. For RPS12, the increase in editing initiation resulted in a 4-fold increase in the number of fully edited and translatable mRNAs.

Of the three transcripts characterized, the essential RPS12 showed the most robust editing progression. Editing of RPS12 is relatively linear, with only a few minor branching alternatives. For this mRNA, the first start codon found on the fully edited transcripts consistently translated into the canonical RPS12 open reading frame and we found no evidence of transcripts that access the alternative reading frame. The most prominent alternatively edited branch was only observed in the EATRO 164 cell line and causes a frame-shift that extends the reading frame at the 3' end (Figure 14 B2<sup>e</sup>). Interestingly, this same alternative edit was previously described by Simpson et al. in the 29-13 strain, which shows that this edit is not an isolated occurrence in the EATRO 164 strain [148]. They observed the alternative at a low abundance, which agrees with our observations as well. Their analysis identified a large

amount of variance at the 5' end, with 5.7% of transcripts being translatable in the canonical ORF. We did observe sloppy editing at the 5' end, but found two primary forms 5' end RPS12 editing. Interestingly, this data indicates that the 29-13 strain has a similar editing efficiency to the TREU 667 strain and EATRO 164 strain grown in SDM80.

In contrast to RPS12, we found distinct evidence that ND7 5' is dual-coding. In both cell lines, alternative editing by different terminal gRNA variants resulted in transcripts with either RF1 (the canonical ND7) or RF3 (a putative metabolite transporter) linked to the first AUG [141]. Interestingly, ND7 5' has also been sequenced in 29-13 cells [148]. While that study did not directly state evidence of dual-coding, they did indicate that a large proportion of the fully edited ND7 5' transcripts had a single nucleotide difference in the 5' UTR. This difference could very well be the same difference we observe in E2 transcripts that links an upstream AUG to the ARF. This suggests that the ability to access the two different reading frames is maintained across a number of different trypanosome cell lines.

While fully edited ND7 5' transcripts were found in both cell lines, a major difference was observed in the efficiency of the editing process. In TREU 667 cells, over 79% of ND7 5' transcripts had initiated the editing process and a full 9.7% are fully edited. In contrast, EATRO 164 cells grown under the same conditions (SDM79) had only 52.4% transcripts that initiated editing and a scant 0.2% fully edited. Growth of EATRO 164 cells in SDM80 did substantially increase the number of transcripts that had initiated RNA editing (72.6%), however, no corresponding increase in fully edited transcripts was observed. The major differences in editing efficiency appear to be due to both the use of alternative gRNAs that could disrupt the editing cascade and well as a gRNA mutation that affected the ability of the guide RNA to

efficiently anchor. Surprisingly, the gRNAs that disrupt editing in the EATRO cell line are also present in the TREU gRNA transcriptome. It is unclear why we see evidence of their use in only the EATRO cells. It may be that in the TREU cells, these gRNAs are more efficiently used in a different editing pathway. A full understanding of gRNA selection and use will require the characterization of the entire edited transcriptome. In addition to the large decrease in the efficiency of ND7 5' editing observed in the EATRO cells, we also saw a distinct shift in the number of fully edited transcripts that translate in RF3, the alternative open reading frame. This alternative protein has been previously predicted to be a metabolite transporter as it shares distant homology with a bacterial sugar transporter, SemiSWEET [141].

The most pronounced differences between the two cell lines was observed for the CR3 transcript. In both cell lines, CR3 utilizes a much more complicated editing pathway than either ND7 5' or RPS12 and the overall efficiency of the editing process is very low. Surprisingly, the number of CR3 transcripts that initiate RNA editing is comparable to the percentage observed for RPS12. However, editing by the initiating gRNA appears to be very inefficient. In TREU cells, while 15.1 % of the transcripts initiate editing, only 2.2% are fully edited through the first editing block. A similar drop is also observed in the EATRO cells. The identified gRNA population that initiates editing does not contain any mismatched base pairs and it is unclear why full editing by this gRNA is so inefficient. The canonical CR3 is a putative NADH Dehydrogenase complex I member (ND4L) [120]. Editing of Complex 1 members does appear to be developmentally regulated, with full editing only observed in the Bloodstream stage. [5,27–30] . It may be that editing is stalled right after initiation by a transcript specific mechanism. However, a small percentage of transcripts are still edited. Transcripts edited to the canonical

CR3 sequence were only observed in the TREU cell line. In this cell line, four different 3' editing pathways converge to an internal consensus sequence which then diverges again near the 5' end, generating a variety of different proteins in two different reading frames. In the EATRO cell line, editing initiates with the same 3' gRNAs, but diverges at the internal consensus sequence when they employ a different set of gRNAs for full editing. Because of how different their editing pathways are, the TREU and EATRO protein products cannot be directly compared. Searches were run on various different databases in order to determine the putative functions of the many CR3 proteins. Unfortunately, these searches yielded few significant results, with most proteins only sharing homology with the transmembrane domains of many different proteins. The very small percentage of CR3 transcripts that undergo full editing suggests that the protein products may not be made or utilized in this stage of the parasite life cycle. However, we hypothesize that the ability to alternatively edit transcripts may be an important evolutionary mechanism to maintain genetic plasticity.

The dual host life cycle of *T. brucei* leaves it vulnerable to genetic drift especially in regards to the mitochondrial ETC genes that are not always under selection. Previously, we proposed a mechanism that would contribute to the drift robustness of these mitochondrial genes. By overlapping ETC genes not under selection in the bloodstream stage with genes that are under selection during this stage of the life cycle, the accumulation of mutations can be prevented [141]. These overlapped genes share most gRNAs, and this strategy ensures that almost all of the genetic material is protected. We also hypothesize that the sequential nature of gRNA use and the sensitivity of the RNA editing process to both mRNA and gRNA mutations can also protect against genetic drift by increasing the deleterious effects of the mutations

(LaBar and Adami 2017). Increasing the lethality of mutations would insure that deleterious mutations are purged from the population during long periods of growth in the mammalian host. While the process of RNA editing may help weed out mutations by making them lethal, it would also prevent the population from generating beneficial mutations as well. This strategy leaves organisms no options for evolving. We suggest that alternative edits, such as those seen in the CR3 and others previously observed, editing pathways generate protein diversity without compromising the genetic information found within the genome [68].

Our analysis revealed a number of details about the larger mechanisms of gRNA selection in RNA editing. We found a surprising number of gRNAs with that had been identified to edit two different mRNA sequences. While some of these promiscuous gRNAs appear to be unproductive, generating dead end branches on their editing pathways, many appear to be productively used, as with the editing pathways of CR3 in the EATRO 164 cell line. The majority of these editing pathways are directed by promiscuous gRNAs and these pathways generate translatable transcripts. Interestingly, most of these gRNAs were identified to edit members of complex I, particularly ND8 and the 3' editing domain of ND7. It may be possible that this is another mechanism of increasing the drift robustness of *T. brucei* by giving these gRNAs multiple functions. Promiscuous gRNAs have been previously identified editing *L. tarentolae* RPS12 and ND3, however, they were not shown to be producing translatable transcripts [156,157].

In addition to this, we determined that RNA editing is not strictly sequential. While overall, editing proceeds from the 3' to 5' across of the editing domain, we have found evidence that shows that gRNAs can overwrite editing that has been previously generated.

These observations show that the RNA editing system can tolerate some amount of abnormal editing, despite the fragility of the system as a whole. Interestingly, another study that examined the use of alternative RPS12 gRNAs previously identified found that the three gRNAs in question were being utilized, and that despite not generating the conventional editing patterns, a very small number of transcripts returned to the canonical editing pattern after the alternatives [148]. This last observation may be another instance of gRNA overwriting. Non-sequential editing is also suggested by our proposal that the poly-U tail can be used to direct editing. We showed that deletions in partially edited transcripts were more common in regions close to the U tail of the preceding gRNA (Figure 24). Interestingly, in some cases, an upstream deletion site would be completely deleted while a downstream site would not be deleted. This phenomenon was also observed in another study, in an *in vitro* editing system [158].

In this analysis we determined the editing pathways of RPS12, ND7 5', and CR3 in EATRO 164 and TREU 667 cells, and we found evidence of dual-coding in ND7 and CR3. We also showed that editing patterns can vary quite significantly between cell lines and based on available energy sources. Using the gRNA transcriptomes to validate alternative edits was vital to completing this project. In light of the extreme variation we observed in the CR3 editing pathway, we believe that in order to fully understand the dynamics of the editing pathways as a whole, we need to sequence all of the edited mRNAs and gRNAs in multiple cell lines under multiple conditions.

## Acknowledgements

We thank the Ken Stuart Lab for the trypanosome cell lines. We would also like to thank Hanyou Pan for his contributions to the data analysis of ND7 5'.

# CHAPTER 5: CLUSTER CLASSIFICATION OF UNKNOWN GRNAS

## REVEALS THE ROBUSTNESS OF THE RNA EDITING SYSTEM

### Abstract

The RNA editing system of *Trypanosoma brucei* uses small RNAs called guide RNAs to direct the insertion and deletion of mRNAs. Thousands of gRNAs are used in this system to render twelve mitochondrial mRNAs translatable. These gRNAs edit sequentially, with each gRNA generating the anchor binding sequence for the next gRNA. This means that the system is inherently fragile. gRNA transcriptomes have been generated and gRNAs were identified based on their complementarity to previously described edited mRNAs. This method was effective in identifying many gRNAs, but we found that many were still unidentified. To determine if these gRNAs were nonfunctional mutants or potentially had undescribed functions, we grouped all gRNAs into clusters, where each cluster had significant sequence conservation, using our new program ACORNS. This showed that most unidentified gRNAs were not related to any functionally known gRNAs and could be generating unobserved alternative editing patterns.

Recently, the editing pathways of three genes, RPS12, ND7 5' and CR3 were described in detail. Each gene had branches in their editing pathways, but it was not always clear based on gRNA abundance data alone why one branch of the pathway was more expressed than another. Using the defined gRNA clusters, we screened all related members editing these three genes against their targets to determine what proportion of each family was able to productively edit, using another new program GUIDE. We found that most of these unidentified gRNAs were predicted to disrupt RNA editing. However, using this information in combination with the

mRNA data for these three genes, we found that these mutations are highly tolerated by the editing system. We also determined that gRNA abundance does not correlate with the mRNA editing preference. We also observed many gRNA populations of high abundance that were apparently not used to edit. In most cases these populations had no issues in their mRNA alignments but were seen to be only used in one cell line and not the other. Currently, the complete mechanism of gRNA selection is unknown.

## Introduction

The RNA editing system of the trypanosomes is a unique and complex system. It requires the use of two genetic components, the protein coding genes whose transcripts require editing, and the small RNA genes encoding the guide RNAs (gRNAs) that direct the specific edits [7]. The genes that require editing are all mitochondrially encoded and are either associated with the electron transport chain or the mitochondrial ribosome. In this RNA editing process, mRNA transcripts can be dramatically altered by the insertion and deletion of uridines [4,6]. These editing events are catalyzed by a multi-subunit editosome complex [140]. Currently, over 47 proteins have been identified that are involved in the cleavage, uridylyl addition or deletion and subsequent re-ligation events that are required for every nucleotide change [140]. Formation of these specialized complexes as well as performing the hundreds of edits required is highly energetically expensive. The RNA editing process is also incredibly fragile to mutations, as it is sequentially dependent. RNA editing starts at the 3' end of the mRNA, and each gRNA generates the anchor for the next gRNA. A mutation in any gRNA could disrupt formation of the next anchor, halting the editing process and with it abort the expression of the protein.

It has been hypothesized that this energy intensive process evolved in response to the parasites complex life cycle [75,76,141]. A dixenous kinetoplastid, *Trypanosoma brucei*, undergoes substantial shifts in its energy sources over its life cycle, requiring extensive regulation of metabolic pathways. In the nutrient deprived environment of the insect, it relies primarily on metabolizing amino acids to drive the Krebs cycle and Oxidative Phosphorylation. Once transferred to the glucose rich bloodstream of its mammalian host however, it shifts from relying on mitochondrial respiration, and moves to using glycolysis alone. In addition, the exclusively bloodstream nature of *T. brucei* in the mammal host requires that they replicate continuously, escaping the host's immune response by periodically switching their surface glycoproteins [12]. These conditions should make the genes involved in mitochondrial respiration particularly susceptible to genetic drift. Recently, we showed that as many as six of the mitochondrial transcripts may be dual coding and that alternative editing near the 5'-end of the transcript can be used to access both open reading frames (ORFs). We hypothesized, that the alternative ORF may give the parasites a selective growth advantage in the mammalian host. This would significantly increase the deleterious effects of mutations and protect against genetic drift. The linking of an ETC gene with an essential gene has been previously shown for cytochrome oxidase III (COIII). For this transcript, alternative editing by one gRNA can link an ORF found in the unedited 5' sequence with the trans-membrane domains found in the fully edited carboxyl-end of the protein [67]. This alternative protein, AEP-1, is involved in mt DNA maintenance and appears to be essential during bloodstream growth [68]. The detection of such a large number of dual-coding genes (seven including COIII) in a genome that encodes only 17 proteins, suggests that this is an important mechanism to protect the integrity of the ETC

genes that are required only in the insect. Thus, the ability to overlay genetic information maybe a protective strategy made possible by RNA editing [99,141]. The changes in metabolism are reflected in other changes in the pattern of RNA editing found the different life cycle stages in *T. brucei* [5,24–30]. For instance, many of the transcripts encoding members of NADH dehydrogenase (Complex I) are only fully edited in the bloodstream stage. In contrast, cytochrome oxidase II (COII, Complex 3) and cytochrome B (Complex 3) are both preferentially edited in the insect stage. Recently, we even observed a shift in the editing pattern found in insect stage trypanosomes if the growth medium is switched from glucose rich to the glucose depleted. This suggests that editing can respond to its environment (Chapter 4). In addition, multiple studies have shown that different editing patterns can be observed for transcripts of the same mRNA, with edits being directed by multiple alternative gRNAs. The most dramatic example of this is seen in the putative NADH dehydrogenase subunit 4L, known as C-rich region 3 (CR3) [120]. For transcripts of this gene, two different sets of gRNAs were used in two different cell lines of *T. brucei*. Both editing pathways were complex and possessed multiple branches that were fully translatable. Curiously, the gRNAs required to generate the editing pathways of the two different cell lines were present in both cell lines, but apparently not selected for use for unknown reasons.

The gRNAs that direct the RNA editing process are stored on 1 kb circular DNA molecules called minicircles. These minicircles make up the bulk of the mitochondrial DNA of *T. brucei*, with up to 10,000 minicircles per cell, and there are an estimated ~1200 different sequence classes of minicircle at varying abundances [8,33]. gRNA transcriptomes have been generated for the insect and bloodstream stage trypanosomes [9,119]. Analyses of these

transcriptomes identified over 600 gRNA populations (gRNAs that edit the same region of the same gene) involved in editing the mitochondrial mRNA transcripts. These populations typically contained multiple different sequence classes of gRNAs (major classes). The sequence differences observed are all R to R or Y to Y changes, and since both A:U and G:U base pairs occur in the editing process, the multiple sequence classes can all guide the generation of the same mRNA sequence. These data also show extreme quantitative differences between the different gRNA populations, with population sizes ranging from <10 to >300,000 reads. We had postulated that regions with very low gRNA coverage are areas with mRNA sequence variations and that the editing of these regions may be directed by gRNAs not identified due to internal mismatches with conventional sequence. In addition, we also identified some abundant, mutated gRNAs that could potentially introduce frameshifts or create sequence that disrupts the upstream anchor binding site. This was surprising, as we had hypothesized that these mutations would be selected against due to the fragile nature of the editing process. Because of the high stringency of our initial screen, only gRNAs with mismatch mutations biased towards the 3' or 5' ends of the gRNAs were initially identified and it was unclear how prevalent mutated gRNAs that could disrupt editing were. Analyses of our gRNA libraries did indicate that they contain millions of “unidentified” reads that had key guide RNA characteristics (characteristic transcription start sites, U-tail and ~ length).

In this manuscript, we describe the development of two new pipelines that allows us to begin to characterize these “unidentified” transcripts in order to determine what impact they might have on RNA editing efficiency. ACORNS (Assemble Clusters Of Related Nucleotide Sequences), allows us to identify and classify gRNA sequences. A second program, GUIDE (gRNA

Uridine Insertion/Deletion Editor), simulates the RNA editing process and allows us to predict the effect of the gRNA mutation on the sequencing process. Using the CR3, RPS12 and ND7 5' mRNA transcriptome data, we identified the gRNA populations that can edit or disrupt editing of these genes. We found a surprisingly high toleration of mismatches and gaps in gRNA/mRNA alignments, as well as many instances where the most abundant gRNAs present were not those preferentially used. A deeper look at the gRNA transcriptomes also revealed the presence of surprisingly many unidentified and functionally unknown gRNAs, suggesting that more work is needed to discover their roles in the RNA editing system.

## **Materials and Methods**

### **Cluster analysis of related gRNAs by ACORNS**

In order to determine the relationships of previously identified and unknown gRNAs, a new program was created called ACORNS (Assemble Clusters of Related Nucleotide Sequences). This program functions to identify putative gRNAs, determine relationships between gRNAs, and group related gRNAs into clusters.

**Identification of putative gRNAs.** As described previously, identical sequences were collapsed and sequences without four consecutive Ts were filtered out, indicating the lack of a poly-U tail [9]. To perform this analysis, a new program was generated, called Assemble Clusters Of Related Nucleotide Sequences (ACORNS). ACORNS first filtered out all putative gRNAs, based on two criteria: having 40 nucleotides prior to the start of the poly-U tail, and having a transcription start site that either matches one of the top twenty most common six nucleotide gRNA transcription start sites, or is one mutation away from one of these sites [9,119]. Maxicircle edited and unedited sequences as well as ribosomal RNAs were also filtered

from the sequence file. Identical sequences were collapsed, but retained their overall total read abundance. Additionally, if sequences were identical with the exception of the start position of the poly-U tail, or 5' end location of the gRNA, they were still consolidated, keep the sequence of the most abundant transcript.

**Identification of gRNA families.** A pair of related gRNAs are defined as two gRNAs that only differ by a single substitution, insertion or deletion mutation. Once putative gRNAs were identified, ACORNS aligned each gRNA with every other gRNA to determine which gRNAs were related. Alignments were scored using the python-levenshtein package [159]. In order to prevent errors due to 5' exonuclease activity or difference in poly-U site, overhanging nucleotides on either end of the alignment were not counted as mismatches.

ACORNS then grouped gRNAs into families based on their relationships. Each family was grouped together by starting with the most abundant transcript, and including any other gRNAs related to that transcript, as well as any gRNAs related to those gRNAs, and so on. Once gRNAs were grouped into families, sequences were trimmed to the same length and any new identical sequences were collapsed.

Visualization of these related gRNA clusters was performed by a subprogram, and clusters were visualized in two different ways, with color coding based on gRNA identity or gRNA abundance. Identity was defined as the gene known to be edited by previously identified gRNAs, with other all other gRNAs labeled as “unknown”. Color coding by abundance was based on a log scale, and each scale created for each cluster, with the least abundant gRNA being colored purple and the most abundant gRNA being colored red, and all other gRNAs being scaled accordingly.

## Prediction of RNA editing by GUIDE

A reference library of all edited forms previously described of RPS12, ND7 5', and CR3 were collected and annotated with the editing states in each transcript (Chapter 4). For this analysis we generated a new piece of software known as the gRNA uridine insertion/deletion editor (GUIDE). GUIDE uses the outputs of ACORNS, and for each gene, GUIDE pulls the families of gRNAs that have members that had been previously identified to edit the given gene. Editing of each member of each family was then simulated. For each family, the template generated by the previously identified members of that family were used. If gRNAs generating different forms of the same mRNA were in the same family, all appropriate templates were tested on each family member, and the predicted edit with the longest alignment with the gRNA was saved.

In order to predict the edits for each gRNA, GUIDE determines the most likely anchor binding location on the template. Anchor sequences were required to use Watson-Crick base pairs only and be located in the first twenty nucleotides (nt) of the gRNA. The first twenty nt of each were scanned along the template, and the longest consecutive stretch of Watson-Crick base pairs was identified. Additionally, a second anchor was identified, using a weighted score, where each G:C base pair was worth two points and A:U base pairs were worth one point. Editing from both anchors were tested and the anchor that generated the longest edited sequence was saved.

The rules of editing used by the program are the standard uridine insertion/deletion rules observed in the editing system. If a nucleotide from the gRNA is aligned to a nucleotide from the mRNA that is “illegal” (anything other than Watson-Crick or G:U base pairs), a uridine

will either be inserted or deleted to resolve the mismatch. If this cannot resolve the mismatch, editing ends. Once edits were predicted, they were classified based on their effect on the editing process as a whole and their effect on the predicted protein.

## **Results**

In our characterization of the EATRO 164 gRNA transcriptome, we identified over 3 million reads and ~64,000 unique gRNA sequences capable of generating conventional editing patterns [9]. These gRNAs were identified by finding the longest common substring to conventionally edited mRNAs and retaining only those scoring 45 or more (Watson-Crick base pairs = 2; G:U base pairs = 1). In these studies, we found that lowering the stringency and/or allowing for mismatches lead to the identification of thousands of gRNAs with characteristics suggesting that they were misaligned. While this initial analysis did identify an almost full cohort of guide RNAs, there were still millions of reads in the transcriptome that could not be classified. To determine how many of these reads were unidentified gRNAs for possible alternative sequence or mutated conventional gRNAs, we developed a new pipeline. ACORNS (Assemble Clusters Of Related Nucleotide Sequences) has a number of features that allows it to identify and classify putative guide RNAs:

1. It filters out all contaminating nuclear and maxicircle sequences.
2. It Identifies putative gRNAs based on three criteria: presence of a U-tail (defined as 4 consecutive U-residues), length (40 nts prior to last stretch of U-residues) and transcription start site (based on sequences defined in the previous analyses) [9,119].

3. It identifies related gRNAs by scoring the best alignment of each gRNA against all other putative gRNAs in the library. gRNAs with a single mismatch or gap are then classified as related.

4. It clusters related gRNAs by starting with the most abundant transcript, grouping all relatives of that transcript into the cluster, then adding all relatives of any relative to the cluster.

5. Using a subprogram, the clusters of related guide RNAs can be visualized, revealing the relationships between the previous identified gRNAs and unknown gRNAs.

ACORNS analyses were done for both the TREU 667 and EATRO 164 procytic gRNA transcriptomes, leading to the identification of 1256 clusters in TREU and 1168 clusters in EATRO 164 (Tables 11 and 12). This allowed us to identify all of the gRNAs that were distinctly related to a conventional gRNA but had undergone a mutational event. In addition, we also identified a number of clusters that had no previously identified members. The clusters identified varied greatly in size, with the largest cluster containing over 3400 sequence members with a total of 1,377,190 transcript reads. These large clusters (1000+ members) were actually quite rare, with only 10 and 4 clusters of this size identified in the TREU and EATRO gRNA libraries, respectively (Table 13). Interestingly, the majority of clusters were quite small, containing fewer than 25 sequence members and an average transcript number of approximately 250 reads (Table 13). The majority of the small clusters were “unidentified”, in that they contained no previously identified gRNA member (Table 13). An increase in the size of the cluster, also increased the probability that they contained a previously identified gRNA and that cluster members could then be identified.

**Table 11. Summary of ACORNS Results.** gRNAs are classified as previously identified, related to a previously identified gRNA based on sequence similarity or unrelated to any previously identified gRNAs.

	TREU 667 Procytic	EATRO 164 Procytic
Initial Reads	15,251,292	11,387,683
Final Reads after ACORNS Step 2.	11,199,364	9,049,005
Previously Identified Reads	5,121,216	4,413,142
Previously Unidentified tagged as Related	2,636,714	2,072,099
Previously Unidentified tagged as Unrelated	3,441,434	2,563,764

**Table 12. Cluster summary.** Clusters were defined as a group of 10 or more related gRNAs. Cluster characteristics describe the percent of the cluster members that had been previously identified.

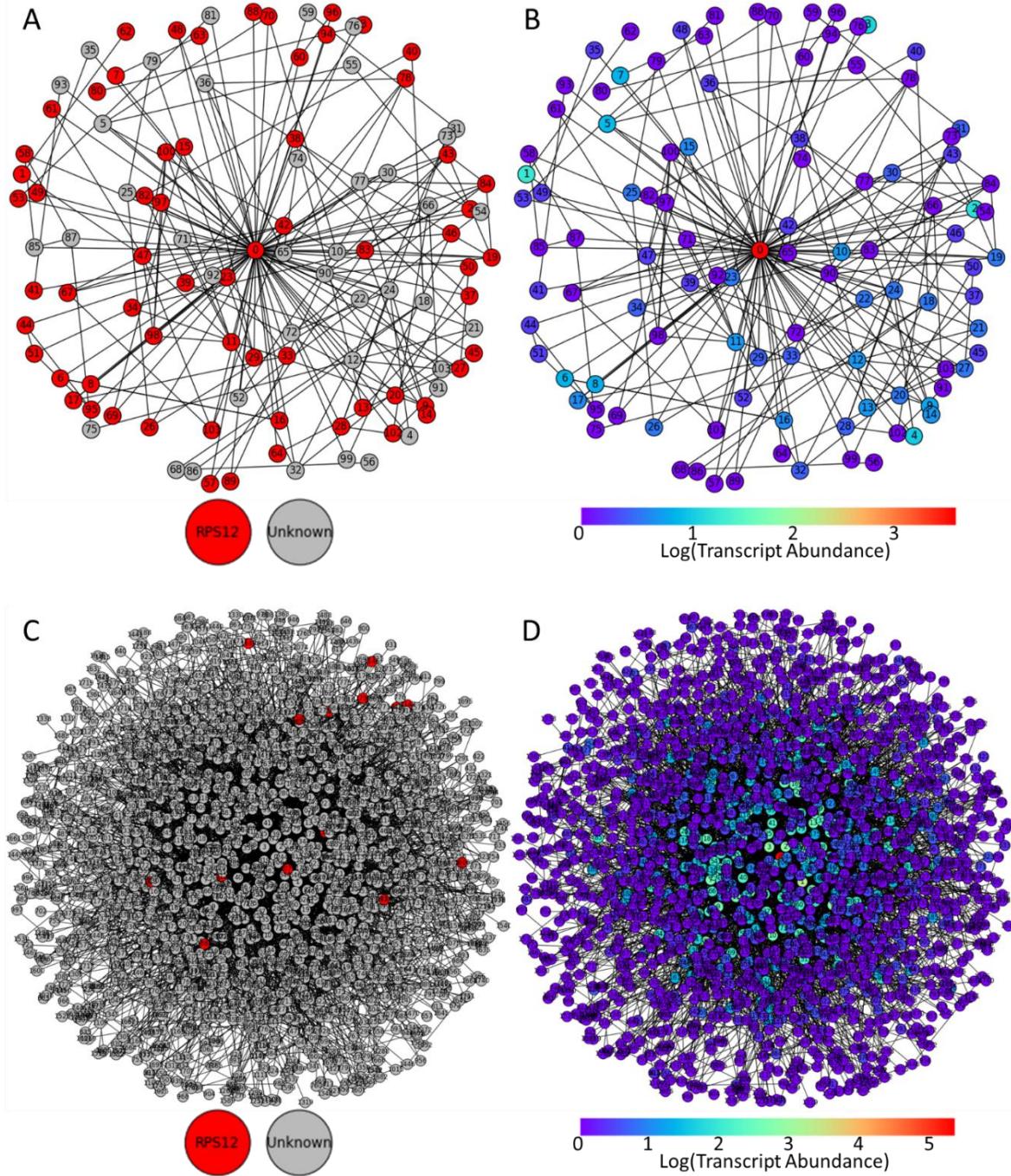
Cluster Characteristics	TREU 667 Procytic		EATRO 164 Procytic	
	# of Clusters	Reads	# of Clusters	Reads
≥95% Previously Identified	238	4,232,188	191	3,281,190
>0% Previously Identified	267	3,488,118	215	3,180,019
0% Previously Identified	751	3,289,359	762	2,401,148
Unclustered	NA	189,699	NA	186,648

**Table 13. Cluster size summary.** Cluster size is determined by the number of unique sequence classes in a cluster. % Unidentified clusters is the percentage of all clusters in each bin that are completely unidentified.

Cluster Size	# of clusters identified		Total # Reads		% Unidentified clusters	
	TREU 667	EATRO 164	TREU 667	EATRO 164	TREU 667	EATRO 164
10 - 24	518	503	137,413	134,707	68%	77%
25 - 49	291	265	240,804	211,154	63%	64%
50 - 99	185	184	427,029	425,535	55%	61%
100 - 249	173	139	1,575,518	1,398,488	45%	54%
250 - 499	60	47	2,167,031	2,153,210	45%	26%
500 - 999	19	26	2,354,552	3,379,901	21%	27%
1000 +	10	4	4,107,318	1,159,362	30%	25%

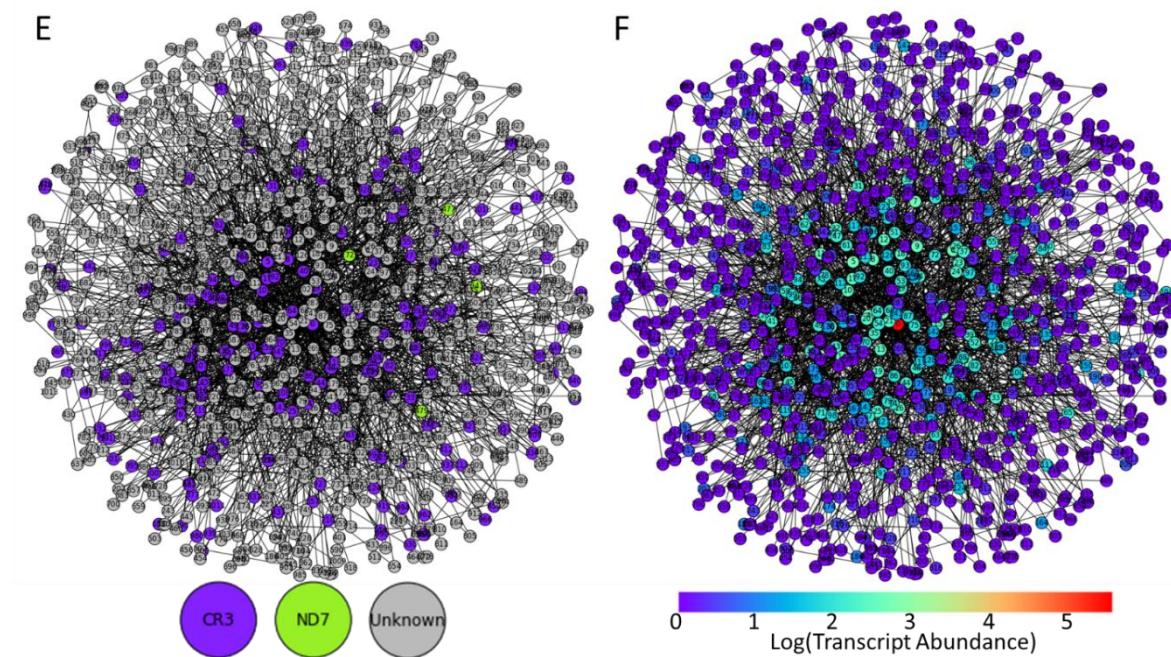
Figure 25 shows an example of three different clusters assembled by the ACORNS program. The visualization program represents each individual gRNA sequence as a dot, with lines connecting dots representing relationships between the individual sequences. Each member is numbered in order of abundance of transcript reads (0 = most abundant). gRNAs can be further characterized within the cluster by using color to identify either transcript specific gRNAs (Figure 25 A, C and E), or to indicate Log transcript abundance (Figure 25 B, D and F). The

first cluster shown illustrates a relatively small cluster with 104 different gRNA sequence members (Figure 25 A and B). In this visualization, it is clear that most of the members were previously identified (red dots), including the most abundant member (dot 0). Figure 25B shows the transcript abundance of each cluster member. This cluster is very characteristic of most of the clusters we identified, in that it has a central very abundant gRNA (red dot in center, >5,000 reads) with most other cluster members identified having very few reads (purple = 1 read in the transcriptome). 98.1% of the transcripts in this cluster were previously identified. Figure 25 C and D illustrate a large cluster with 1801 sequence members. In contrast to the first cluster, most of the identified sequence members were previously unidentified (gray dots), a few of the members however, had been previously tagged as RPS12 specific gRNAs (red dots). The previously identified RPS12 gRNAs were tagged as editing the Block H region of RPS12. The originally identified transcripts were rare, with only 255 reads found in the EATRO 164 library. The ACORNS analysis allowed us to identify an additional 276,358 reads that cover this region. However, most of these newly identified gRNAs, including the most abundant sequence class, have mutations that affect their ability to correctly edit the mRNA transcript.



**Figure 25. Example clusters of related gRNAs generated by ACORNS from the EATRO 164 PC gRNA transcriptome.** Each dot represents an individual gRNA sequence and lines connecting dots represent relationships between gRNAs. Each cluster is shown with two different color schemes; editing gene identity (A, C and E) and gRNA read abundance (B, D and F). The first two clusters (A,B and C,D) both contain RPS12 identified gRNAs (red dots) and previously unidentified gRNAs (gray dots). The third cluster (E and F) contains a majority of previously unidentified gRNAs (grey), but also contains CR3 identified gRNAs (purple) and a few gRNAs that had been tagged for a dead-end ND7 5' alternative editing pattern (green).

**Figure 25 (cont'd).**



The final cluster shown also illustrates a large cluster with 1039 members. (Figure 25 E and F). For this cluster, the most abundant member (member 0) contained over >400,000 reads and had been previously identified as an alternative CR3 gRNA ( $gC^{e80}$ ). The ACORN cluster analysis allowed us to identify an additional 829 sequences and over 12,700 reads as being related to this alternative gRNA. Interestingly, in the visualization of this cluster, we noted that 4 of the members were previously tagged as involved in the generation of a disruptive (dead-end) alternative edit of ND7 5'. An examination of these transcripts indicate that they are in fact a specific mutant subclass of the gRNA cluster that are now capable of anchoring and creating a misedit in ND7. This cluster was the only example of a cluster containing related gRNAs with different targets (distinct from promiscuous gRNAs that all have multiple targets), and could be an example of how alternative editing originates.

These data analyses suggest that both procyclic libraries have large numbers of unidentified gRNAs of unknown function. It may be that these gRNAs are directing alternative editing events that have not yet been characterized. The full mRNA transcriptome has not been sequenced and the limited amount of mRNA sequence available does suggest an abundance of alternative editing (Chapters 3 and 4) [9,67,68,99,141,148,160]. Surprisingly, both libraries also contained large numbers of mutated conventional gRNAs. This was unexpected, due to the fragile nature of the RNA editing process. We had hypothesized, that the sequential nature of the RNA editing process, would decrease the tolerance for mutations in both the gRNA and the mRNA genes (Chapter 4) [9]. In order to determine how these mutated gRNAs might influence the RNA editing process, we developed a second program (GUIDE, gRNA Uridine Insertion/Deletion Editor) that simulates the RNA editing process. This program takes the fully edited mRNA templates for the identified gRNAs in each cluster, finds the best anchor for each gRNA and then simulates editing based on the conventional base-pairing rules (A:U and G:U pairing both allowed). For each gRNA, the anchor length, length of complementarity and the number of sites showing non-conventional editing are determined. This allows the classification of the gRNA into different bins: 1) low quality anchor; 2) does not fully edit; 3) conventionally edits and 4) alternatively edits. Guide RNAs that generate alternatively edited sequence were then further classified as 1) Disruptive edits (does not generate the anchor for the next gRNA); 2) Frameshift edits (does not disrupt editing but generates a frameshift); 3) missense editing (does not disrupt editing but generates a missense mutation). We note that in the GUIDE program, we initially classified C:A base pairs as disruptive (stopping the editing process). However, because C:A base pairs have been previously identified within known gRNA

alignments, we did sort gRNAs containing a C:A base pair into their own bin [9,119]. These analyses were done for all gRNA clusters identified for RPS12, ND7 5' and CR3. These 3 genes were chosen because the mRNA transcriptomes and full editing pathways had been previously characterized (Chapter 4).

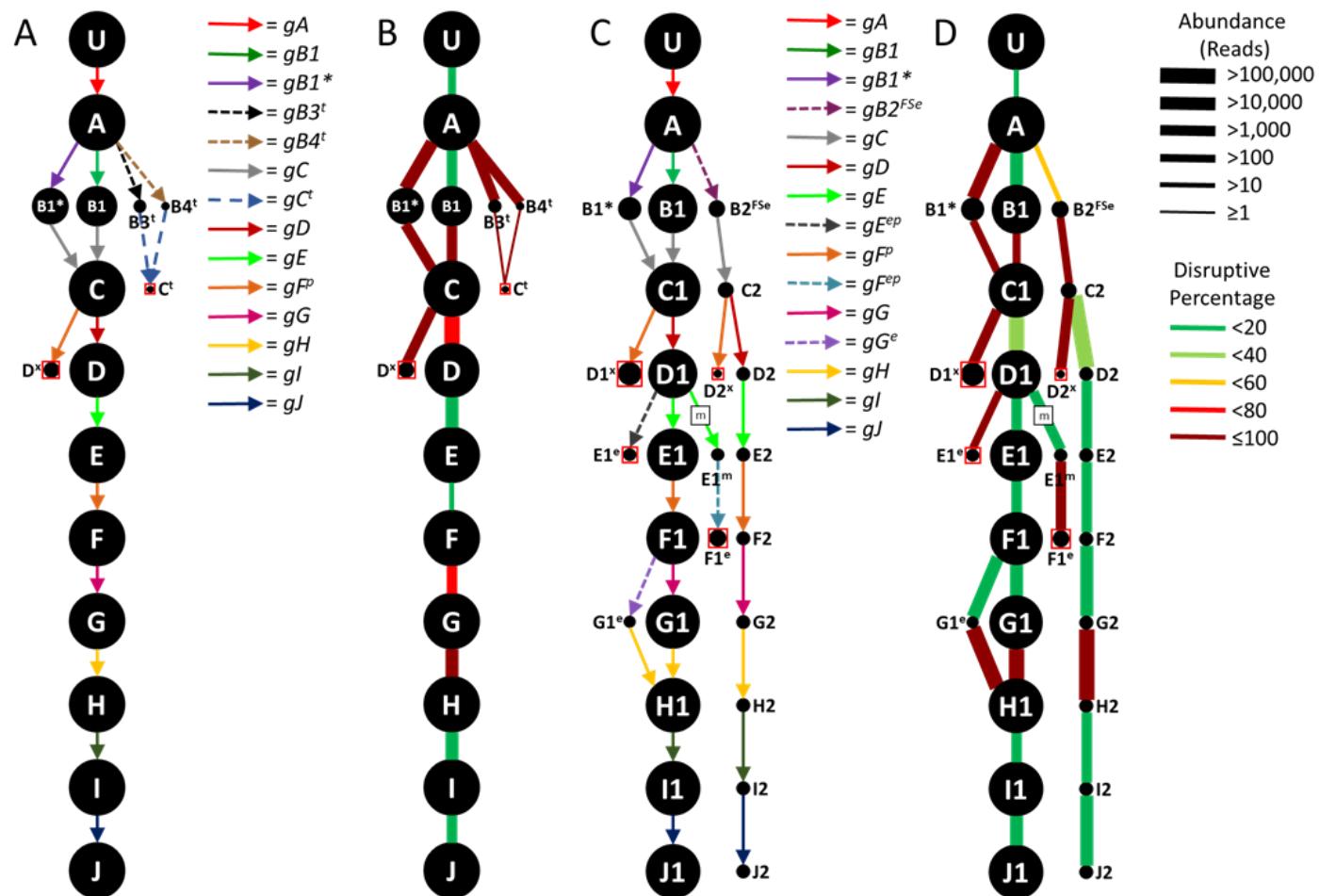
Of the unidentified gRNAs that are related to previously identified gRNAs, 942,397 and 958,140 reads were related to RPS12, ND7 5' or CR3 editing gRNAs in TREU and EATRO cells respectively. Of these newly identified gRNAs, 92.6% in TREU and 95.5% in EATRO were predicted to be incapable of fully editing, to fail to create anchor for the next gRNA or to cause a frameshift. This suggests that a large number of gRNAs could be disruptive, however, because of the large differences in population sizes for the different guide RNAs, it is important to evaluate the percent of disruptive editing at the population level. This data could indicate how tolerant the RNA editing process is to mutations in the gRNA population.

## RPS12

Editing of the essential RPS12 transcript was relatively straight forward with a limited number of alternative edits (Chapter 4, Figure 26). The main editing pathway involved 10 gRNA populations (gA – gJ). Analyses of these populations still show a large variation in the abundance of the different populations even after the cluster analyses (Figure 26, Table 14). Surprisingly, we saw no correlation between the abundance of the gRNA populations and RNA editing efficiency. For example, in the TREU cell line, the B1 and B1\* mRNA transcripts are equally abundant (Figure 26A). However, the gB1\* gRNA is almost 50-fold more abundant than gB1 (Figure 26B and Table 14). In contrast, in the EATRO cell line, the B1 mRNA transcript is significantly more abundant (5x), while the gB1 and gB1\* guide RNAs are approximately equal

in abundance. Interestingly, the only gRNA identified that can generate the B1\* edits requires a gap in the alignment with the mRNA to create the correct sequence. Our GUIDE program bins this gRNA into the “disrupts editing” bin (brown) and it is unclear why this gap is tolerated. We do note that there are three G:C base pairs surrounding the gap that would help stabilize this alignment. We hypothesize however, that the gap in the alignment does in fact affect the ability of this gRNA to efficiently edit. In the TREU cells, the abundance of gB1\* may contribute to its editing efficiency compensating for the gap and increasing the amount of B1\* mRNA observed. This pattern appears to be repeated in EATRO, where the B1 mRNA is almost five times as abundant as the B1\* mRNA, but the gRNAs are nearly equally abundant. Interestingly, in both transcriptomes the gRNAs responsible for generating the B1 edits both require C:A base pairs for the gRNAs to function correctly. This was the first of many instances observed where editing is impossible without the use of C:A base pairs.

A significant number of the gRNA populations required for the editing of RPS12 did have large numbers of mutated gRNAs that were predicted to stop (mismatch cannot be resolved by the insertion or deletion of a U-residue) or disrupt editing (does not generate the correct anchor sequence for the next gRNA). For example, while the gC populations had a small number of gRNAs with perfect alignment to the canonical sequence (1.6% in TREU and 5.9% of the gC populations in EATRO), most of these gRNAs contain an illegal base pair (G:A or G:G). Similarly, more than 70% of the gD population in the TREU cell line have a U:U mis-match in the middle of the alignment and both the gG and gH populations have a majority of gRNAs that contain a single mismatch in the best alignment with their editing block.(Figures 27 and 28). Surprisingly, when the GUIDE program simulated the edits generated by the mutants of the gD



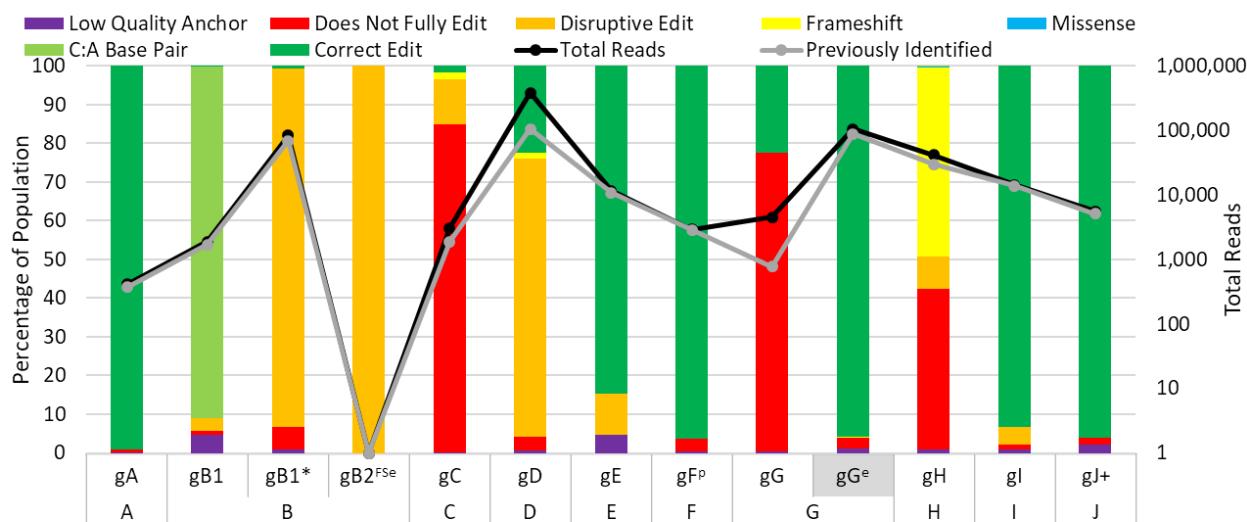
**Figure 26. Observed RPS12 editing pathways in the TREU 667 cell line (A and B) the EATRO cell line (C and D).** U = unedited transcripts. Dot sizes are proportional to the percent of block level edited transcripts using the gRNA indicated. Colored arrows indicate the gRNA population used (A and C). Dashed arrows represent gRNA populations used in only one cell line (superscript 'e' or 't'). gRNA names with superscript 'p' represent promiscuous gRNAs. Dots enclosed by a red box represent end point mRNAs with no AUG start codon. Lines connecting dots (B and D) indicate gRNA population size and functionality. Disruptive gRNAs were considered to be those that were predicted by GUIDE to be unable to complete editing (excluding those that only required a C:A base pair to finish editing), unable to generate the anchor for the next productive gRNA, or generated a frameshift mutation.

**Table 14. gRNA population analyses for RPS12.**

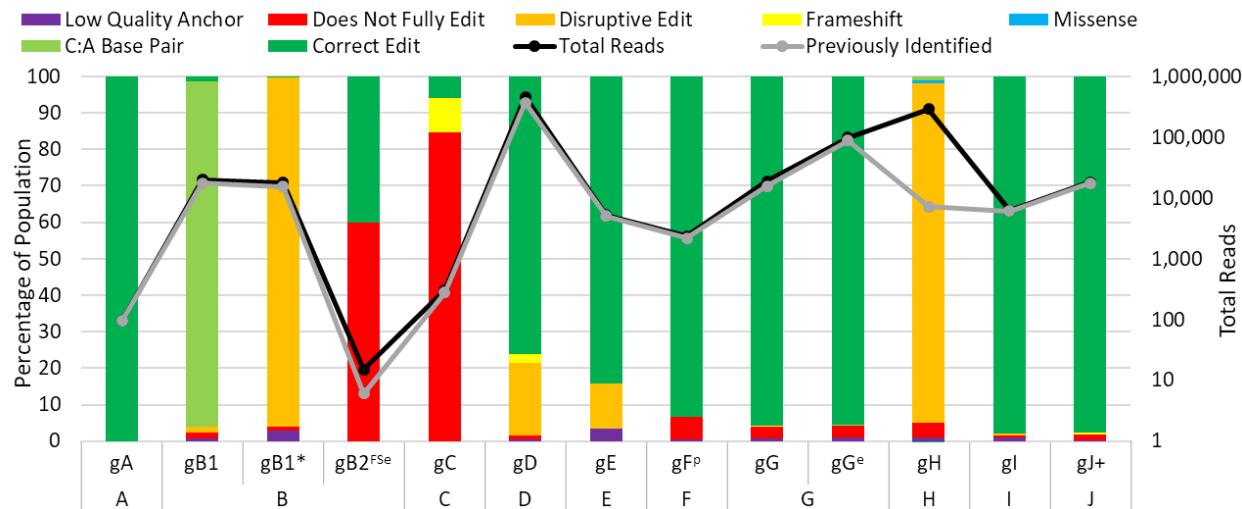
Editing Block	Population	TREU 667 gRNA Reads	TREU 667 mRNA Reads	EATRO 164 gRNA Reads	EATRO 164 mRNA Reads
A	gA	416	69,923	95	38,946
B	gB1	1,901	27,805	19,878	28,494
	gB1*	84,459	23,318	17,975	5,880
	gB3 <sup>t</sup>	1,692	2,463	0	0
	gB4 <sup>t</sup>	2,784	946	0	0
	gB2 <sup>FSe</sup>	1	0	15	2,820
C	gC	3,081	46,436	304	35,255
	gC <sup>t</sup>	6	577	0	0
D	gD	387,246	35,563	460,928	20,067
	gF <sup>p</sup> (Dx edit)	2,877	2,724	2,371	4,617
E	gE	11,593	35,046	5,269	17,723
	gE <sup>ep</sup>	1,162	0	178	879
F	gF <sup>p</sup>	2,940	29,838	2,371	10,818
	gF <sup>ep</sup>	0	0	1,843	955
G	gG	4,505	28,973	19,028	11,163
	gG <sup>e</sup>	105,217	0	99,153	468
H	gH	41,824	29,331	287,393	10,843
I	gI	14,251	26,417	6,136	9,747
J	gJ+	5,573	17,870	18,094	7,665

population, it predicted that the mutants would generate an alternative sequence that lacks the anchor binding site for the gE population (Figures 27 and 28). However, this alternative sequence was not identified in the RPS12 mRNAs of the TREU or EATRO cells (Chapter 4). This suggests that the alignment error is tolerated in the generation of the canonical sequence. For most of the mismatches detected, the error in the alignments are all immediately flanked by multiple G:C base pairs. It may be that the presence of multiple stable G:C pairs allows the editing machinery to tolerate these mismatches. The EATRO gH gRNAs are the exception. This population contains a majority of transcripts with a single point mutation near the end of the gRNA that should prevent it from generating the full anchor sequence for gI. There are no

stabilizing G:C pairs that flank this mismatch. The other notable difference between the TREU and EATRO editing patterns were found in an alternative edit ( $G^e$ ) observed in the EATRO cell line only. We note that the  $gG^e$  guide RNA is very abundant in both transcriptomes (>100,000 reads in TREU cells and >90,000 reads in EATRO cells). However, we found evidence of its use in only the EATRO cell line. In summary, we identified multiple populations of gRNAs that were predicted to be nonfunctional or disruptive in the RPS12 editing pathway. However, these predictions are contradicted by the observed RPS12 mRNA data, suggesting that many of these mismatches that we had previously considered to render gRNAs nonfunctional or disruptive are tolerated by the RNA editing system.



**Figure 27. Analysis of functionality and abundance of productive gRNAs populations that edit RPS12 in TREU 667 cells.** The functionality of each subpopulation is shown as a bar, with percentage shown on the left y-axis, and subpopulation abundance before and after identification of gRNA relatives is shown on the right y-axis. gRNAs labeled as ‘Disruptive Edit’ failed to generate the anchor for the subsequent gRNA, and gRNAs labeled as ‘C:A base pair’ required a C:A base pair to be tolerated for the editing to be completed correctly. gRNAs with shaded names were found in the TREU 667 gRNA transcriptome but were not used in editing the TREU 667 RPS12 mRNAs, despite being used in EATRO 164 cells.

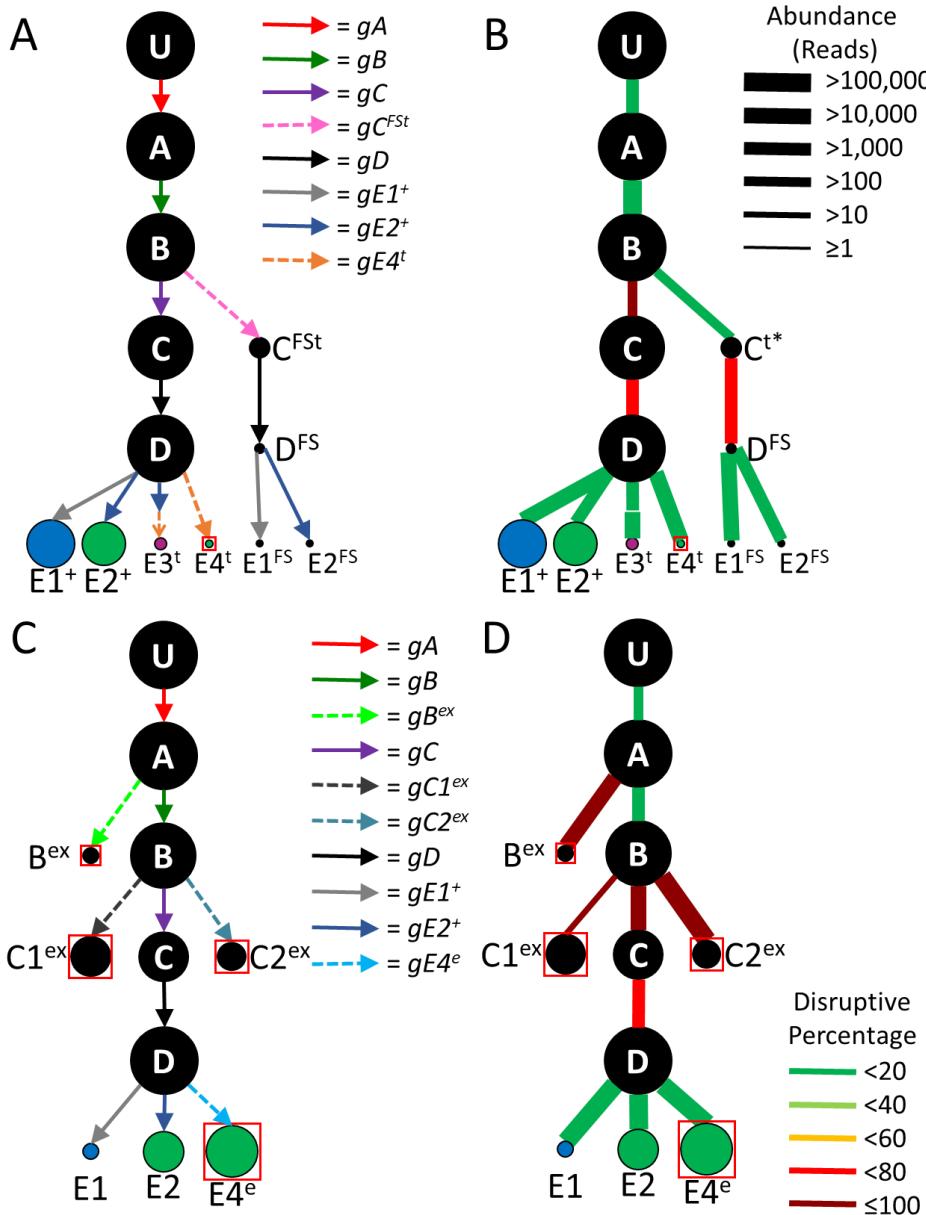


**Figure 28. Analysis of functionality and abundance of productive gRNAs that edit RPS12 in EATRO 164 cells.** For description of axes and gRNA functionality labels, please see Figure 27.

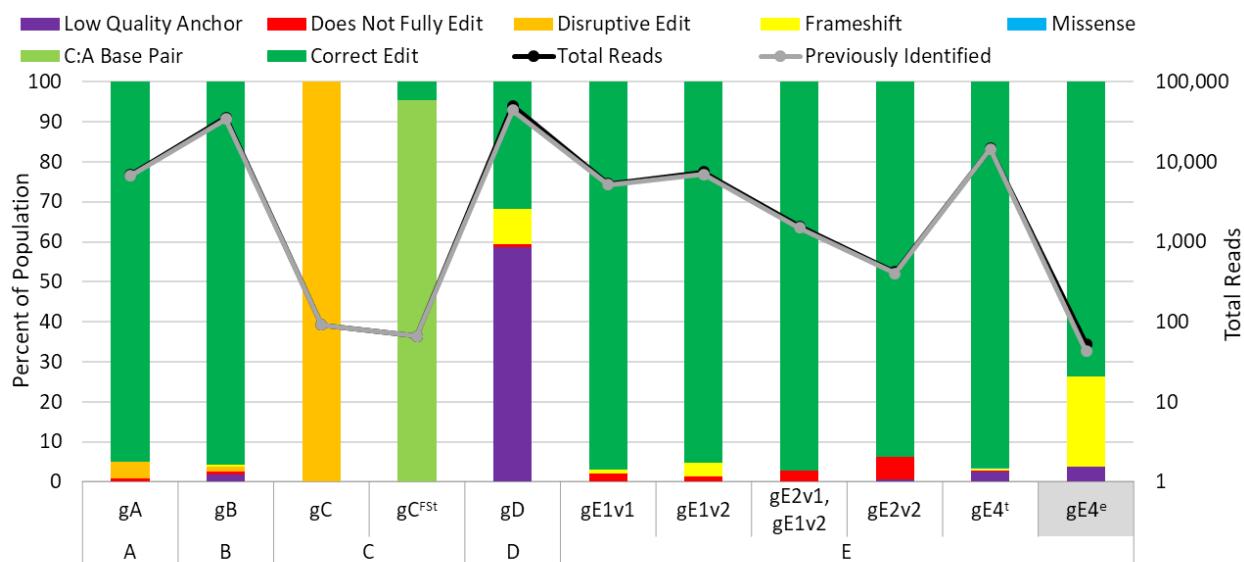
## ND7

As with RPS12, the ACORNS program identified gRNAs related to those previously identified to generate the ND7 5' editing pathways (Figure 29). The ND7 5' editing domain contains five editing block levels, and the editing pathway is relatively straight forward until the final editing block, where alternative editing generates transcripts translatable in multiple reading frames (Chapter 4). In full editing of this transcript, two of the 5 gRNA populations appear to be problematic; gC and gD (Figure 29). In both cell lines, the gC population requires both C:A base pairs as well as multiple mismatches or gaps in the best alignments with the canonical sequence (Figure 30). These multiple alignment errors do not appear to be well tolerated by the editing system, as a severe decrease in editing efficiency was observed in both cell lines from the B to C block level (25.8% in TREU cells and 32.4% in EATRO cells) (Chapter 4). For the ND7 5' gD guide RNAs, both cell lines have a majority of gRNA transcripts with a base pair mismatch in their anchor. In the EATRO cells, almost all gRNAs in the gD population either

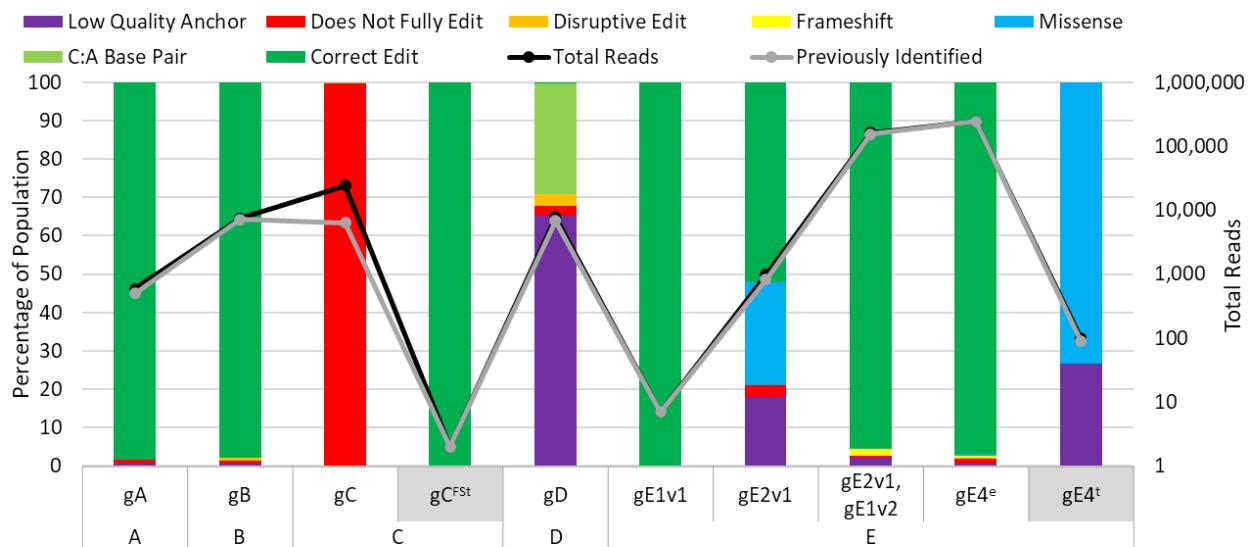
require a C:A or A:A base pair in the anchor (Figure 31). While these mismatches are surrounded by G:C base pairs, editing efficiency does appear to suffer with a 13.0% drop in editing efficiency and only 0.4% of transcripts completing the D block level (Table 15). In contrast, in TREU cells, while most gRNAs have a A:A mismatch in the anchor, ~31% of the population has the ability to form a conventional Watson-Crick anchor (Figure 30). The drop in editing efficiency is less (7.2%) than observed in the EATRO cells, and a full 11.7% of transcripts complete D block level editing. In these two problematic gRNA populations of ND7 5' editing, we begin to see what the limits of the RNA editing system are. We find that multiple mismatches as well as mismatches in the anchor binding region of a gRNA appear to have severe impacts on editing efficiency. Surprisingly, however, these aberrant gRNAs do not appear to halt editing altogether.



**Figure 29. Observed ND7 5' editing pathways in TREU 667 cell line (A and B) and the EATRO 164 cell line (C and D).** For descriptions of dots and arrows, please see Figure 26. Dashed arrows with open heads represent a hypothetical rewrite. + indicates that more than one mRNA form was condensed into this circle to simplify the figure. Condensed forms encode largely the same amino acid sequence with only small variants. Terminal dots are colored blue for reading frame 1, magenta for reading frame 2, or green for reading frame 3. Boxed green dots have no functional start codon but are translatable into reading frame 3 with the use of an alternative start codon (UUG). Lines connecting dots (B and D) indicate gRNA population size and functionality. Disruptive gRNAs were considered to be those that were predicted by GUIDE to be unable to complete editing (excluding those that only required a C:A base pair to finish editing), unable to generate the anchor for the next productive gRNA, or generated a frameshift mutation.



**Figure 30. Analysis of functionality and abundance of productive gRNAs that edit ND7 5' in TREU 667 cells.** For description of axes and gRNA functionality labels, please see Figure 27. gRNAs with shaded names were found in the TREU 667 gRNA transcriptome but were not found to be utilized in editing the TREU 667 ND7 5' mRNAs, despite being utilized in EATRO 164 cells. The population labeled gE2v1,gE1v2 contained members that generated both sequence patterns and could not be separated.



**Figure 31. Analysis of functionality and abundance of productive gRNAs that edit ND7 5' in EATRO 164 cells.** For description of axes and gRNA functionality labels, please see Figure 27. gRNAs with shaded names were found in the EATRO 164 gRNA transcriptome but were not found to be utilized in editing the EATRO 164 ND7 5' mRNAs, despite being utilized in TREU 667 cells. The population labeled gE2v1,gE1v2 contained members that generated both sequence patterns and could not be separated.

**Table 15. gRNA population analysis for ND7 5'**

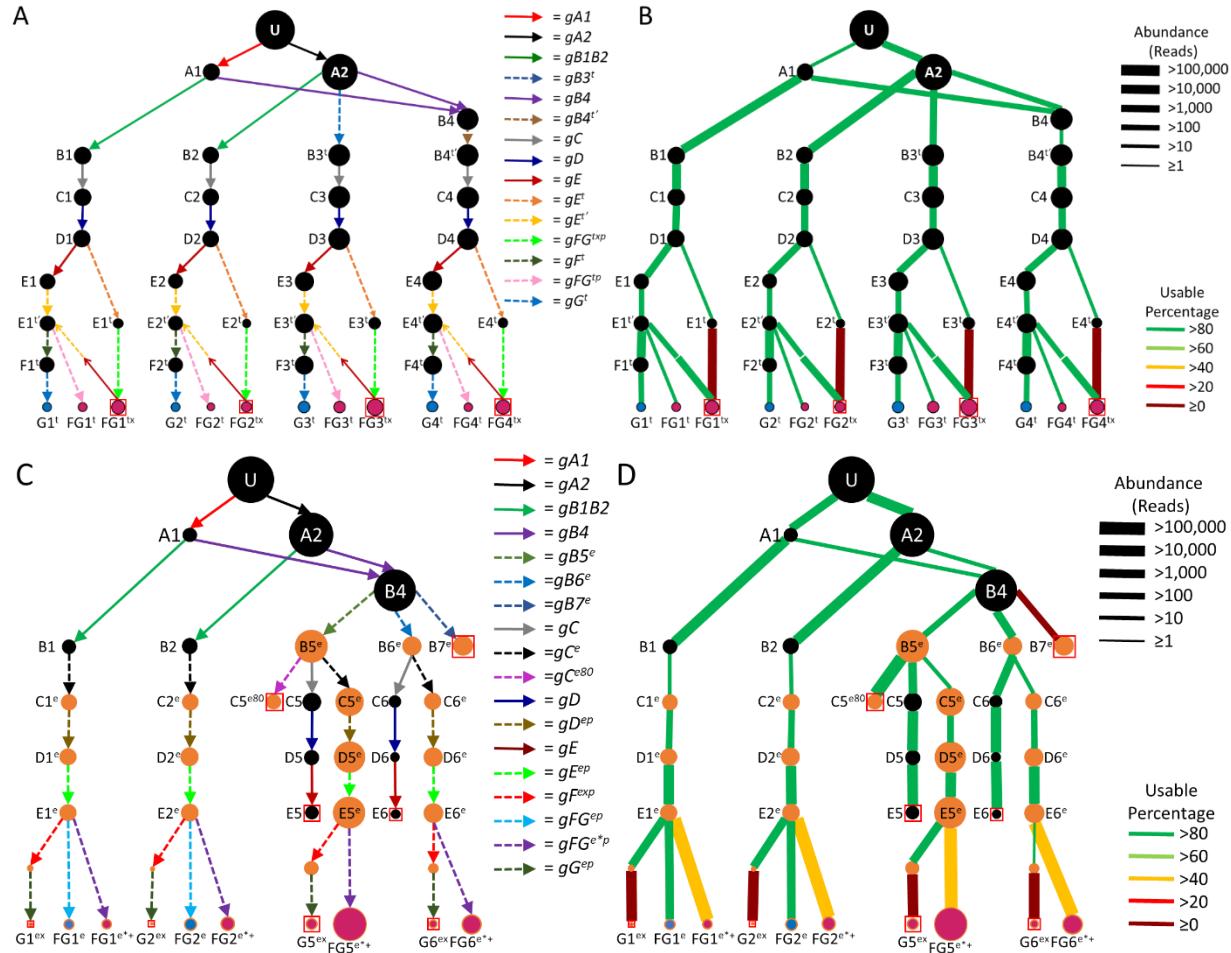
Editing Block	Population	TREU 667 gRNA Reads	TREU 667 mRNA Reads	EATRO 164 gRNA Reads	EATRO 164 mRNA Reads
A	gA	6,819	523,128	589	430,266
B	gB	35,890	510,653	7,478	395,621
	gB <sup>ex</sup>	113	0	21,775	23,440
C	gC	92	196,504	24,332	67,080
	gC <sup>Fst</sup>	66	19,717	2	0
	gC1 <sup>ex</sup>	161	0	94	38,479
	gC2 <sup>ex</sup>	538,052	0	434,196	21,471
D	gD	49,883	133,558	7,781	3,552
E	gE1+/gE2+	14,804	106,379	163,481	942
	gE4 <sup>t</sup>	14,772	5,651	98	0
	gE4 <sup>e</sup>	53	0	246,486	1,301

### CR3

The editing pathways of CR3 are significantly different from those of RPS12 and ND7 5' (Figure 32, Chapter 4). The most obvious difference is the fact that the pathways are highly branched, generating multiple distinctly different mRNA products. The other key difference is that most of the editing pathways defined in the TREU and EATRO cells are not shared. The two cell lines appear to use almost completely different sets of gRNAs to edit CR3, with the exception of some edits at the very 3' end of the transcript. What was most surprising about these data was the fact that most of the gRNAs necessary to generate both the TREU and EATRO CR3 editing pathways were found in both gRNA transcriptomes in similar relative abundances (Table 16, Chapter 4). In addition, similar to the RPS12 and ND7 5' data, we saw very little correlation between the abundance of the gRNA population and the corresponding abundance of the mRNA transcript generated. For example, in the EATRO cell line, the gB1B2 population (38,663 reads) was significantly more abundant than the gB4 guide RNA population (100 reads). Nevertheless, the number of mRNAs with the B4 editing pattern was 5-fold higher

than the B1B2 mRNAs. The EATRO B4 mRNAs can be further edited by three different gRNAs, gB5<sup>e</sup>, gB6<sup>e</sup> and gB7<sup>e</sup>. Again, while gB6<sup>e</sup> is the most abundant of the three, the predominant mRNA found is the B5<sup>e</sup> transcript (Table 16). In addition, while the gB5<sup>e</sup> guide RNA is also abundant in TREU cells, we found no corresponding B5<sup>e</sup> mRNAs in this cell line. The most significant divergence of the EATRO CR3 editing pattern occurs at the B to C editing block transition. In TREU, all of the different 3' editing patterns converge and are further edited by the gC guide RNA population. In EATRO cells, the gC population is much less abundant (1590 reads instead of >17,000). Correspondingly, very few mRNA transcripts were observed that used a gC guide. Instead, the bulk of the B-level transcripts were further edited by gC<sup>e</sup> (a rare gRNA with only 43 reads detected). The C-block gRNAs also contained the most abundant gRNA detected for CR3 editing, gC<sup>e80</sup>. This gRNA has close to 500,000 reads in both cell line transcriptomes. It can anchor to all of the B block mRNAs in both the TREU and EATRO pathways but was only observed editing the B5e template in the EATRO cells. Despite the overwhelming abundance of the gC<sup>e80</sup> guide RNA population, the corresponding mRNA had the smallest number of detected reads.

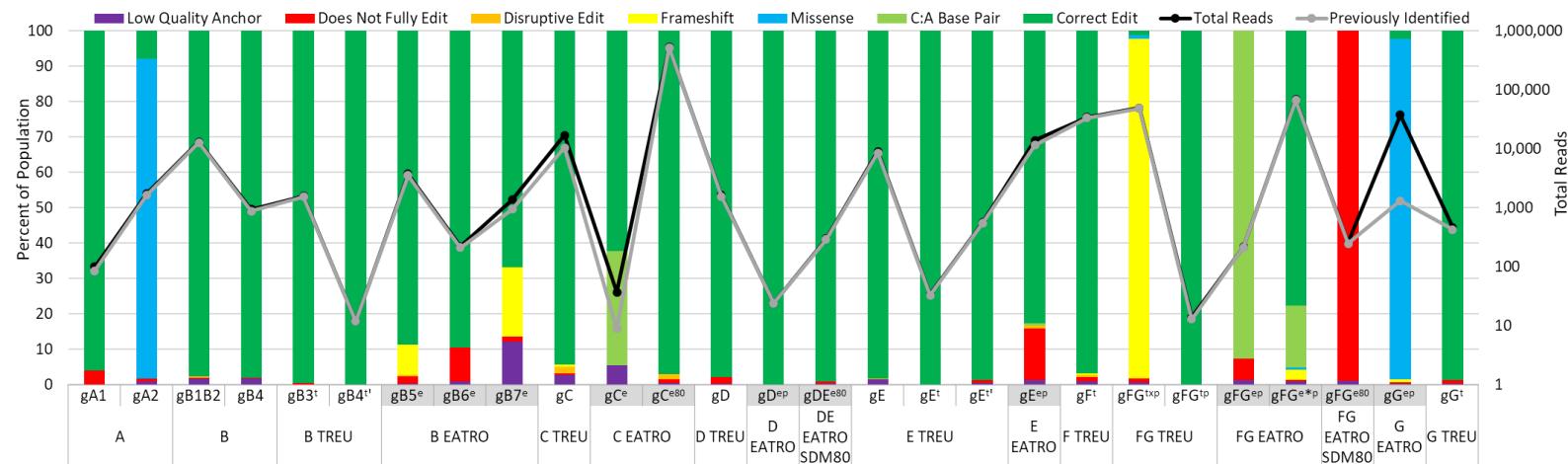
Surprisingly, there were very few CR3 gRNA populations in either cell line that had sequence members that could disrupt editing (Figures 33 and 34). In the TREU cell line, the FG<sup>tx</sup> 5' end pattern of mRNA editing is generated by a promiscuous gRNA, gFG<sup>txp</sup> (originally identified as a CR4 gRNA). This gRNA population is predicted to cause a frameshift just upstream of the anchor region of gFG<sup>txp</sup>. However, this frameshift is not observed in the mRNA population, and it appears that this gap is tolerated like those observed in RPS12.



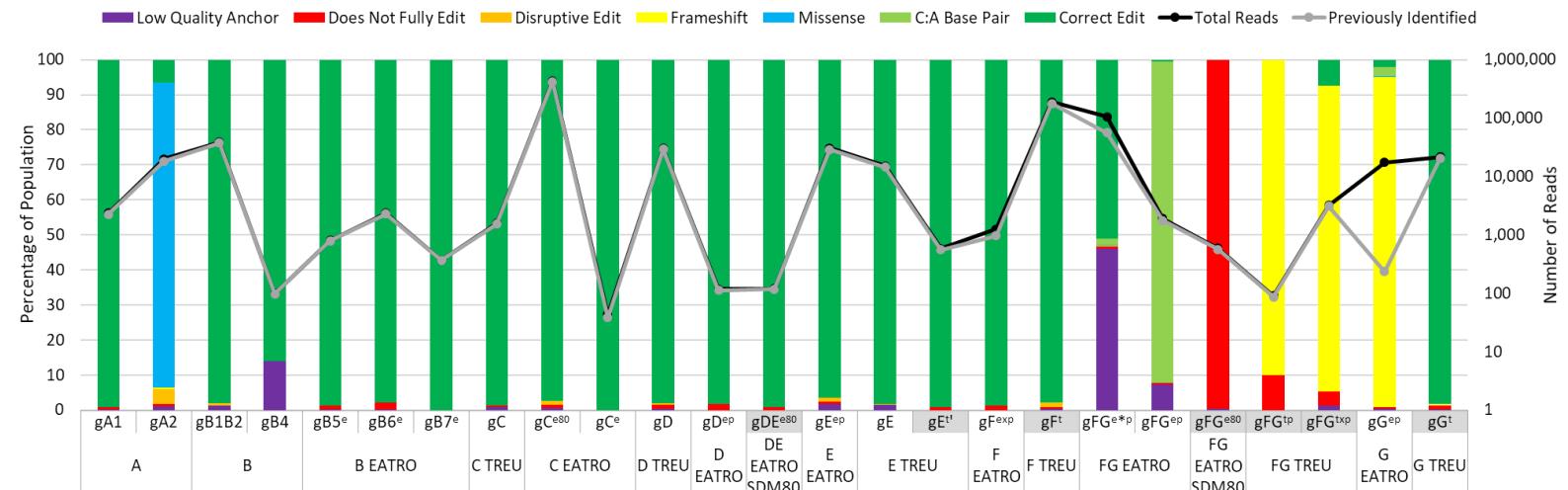
**Figure 32. Observed CR3 editing pathways in TREU 667 (A and B) and EATRO 164 (C and D) cell lines.** For descriptions of dots and arrows, please see Figure 26. Dashed arrows with open heads represent a hypothetical rewrite. + indicates that more than one mRNA form was condensed into this circle to simplify the figure. Condensed forms encode largely the same amino acid sequence with only small variants. Terminal dots are colored blue for reading frame 1, or magenta for reading frame 2. Boxed magenta dots have no functional start codon but are translatable into reading frame 2 with the use of an alternative start codon (UUG). Lines connecting dots (B and D) indicate gRNA population size and functionality. Usable gRNAs were considered to be those that correctly edit, utilize C:A base pairs to edit, or generate only small missense mutations.

**Table 16. CR3 gRNA population analysis**

Editing Block	Population	TREU 667 gRNA Reads	TREU 667 mRNA Reads	EATRO 164 gRNA Reads	EATRO 164 mRNA Reads
A	gA1	100	5,907	2,368	534
	gA2	1,729	28,654	19,947	5,934
B	gB1B2	13,278	10,776	38,663	757
	gB3 <sup>t</sup>	1,598	9,158	0	0
	gB4	945	9,536	100	3,720
	gB4 <sup>t'</sup>	12	8,692	0	0
	gB5 <sup>e</sup>	3,800	0	803	2,025
	gB6 <sup>e</sup>	221	0	2,359	689
	gB7 <sup>e</sup>	1,376	0	369	744
C	gC	17,045	27,646	1,590	616
	gC <sup>e</sup>	37	0	43	2,079
	gC <sup>e80</sup>	538,230	0	434,196	270
D	gD	1,649	26,287	30,125	280
	gD <sup>ep</sup>	24	0	118	1,981
DE	gDE <sup>e80</sup>	301	0	120	0
E	gE	9,135	15,016	15,242	204
	gEt <sup>t</sup>	567	14,687	582	0
	gEt <sup>t</sup>	34	5,099	0	0
	gE <sup>ep</sup>	14,062	0	30,982	1,864
F	gF <sup>t</sup>	34,666	8,073	189,828	0
	gF <sup>exp</sup>	0	0	1,265	112
FG	gFG <sup>tp</sup>	14	1,035	90	0
	gFG <sup>txp</sup>	49,892	2,172	3,220	0
	gFG <sup>ep</sup>	219	0	1,912	79
	gFG <sup>e*p</sup>	69,138	0	104,747	691
	gFG <sup>e80</sup>	255	0	589	0
G	gG <sup>t</sup>	466	1,231	21,788	0
	gG <sup>ep</sup>	38,001	0	17,318	68



**Figure 33. Analysis of functionality and abundance of gRNA subpopulations that edit CR3 in TREU 667 cells.** For description of axes and gRNA functionality labels, please see Figure 27. gRNAs with shaded names were found in the TREU 667 gRNA transcriptome but were not found to be utilized in editing the TREU 667 CR3 mRNAs, despite being utilized in EATRO 164 cells.



**Figure 34. Analysis of functionality and abundance of gRNA subpopulations that edit CR3 in EATRO 164 cells.** For description of axes and gRNA functionality labels, please see Figure 27. gRNAs with shaded names were found in the EATRO 164 gRNA transcriptome but were not found to be utilized in editing the EATRO 164 CR3 mRNAs, despite being utilized in TREU 667 cells.

In the EATRO cell line, only two terminal gRNA populations have predicted editing problems. About half of the gFG<sup>e\*p</sup> subpopulation has a low-quality anchor (<5 nt), however the subpopulation does have >53,000 reads that have higher quality anchors, so this mutation may be tolerated. The gG<sup>e\*p</sup> population is also interesting, in that it is predicted to generate a frameshift when editing was simulated by our GUIDE program (Figure 34). This frameshift is not observed in the mRNA transcriptome data however because editing appears to stop before reaching the last two sites (Chapter 4).

## **Discussion**

ACORNS and GUIDE are two new tools that can aide in the understanding of the complex dynamics of the kinetoplastid RNA editing system. The ability of ACORNS to cluster related gRNAs proved to be a powerful mechanism for the identification of gRNAs with mutations that disrupt their alignment to fully edited sequence. These analyses allowed us to identify and characterize nearly 2 million additional gRNA transcripts from our transcriptome libraries. In addition, these analyses identified a large cohort of gRNA clusters that are not involved in directing the sequence changes associated with the known canonical transcripts. More than 25% of the gRNA transcripts found in the EATRO and TREU gRNA transcriptomes have no known functional relatives. This strongly suggests that the coding capacity of the mitochondrial genome is much larger than previously thought.

Full characterization of the known gRNA population was also informative. In our initial gRNA characterization study, we had noted the extreme population differences found between the different identified gRNAs [9]. We had hypothesized that the low copy number gRNAs were an artifact of the high stringency of our initial screen. The cluster analyses suggest that extreme

population size differences do exist between the different gRNAs. In both the TREU and EATRO transcriptomes, ~30 clusters were identified with over 500 sequence members each. These clusters accounted for the bulk of the gRNA transcript reads (6,461,870 reads in TREU and 4,539,263 reads in EATRO). In contrast, over 500 clusters were found that contained fewer than 25 sequence members. These tiny clusters account for only ~135,000 transcript reads. The very large numbers of low copy number gRNAs are suggestive of high plasticity in the gRNA encoding minicircles. Studies in *Leishmania* have suggested that minicircle sequence class frequencies are extremely variable [63,78]. How these huge differences in gRNA abundance influences gRNA selection and use is unclear. In only ~50% of the editing branch points characterized for RPS12, ND7 5' and CR3, did the abundance of the gRNAs involved somewhat correlate with the preferred editing path. However, even when gRNA abundance did align with mRNA abundance, they were often not proportional. In addition, block editing by an abundant gRNA is often followed by editing using a rare gRNA, with no equivalent drop in editing efficiency. In addition, multiple instances were observed, where highly abundant gRNAs did not appear to be used in one cell line. In one example, the gRNAs responsible for generating the B5<sup>e</sup>, B6<sup>e</sup> and B7<sup>e</sup> mRNA forms of RPS12 are present in both cell lines, but only apparently act in the EATRO cell line, despite being all more than ten times more abundant than their only competitor. It may be that protein factors play a predominant role in gRNA selection. One study examined the role of the RNA editing mediator complex (REMC), which is heterogenous and consists of one primary subunit TbRGG2 that formed associations with either MRB8170 or MRB8180 [161]. They showed that depletion of MRB8180 caused global effects on RNA editing, but depletion of MRB8170 had transcript specific effects, substantially increasing the amount of

pre-edited RPS12 transcripts, but not significantly affecting the amount of pre-edited ND7 5' transcripts. This specificity is intriguing, and it is possible that the REMC or other protein factors are involved in gRNA selection.

The identification of all relatives of conventional gRNAs, including those that had undergone a mutational event, also allowed us to more carefully characterize the effect of mutational noise on the RNA editing system. These analyses indicate that the RNA editing system can tolerate more mutational noise than we had originally hypothesized. Surprisingly, a large number of the gRNA populations we characterized had base pair mismatches with their best aligned mRNA transcripts. In this study, we saw that while most of the mismatches were specifically C:A base pairs, almost every other mismatch was also observed. Even gaps in the alignment of the gRNA to the mRNA appeared to be tolerated. In almost all cases however, the mismatch base pair and/or gap appeared to be stabilized by multiple flanking G:C base pairs. While we cannot rule out the possibility that rare, perfect match gRNAs do exist for these regions, it may be that the incompletely base paired interaction is the most stable structure possible, hence is generated by that gRNA [162]. While a minimum number of non-paired nucleotides does appear to be tolerated within the guiding region, mismatches or non-Watson-Crick base pairs within the anchor region do not appear to be tolerated, greatly affecting the efficiency of the RNA editing process [119].

In a large number of the gRNA populations containing alignment mismatches, our GUIDE program predicted that the mismatch would drive the generation of an alternative edit. However, these alternative edits were not observed in the mRNA transcriptome data. These data contradict one of the most prominent models of RNA editing progression, known as the

“mismatch recognition” model [33]. In this model, when the gRNA/mRNA duplex initially forms, the editosome proceeds to edit beginning at the first mismatch site closest to the anchor binding region. Once this mismatch is resolved either by the insertion or deletion of a uridine, the next mismatch site is edited, and no sites further will be edited until the sites nearest the anchor binding region are resolved. When GUIDE predicts editing patterns, it follows this model, predicting each editing site, moving from the anchor binding region towards the poly-U tail. If editing did follow this strict mismatch recognition model, the alternative sequences predicted by GUIDE should have been observed in the mRNA transcriptome data. An alternative model suggests that RNA editing occurs via a more “dynamic interaction” [162]. This model proposes that when a gRNA/mRNA duplex forms, the editosome targets regions of the duplex with low thermodynamic stability and edits those regions. As editing progresses, the duplex realigns, changing the targets of the editing system. These cycles of progressive realignment proceed until the gRNA/mRNA duplex reaches maximum stability. In this way, RNA editing does not necessarily proceed in a strict 3’ to 5’ directional manner. This model suggests that mismatches and gaps in alignments can be tolerated because they do not significantly impact the stability of the final gRNA/mRNA duplex. Supporting this, are the frequent observations of neighboring G:C base pairs, which would substantially enhance the stability of these mismatches. The “dynamic interaction” model is further supported based on the existence of “junction regions” in partially edited mRNAs [161–165,156,166,148,160]. These regions adjoin the unedited and fully edited regions of a partially edited mRNA, but do not match either the unedited or fully edited sequence. Junction regions vary significantly across partially edited transcripts, possessing no consensus sequence. These regions can vary in size

and depletion of different protein factors affects their occurrence and length, but the presence of junction regions remains ubiquitous across all partially edited mRNAs [161–165,156,166,148,160]. The mismatch recognition model reconciles the presence of junction regions as areas of mis-editing (utilization of the wrong gRNA or a misaligned gRNA), hence all junction regions would have a gRNA capable of generating the sequence. In our characterization of the editing pathways of RPS12, ND7 5' and CR3, highly abundant mRNA sequences were screened against the gRNA transcriptomes at low stringency levels in order to identify true alternative edits. While we were able to identify and number of gRNAs that could direct alternative edits, a large number of “junction sequences” were identified that do not match any gRNA in our databases. If the “dynamic interaction” model is correct, these multiple variable junction sequences could be generated by the same gRNA during the editing process. Alternative base pairs have also been shown to be tolerated to different extents in an *in vitro* gRNA directed deletion assay [151]. In this study, substitutions were made of the nucleotides immediately upstream of a deletion site. This study found that when the base pair upstream of the deletion site was C:A, C:U or a C:C, the site was still found deleted in the mRNAs. Deletions were not observed when the base pair was a G:A or G:G.

Another facet of gRNA utilization is the existence of promiscuous gRNAs, gRNAs editing more than one target. In this analysis we showed many populations that edit more than one gene, and one population editing the same gene in two different locations (gF<sup>P</sup> of RPS12). Interestingly, most of the promiscuous populations were found in the CR3 data set, and these promiscuous gRNAs were the only productive promiscuous gRNAs identified. No promiscuous populations were found editing ND7 5', and the only promiscuous gRNAs found to edit RPS12

lead to editing pathway dead ends. Predictions made by GUIDE indicate that many of the promiscuous gRNAs generating the CR3 EATRO specific pathways should also be able to productively edit in TREU. Because these are promiscuous gRNA, it is possible that their availability is impacted by their use in editing other transcripts. A full mRNA transcriptome would allow a full analysis of the global impacts of editing on any particular gRNA cluster.

This study revealed the surprising amount of noise and errors that are tolerated in the RNA editing system of *Trypanosoma brucei*. Some questions still remain, such as the functions of the large proportion of unknown gRNAs, how gRNAs are selected, and what is the true extent of gRNA promiscuity. In order to answer these questions, we believe that a full deep sequence of all edited mRNAs paired with gRNA transcriptomes of multiple cell lines would shed more light on this complicated situation.

### **Acknowledgements**

We would like to thank the Ken Stuart Lab for the trypanosome cell lines and Chris Adami for his assistance in the conceptualization of this work.

# CHAPTER 6: SUMMARY AND DISCUSSION

## Introduction

*Trypanosoma brucei* is one of the few organisms that utilizes the kinetoplastid RNA editing system. This system seems unnecessarily complex, using two genetic components to generate one fully functional product. Moreover, this system is prone to malfunction; each gRNA generates the anchoring region for the next gRNA, making the system sequentially dependent. This makes the mutation or loss of any gRNA along the editing pathway extremely detrimental, especially considering that two of the twelve edited genes are essential [17,100]. This problem is made worse by the fact that some gRNAs are incredibly rare, and during replication, the 5,000-10,000 minicircles encoding the gRNAs are divided asymmetrically, making minicircle loss not only possible, but routine [8,167]. This system should not work, but it does. Kinetoplastids are some of the most successful parasites on Earth, infecting insects, plants, mammals, fish, birds, and reptiles [168].

The study of this editing system inevitably brings up questions of how this system evolved and how it continues to be maintained. The concept of drift robustness begins to explain this surprising amount of fragile complexity [80]. Drift robustness is a form of genetic robustness that allows an organism to be protected from extreme events of genetic drift by making mutations either neutral or lethal. *T. brucei* frequently undergoes population bottlenecks throughout its life cycle and as its mitochondria is completely asexual, it seems particularly prone to genetic drift [22]. The fragility of the RNA editing system seems to coincide remarkably well with the idea of drift robustness. By rendering mutations or loss of

minicircles lethal, this would prevent accumulation of slightly deleterious mutants in the population. But such a system should suffer from another disadvantage. In a system where mutations are lethal or neutral, how can the organisms continue to evolve?

In the first examination of the gRNA transcriptome, many gRNAs were identified that were capable of generating alternative edits [9]. These edits ranged from having no effect on the protein sequence to causing a frameshift and altering a large portion of the protein. These findings sparked this project, which sought to understand the impact of alternative editing on genetic integrity, developmental regulation, protein diversity and editing efficiency. These goals were accomplished through the generation of the bloodstream gRNA transcriptome and comparative analysis of it with the insect stage gRNA transcriptome, analysis of dual-coding genes that utilize alternative edits to access multiple reading frames, the generation of libraries of putative dual-coding mRNAs at different states of editing, and analysis of gRNA population diversity and the impact of that diversity on the editing system as a whole.

## **Summary of Chapter 2**

This chapter characterized the gRNA transcriptome of EATRO 164 bloodstream stage *Trypanosoma brucei*, and compared it to the gRNA transcriptome of procyclic stage trypanosomes. As with the procyclic gRNA transcriptome, conventionally accepted fully edited mRNA sequences were used to identify the gRNAs, and a comparison of the two life cycle transcriptomes show a 3.5:1 ratio of procyclic to bloodstream gRNA reads. This ratio varies significantly by gene and by gRNA populations within genes. The variation in the abundance of the initiating gRNAs for each gene, however, displays a trend that correlates with the developmental pattern of edited gene expression. Surprisingly, there were very few gRNAs

found in both transcriptomes, but there were many gRNAs that appeared to be related between transcriptomes. Comparing these related major classes from each transcriptome revealed a median value of ten single nucleotide variations per major class. Nucleotide variations were much less likely to occur in the consecutive Watson-Crick anchor region, indicating a very strong bias against G:U base pairs in this region. In spite of the variation we saw between related gRNAs, we did find several conserved gRNA characteristics, such as transcription start site sequence, length of complementarity, and non-base pairing nucleotides at the 5' and 3' ends of gRNAs. Overall, gRNA coverage of edited mRNAs as well as overlap between adjacent gRNAs was lower in the bloodstream gRNA transcriptome than in the procyclic gRNA transcriptome.

This work indicates that gRNAs are expressed during both life cycle stages, and that the differences in the extent of editing previously reported for different mRNA transcripts are not due to the presence or absence of gRNAs. However, the abundance of the initiating gRNAs may be important in the developmental regulation of RNA editing.

### **Summary of Chapter 3**

In this work, we show that many of the mitochondrially edited mRNAs in *T. brucei* can alter the choice of open reading frame by alternative editing of the 5' end. Dual-coding genes have specific mutational biases, such as an increase of the ratio of nonsynonymous to synonymous mutations, and an increase in mutational frequency overall. Analyses of mutational bias of all mitochondrial genes indicate that six of the pan-edited genes may be dual-coding. These analyses include measuring the conservation of editing patterns between *T. brucei* and *T. vivax* transcripts and frequencies of different types of mutations. These data were

used in a principal component analysis, which showed a distinct difference between alternative reading frames of dual-coding genes and single-coding genes. Discovery of alternative gRNAs reveal that RNA editing can allow access to both reading frames. We predicted the functions of two of these alternative reading frames as small metabolite transmembrane transporters. We hypothesize that dual-coding genes can protect genetic information by overlapping genes that are under selection at different portions of the life cycle.

## **Summary of Chapter 4**

In this study, we analyzed the editing patterns of three putative dual-coding genes, ribosomal protein S12, the 5' editing domain of NADH dehydrogenase subunit 7, and C-rich region 3, and constructed detailed editing pathway maps using mRNA and gRNA transcriptome data. While editing of RPS12 showed only transcripts that produce the canonical RPS12 protein, we did observe a second downstream start codon capable of producing the alternative protein, if selected by the ribosome. In ND7 5' and CR3, we found evidence that both of these transcripts are edited to express protein products in more than one reading frame. Moreover, we found that CR3 has a very complex set of highly branched editing pathways that vary significantly between cell lines, with a different set of gRNAs being used in each cell line, despite both sets of gRNAs being present in both cell lines. We also found that changing the energy source available to cells also alters the editing preferences of both CR3 and ND7 5'. In addition to this, we found evidence that the poly-U tail that is added post transcriptionally to gRNAs may also be used in editing. These findings suggest that these reading frames can be alternatively selected based on the current environment, and that alternative editing may be a way for the trypanosomes to continue evolving this rigid editing system.

## **Summary of Chapter 5**

In the analysis of the gRNA transcriptomes, gRNAs were identified based on complementarity to edited mRNAs. While millions of gRNAs were found using this method, we discovered that many gRNAs were still left unidentified. This high proportion of unidentified gRNAs was shocking, as we had predicted that the RNA editing system should be intolerant of mutations, and the presence of so many unidentified gRNAs meant that many editing pathways had not been characterized, or many mutated gRNAs were also present. To determine the identity and function of these gRNAs, two new programs were created, ACORNS and GUIDE. The first program functions to group related gRNAs into clusters, where each cluster had significant sequence conservation. This analysis showed that more than half of all unidentified gRNAs were not related to any functionally known gRNAs and could be generating uncharacterized alternative editing patterns.

However, there were still many gRNAs that were related to previously identified gRNAs. These gRNAs could be capable of disrupting the editing process or generating small alternative edits that are tolerated. In order to investigate this, our second program, GUIDE examined the gRNAs responsible for generating the editing pathways of RPS12, ND7 5' and CR3, as well as their previously unidentified relatives. Using the defined gRNA clusters, GUIDE screened all members editing these three genes against their targets to determine what proportion of each family was able to productively edit. The initial analyses of these previously unidentified gRNAs revealed that nearly all were predicted to disrupt the editing system, but in our examination of the mRNA data, we found this not to be the case.

By combining the analyses of the GUIDE program and the mRNA transcriptome data, we learned more about the robustness of the RNA editing system. In examining the RPS12 gRNAs we found that single mismatches or gaps were highly tolerated in gRNA alignments, with only very small drops in editing efficiency. In the examination of the ND7 5' gRNAs, however, we identified the limit of this tolerance. Significant drops in editing efficiency were observed when gRNAs either possessed multiple mismatches or gaps, or possessed mismatches that disrupted the anchor binding region. Finally, in our analysis of the CR3 gRNAs, we observed many gRNA populations of high abundance that were apparently not used to edit. We found that RNA editing preference does not correlate positively or negatively with gRNA abundance, and in most cases, the apparently unused gRNA populations had no issues in their mRNA alignments. We predict that these gRNAs may be in use elsewhere in the editing system, and to fully understand this system, we recommend that a full mRNA transcriptome be analyzed.

## **Genetic Integrity**

### **The introduction and maintenance of kinetoplastid RNA editing**

As previously mentioned, the sheer size and complexity of the RNA editing system has left many in search of the answers to how it evolved and how it has been maintained. Many hypotheses have been proposed, such as it being a relic of the old RNA world, being a product of constructive neutral evolution, or that the system co-evolved with G-quadruplex structures that served to protect the genetic information [69–74]. As for how it has been maintained, one prominent theory is that RNA editing is advantageous because it is a mechanism by which an organism can fragment and scatter essential genetic information throughout a genome [75,76]. Because kinetoplast DNA is less stable than chromosomal DNA, and minicircles are frequently

lost due to asymmetric division, this hypothesis suggests that scattering essential guide RNA genes throughout the DNA network would prevent fast growing deletion mutants from outcompeting more metabolically versatile parasites during growth in the mammalian host [76,77].

We propose that the RNA editing system, and its inherent fragility operate as a system to weed out deleterious mutations by making them lethal, as a form of drift robustness. Drift robustness is not adaptive, however, and prevents the population from generating beneficial mutations as well. This strategy leaves organisms no options for evolving. We suggest that alternative edits, such as those seen in CR3 and others previously observed, editing pathways generate this evolution without compromising the rigid conservation of other genes such as the essential RPS12 [68].

Supporting this is the fact that, of the kinetoplastids, *Trypanosoma* have some of the harshest life cycles in terms of maintaining genetic integrity, with the electron transport genes under very relaxed selection in the glucose rich bloodstream stage, and very strict selection imposed in the insect stage. In conjunction with this, we observe that *Trypanosoma brucei*, *Trypanosoma cruzi*, and *Trypanosoma vivax* all maintain more genes that are edited and more extensive editing than their other kinetoplastid counterparts with milder life cycles [4,16,19–21,52–56,58,59]. For example, *Phytomonas serpens* infects important crops and is transmitted by sap feeding bugs. These parasites have glucose readily available in both life cycle stages, and are unique in that they lack a fully functional respiratory electron transport chain, and this species is also missing two edited genes entirely, and pan-edits six genes and partially edits three genes, compared to the nine pan-edited and three partially edited genes in the

*Trypanosoma spp.* [64,65]. For *Leishmania spp.*, all life cycle stages possess an active Krebs cycle and ETC linked to the generation of ATP, but these cells are never restricted from access to glucose [61,62,139]. *Leishmania* do not pan-edit ATPase 6 or ND7, but only partially edit them [63]. These observations suggest that RNA editing provides a larger advantage to organisms with a more complex life cycle.

## Dual-coding and dual-function genes

One oversight that has not been accounted for in the hypotheses that attempt to explain the advantage of kinetoplastid RNA editing is how genetic material that is not under selection is maintained. There has been considerable debate on the necessity of Complex I subunits for either stage of the trypanosome life cycle. Studies using RNAi and knockout cell lines of nuclear-encoded members of Complex I have shown that the complex is unnecessary for survival in either life cycle stage, and in this work we were unable to find complete gRNA coverage of the edited ND subunits in the bloodstream stage, despite the ND subunits generally being more fully edited or only fully edited in the bloodstream stage [5,26–29,111,112]. However, the nuclear encoded Complex I member genes are maintained [42], and the vast majority of the gRNAs required to edit mitochondrially encoded ND subunits were found in both life cycle stages. This evidence together suggests that the Complex I proteins are vulnerable to genetic drift but have somehow been maintained.

In this work, we propose that by overlapping Complex I genes not under strict selection with genes that are under selection, the accumulation of mutations can be prevented. Because these overlapped genes share most gRNAs, and alternative edits only occur in the terminal gRNAs, this strategy ensures that almost all of the genetic material is protected. The genes

predicted to be dual-coding based on our mutational bias analysis include almost all of the edited ND transcripts, excluding only ND8. In the examination of procyclic ND7 5' and CR3 (putative ND4L) we found that these two genes are alternatively edited to produce more than one protein product [120]. This evidence shows that not only are dual-coding genes being utilized in *T. brucei*, but also that RNA editing is facilitating the use of these dual-coding genes, providing the world with a concrete advantage of utilizing this type of RNA editing.

Based on limited sequence homology, we hypothesize that the trypanosome mitochondrial alternative reading frames (ARFs) encode small metabolite transporters that provide a distinct growth advantage to bloodstream form parasites. These proteins would function differently from all other edited proteins, which do not function as transporters. While it has been previously suggested that both alternative editing and dual-function proteins are important mechanisms for expanding the functional diversity of proteins found in trypanosomes, a duplication event could easily alleviate the evolutionary constraints imposed by dual-coding genes [67,97–99]. We maintain that in salivarian trypanosomes, these genes must provide the additional benefit of protecting genetic information in order to continue to be overlapped. Protection of the mitochondrial genome during growth in the mammal would increase the capacity for successful transfer to an insect vector and maximize the parasites long-term survival and spread.

Analyses of other trypanosomes do show that some of the ARFs have intriguing homology to the ARFs identified in *T. brucei* and *T. vivax*. However, most of the ARFs are punctuated with stop codons. It is possible that these genes have since lost their function and are no longer required to be dual-coding due to the reduced selective pressures endured during

their life cycles, or it is possible that these stop codons are removed by alternative editing events.

In addition to the dual-coding genes we have described, another dual-coding gene, COIII, is accessed through alternative editing, by connecting an unedited 5' reading frame with an edited 3' reading frame, and the alternative protein produced, AEP-1, has been shown to be essential [67,68]. In addition, there are known Krebs cycle proteins ( $\alpha$ -ketoglutarate dehydrogenase E2 and  $\alpha$ -ketoglutarate decarboxylase) that have two functions, rendering them protected while they are not under selection as well [97,98]. These data show that trypanosomes are utilizing dual-coding and dual-function genes to protect genetic integrity.

Another set of dual-function genes appear in the gRNAs, as promiscuous gRNAs. The editing pathways of CR3 are littered with gRNAs identified to edit other genes, such as ND8 and the 3' editing domain of ND7, both of which were not predicted to be dual-coding, and are under less protection. We believe that this is yet another mechanism of increasing the drift robustness of *T. brucei* through the use of alternative RNA editing.

## **Developmental Regulation**

In the examination of the bloodstream gRNA transcriptome, we found the abundance of the initiating gRNAs in the procyclic and bloodstream gRNA transcriptome is correlated with the developmental editing patterns of the genes they edit. However, we cannot rule out the possibility that not all of the populations of initiating gRNAs were identified. We identified alternative initiating gRNAs for CR3, and without deep sequencing all pan-edited mRNAs, we can't know that others don't exist.

It has been previously reported that gRNA presence did not correlate with developmental RNA editing patterns in *T. brucei* [50,51]. This, however, was reported on the observation of a very limited number of gRNAs, but we found that for the most part, this held true, with gRNA populations having similar relative abundances across both transcriptomes, with, of course, the exception of the initiating gRNAs.

In our examination of the editing pathways of RPS12, ND7 5', and CR3, we found many editing patterns that were only observed or were more prominent in only one cell line. This was most prevalent in CR3, where the editing pathway of the TREU cells is almost completely different from that of the EATRO cells. Curiously, the gRNAs required for both pathways are present in the transcriptomes of both cell lines in relatively equal abundances. Additionally, another exclusive editing pathway was discovered when the EATRO cells were moved into glucose depleted medium. This pathway produced a transcript that would make a unique protein, totally different from all other CR3 protein products. These gRNAs appear functional in both cell lines but are perhaps being used to edit a different gene when they are not observed to be editing CR3. Indeed, many of the EATRO specific CR3 gRNAs are promiscuous gRNAs known to edit other genes.

To better understand the complexities of gRNA selection, we examined the editing branch points of the RPS12, ND7 5' and CR3 pathways. Examination of editing branch points revealed that the abundance of the gRNAs involved did not correlate with the observed editing preference. This work has raised many questions about how gRNAs are selected and editing pattern preferences are exerted, and more study is needed to understand this system.

## **Protein Diversity**

In the examination of the RPS12 mRNAs, we found one major alternative editing event. This event was observed in the EATRO 164 cell line only, and causes a frame shift that extends the reading frame at the 3' end by nine amino acids. This same alternative edit was previously described in the 29-13 strain as well, which shows that this edit is not an isolated occurrence in the EATRO 164 strain [148]. Frameshifting gRNAs were also identified to edit ATPase 6 in both bloodstream and procyclic EATRO 164 cells. The predicted frameshifts also occur close to the 3' end of the transcript and alters the C terminus of the protein. As the frameshifts occur downstream of the highly conserved amino acid region involved in proton translocation, it may be that this is also tolerated [31].

Edits of the ND7 5' mRNAs generate transcripts that can translate into two reading frames in TREU 667 and EATRO 164 cells. Interestingly, this gene was also sequenced in 29-13 cells [148]. That study indicated that a large proportion of the fully edited ND7 5' transcripts had a single nucleotide difference in the 5' UTR. As the upstream start codon that allows the ARF to be translated is in what is known as the 5' UTR, this difference could be the same alternative edit that generates the ARF transcripts. This suggests that this alternative editing may be widespread.

Like ND7 5', we also detected evidence of dual-coding in CR3. The alternative edits in both cell lines produce multiple variations of the CR3 proteins. While the editing efficiencies of CR3 are very low, possibly preventing the generation of these CR3 proteins, we propose that this mechanism of branched editing pathways is a safe way for the trypanosomes to introduce variation into the rigid system and continue to evolve.

In addition to the high level of alternative editing observed in CR3, in our re-examination of the procyclic EATRO 164 and TREU 667 gRNA transcriptomes, we identified many gRNAs in the existing transcriptomes still have no known function. These gRNAs may be generating alternative editing events, thus further increasing the protein diversity of *T. brucei*.

## **Editing Efficiency**

The editing efficiencies of all three genes we examined were surprisingly low. The efficiencies of TREU cell line mRNAs were all less than ten percent, while the EATRO 164 cells grown in SDM79 were all less than one percent. With CR3 and RPS12, more than 80% of mRNAs were completely unedited, while those numbers were much lower in ND7 5'. There are many factors we observed that had the potential to affect the overall editing efficiency.

## **Mutations and non-canonical base pairs**

Because gRNAs utilize both canonical (Watson-Crick) as well as G:U base-pairing to direct the change in sequence, most transition mutations in the gRNA, would not lead to changes in the mRNA sequence and would not be selected against [33]. In our observations of the gRNA transcriptome, we found a very strong bias against A to G transitions in the anchor regions of the gRNAs, suggesting that G:U base-pairing is not well tolerated in this region. However, we also found many populations of gRNAs where non-canonical base pairs appeared to be tolerated, even in essential genes, ATPase 6 and RPS12 [17,100,107,110]. In the bloodstream database, a gRNA population with a C:U base pair must be tolerated to be able to complete the editing pathway of ATPase 6, and in RPS12, we observed a gRNA population that requires toleration of a gap in the alignment to be capable of fully editing.

In addition to this, RPS12 and ND7 5' possess two highly cytosine rich regions that typically have poor gRNA coverage. These regions both encode conserved amino acids vital to the functions of the proteins. When we deep sequenced these regions of the mRNAs, we had hoped to find alternative sequences that allowed us to identify a more abundant population to carry out these edits, but we only found the previously described editing patterns. Using these sequences and performing a search of the gRNA databases at a low stringency yielded populations with imperfect anchors, requiring noncanonical base pairs to be tolerated. In RPS12, this did not appear to affect editing efficiency, but in ND7 5', the effects were severe. However, the gRNA population identified to edit this region of ND7 5' had multiple mismatches and gaps in both cell lines, whereas the RPS12 populations did not.

These examples are far from isolated incidents. Many other populations were identified where noncanonical base pairs were required to generate fully edited mRNA sequences. In most cases, these alternative base pairs do not seem to affect the editing efficiency. The most prominent mismatch was the C:A base pair, but almost every other mismatch was observed. Even gaps in the alignment of the gRNA to the mRNA seem to be tolerated. Of note, we did observe that most mismatches were flanked by nearby G:C base pairs that may have assisted in stabilizing the alignments. The use of alternative base pairs has been previously shown to be tolerated at different extents, with the use of C:A, C:C and C:U base pairs not completely disrupting editing [151]. This evidence suggests that the RNA editing system will tolerate some amount of illegal base pairs, there is a limit to what will be tolerated.

The use of non-canonical base pairs in RNA editing support the model of editing known as the “dynamic interaction” model [162]. In this model, editing of an mRNA in a gRNA duplex

does not proceed in a site by site fashion, but instead, the editosome targets regions of low stability first, and continues to edit and re-edit the mRNA until a thermodynamically stable gRNA/mRNA duplex is achieved. In this model, non-canonical base pairs may be tolerated if they do not significantly disrupt the stability of the duplex.

## **Overwriting**

In our exploration of partially edited mRNAs, we found that RNA editing is not always strictly sequential. While overall, editing proceeds from the 3' to 5' across of the editing domain, we have found evidence that shows that gRNAs can overwrite editing that has been previously generated. Another study that examined the use of alternative RPS12 gRNAs previously identified found that the three gRNAs in question were being utilized, and that despite not generating the conventional editing patterns, a very small number of transcripts returned to the canonical editing pattern after the alternatives [148]. This last observation may be another instance of gRNA overwriting. This is another source of noise, lowering overall editing efficiency.

## **Future Work**

In this work, we proposed the hypothesis that *T. brucei* utilizes dual-coding genes to protect vulnerable genetic material from genetic drift and accesses the overlapped reading frames through alternative RNA editing. In order to test this hypothesis, we should confirm the presence of the alternative protein products *in vivo* and knock down these products to determine if they provide a benefit to the cells. Knocking down edited transcripts has proved difficult, but not impossible. One study was able to genetically engineer an artificial site-specific RNA endonuclease to target the ATPase 6 edited mRNA [110]. This engineering

requires an eight-nucleotide specific target sequence, and in order to use this technique, this would require finding a target sequence specific to the alternatively edited mRNAs only, and leave intact the mRNAs expressing the canonical ETC proteins, which could easily prove difficult as we have found that these sequences can be quite similar.

Another approach to tackle this problem would be to sequence the dual-coding transcripts from cells grown under a variety of energetic conditions. By varying the conditions and identifying which conditions the alternative proteins were most prominently expressed in, this could help in elucidating the functions of these proteins. Once a baseline for each set of conditions was determined, the artificial site-specific RNA endonucleases could be engineered to target the unedited transcripts for each dual-coding gene and observe the effects on cell growth in the varied conditions. This could give insight into the importance of the two different overlapped proteins at different points in the trypanosome life cycle.

Another direction to pursue is understanding the mechanisms of gRNA selection. We observed many promiscuous gRNAs, and it is possible that even more promiscuity exists in the system. To better understand this, it would be necessary to deep sequence all of the edited mRNAs and have paired gRNA transcriptome data. Ideally, this would be done in multiple cell lines, and under various energetic conditions, as we saw significant variation in our two cell lines and media conditions. This would allow us to see the full picture of alternative mRNA editing and gRNA usage. Then we could begin to determine how much gRNA promiscuity plays a role in gRNA selection, or if other factors are in play.

## **Conclusion**

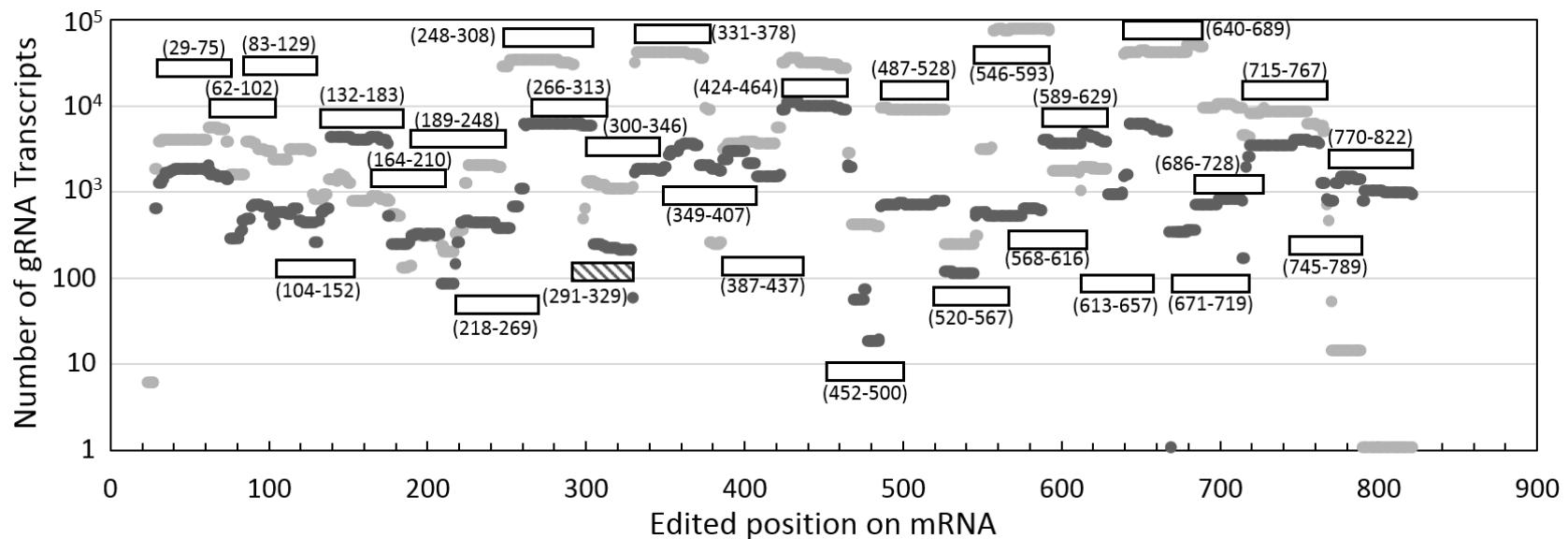
This work found that alternative editing is pervasive in *Trypanosoma brucei*, and brought to light the use of overlapped reading frames, providing another strong reason for the utility and maintenance of RNA editing. This work also showed that the RNA editing system is surprisingly robust and tolerant of mutational noise.

## **APPENDICES**

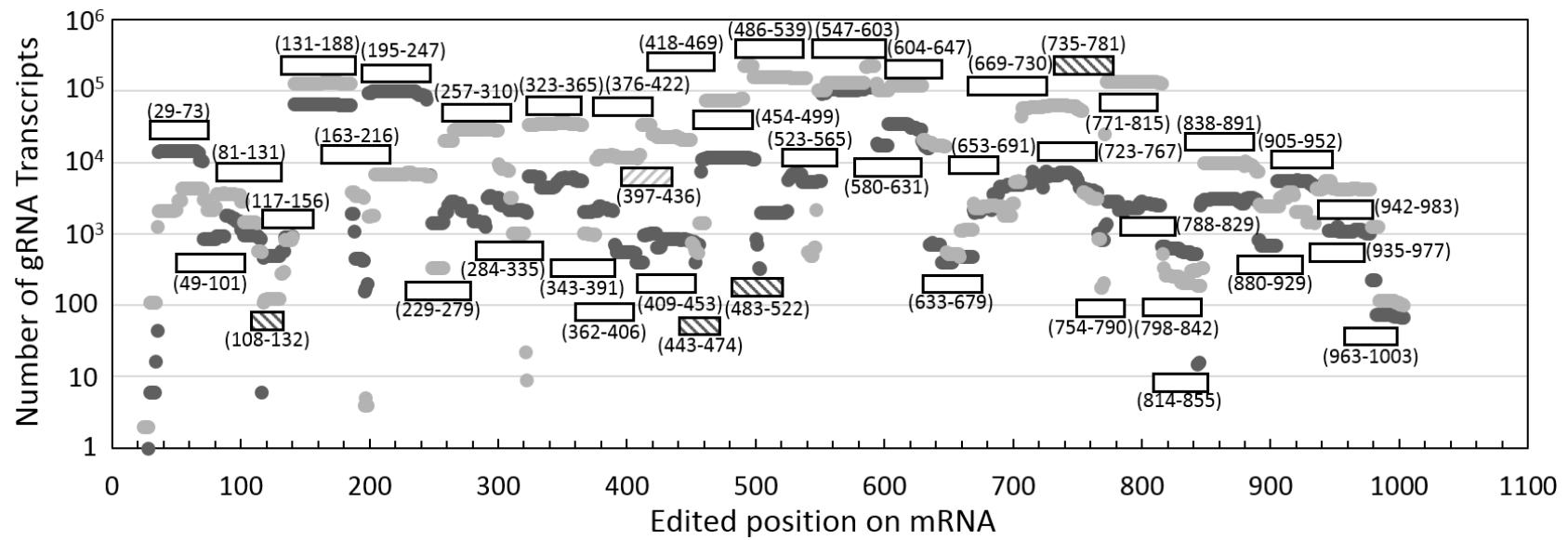
**APPENDIX A. Quantification of the number of identified bloodstream and procyclic gRNA transcripts that cover a respective nucleotide in the fully edited mRNA.**

Bloodstream gRNAs are shown in dark gray and procyclic gRNAs are shown in light gray. Nucleotides and deletion sites were both numbered as edited positions in the mRNA transcripts starting from the 5' end (+1 =0). Boxes indicate the positions of identified populations of gRNAs (coverage ranges shown in parenthesis). Boxes with dark gray or light gray diagonal stripes indicate populations identified only in the bloodstream or procyclic transcriptomes respectively. A. ATPase subunit 6; B. Cytochrome oxidase III; C. C-rich region 3; D. C-rich region 4; E. NADH dehydrogenase subunit 3; F. NADH dehydrogenase subunit 7; G. NADH dehydrogenase subunit 8; H. NADH dehydrogenase subunit 9; I. Ribosomal Protein S12. All individual data points were designated with solid circles. Close overlapping of individual data points generate the observed solid lines.

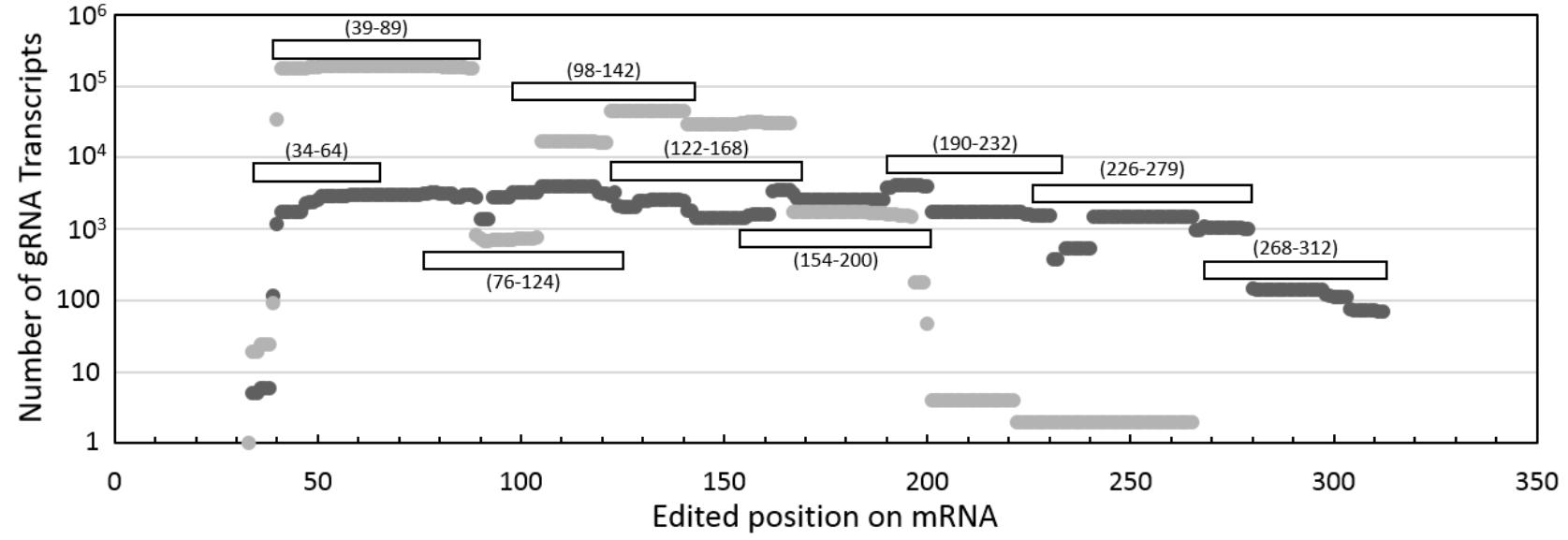
A. ATPase subunit 6



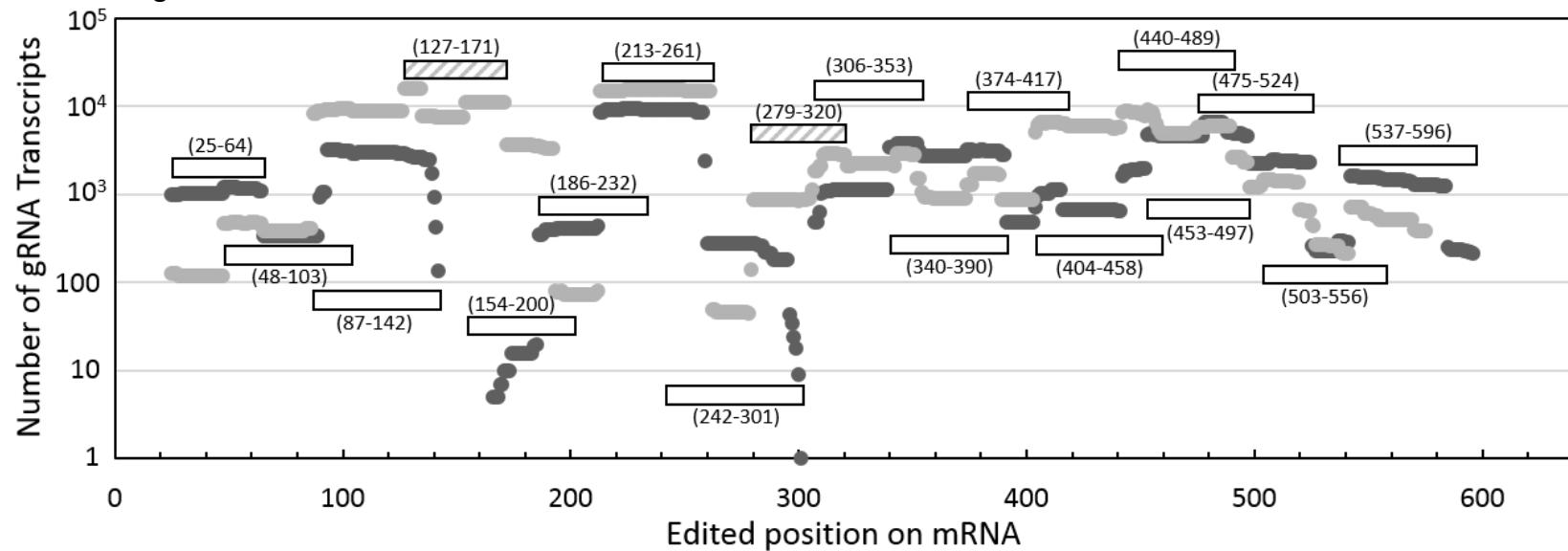
B. Cytochrome oxidase III



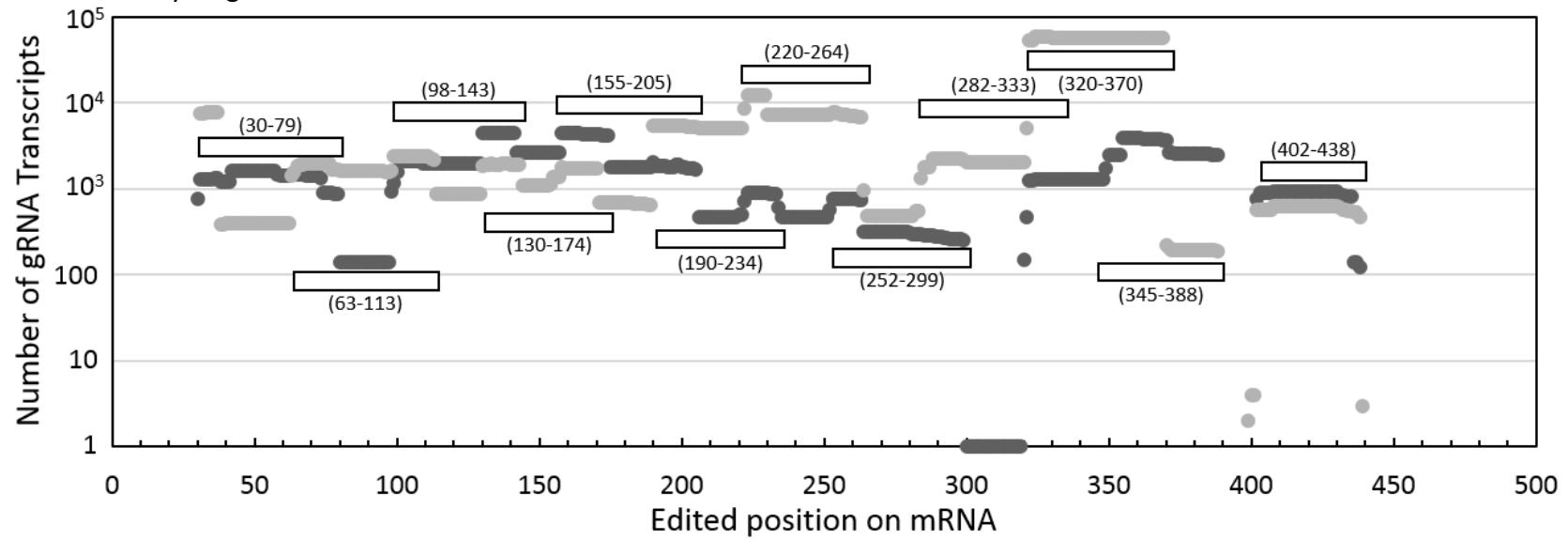
C. C-rich region 3



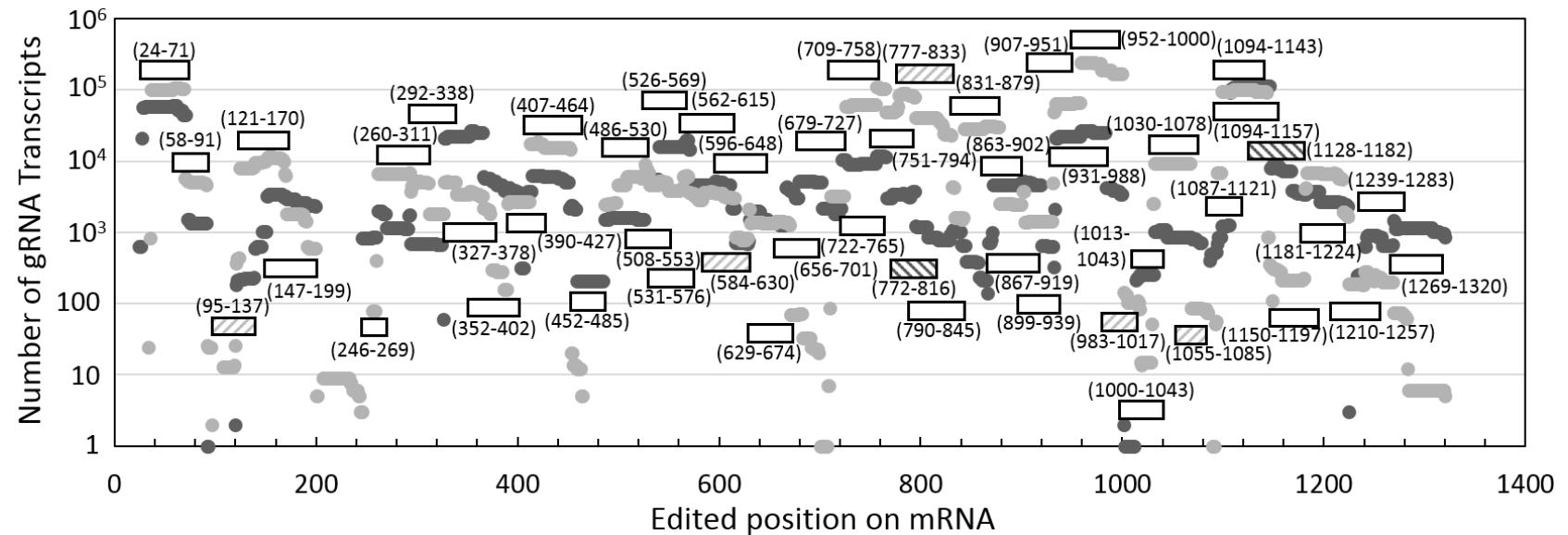
D. C-rich region 4



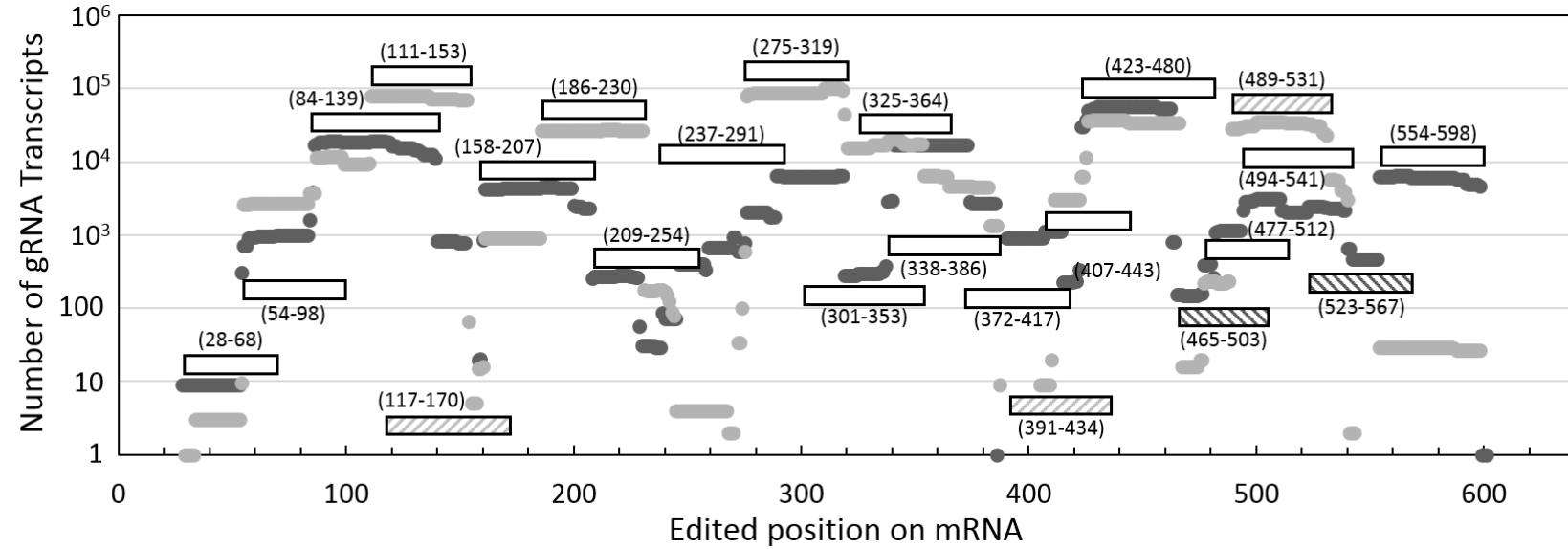
E. NADH dehydrogenase subunit 3



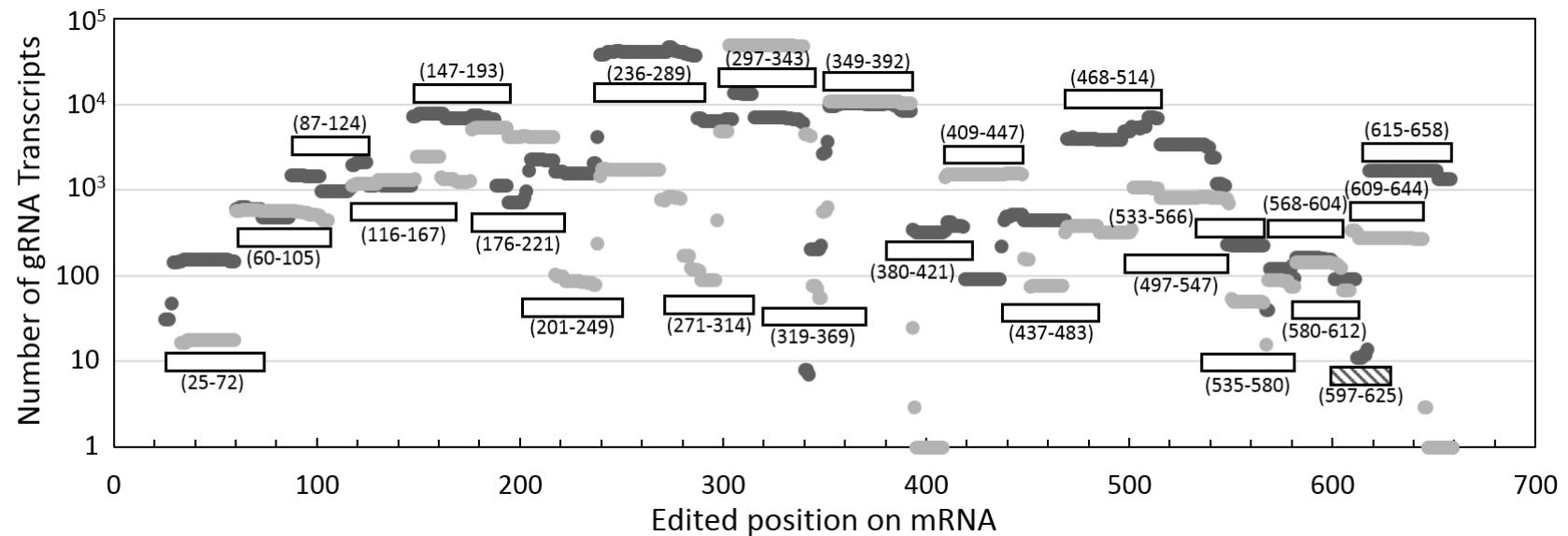
F. NADH dehydrogenase subunit 7



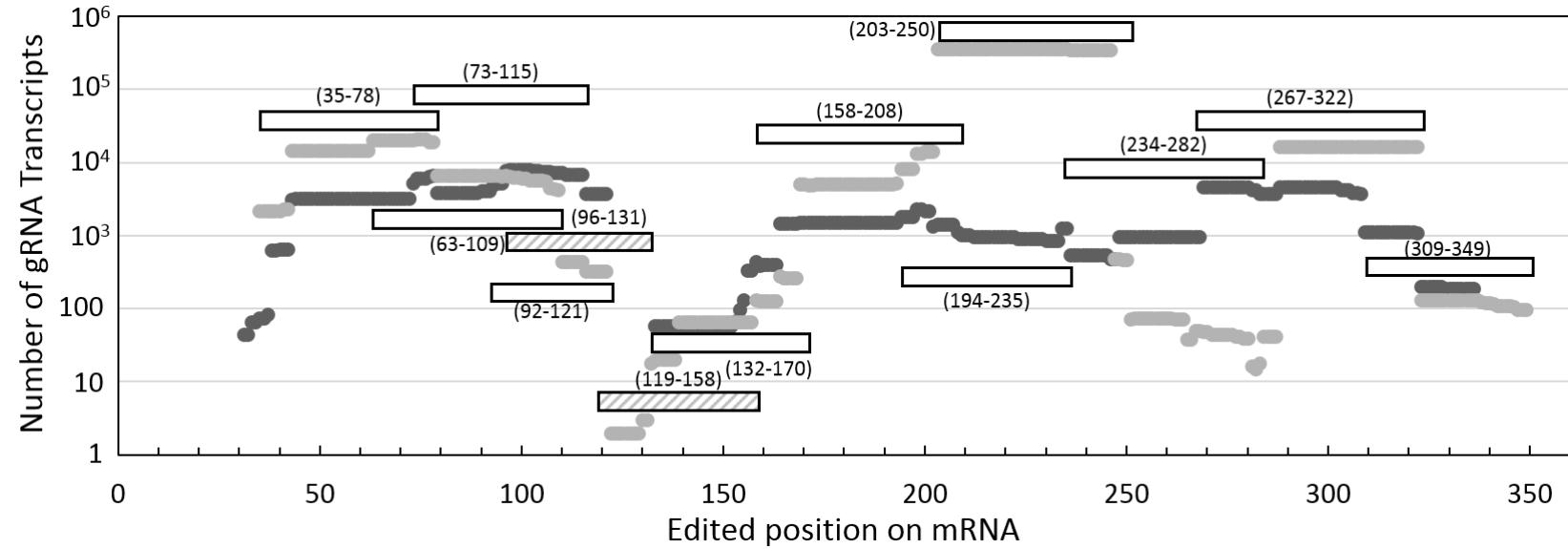
G. NADH dehydrogenase subunit 8



H. NADH dehydrogenase subunit 9



I. Ribosomal Protein S12



## APPENDIX B. Alignment of the mitochondrial fully edited mRNAs and the most abundant gRNAs required for full coverage identified in the bloodstream (blue) and procyclic (gray) life cycle stages.

Conservative mutations between gRNAs are shown in green and mutations that disrupt alignment are shown in red. Lowercase u's indicate uridylates added by editing, asterisks indicate encoded uridylates deleted during editing. Nucleotides and deletion sites in the fully edited mRNA were numbered starting from the 5' end (+1=0). Watson-Crick (|) and G:U (:) base pairs are indicated. Mismatches are indicated by the number sign (#). A) ATPase 6; B) Cytochrome Oxidase III; C) C-Rich Region 3; D) C-Rich Region 4; E) Cytochrome b; F) Maxicircle Unidentified Reading Frame II (Murf II); G) NADH Dehydrogenase Subunit 3; H) NADH Dehydrogenase Subunit 7; I) NADH Dehydrogenase Subunit 8; J) NADH Dehydrogenase Subunit 9; K) Ribosomal Protein S12.

A) ATPase subunit 6

0 10 20 30 40 50 60 70 80 90  
 AAAAAUAAGUAUUUUGAUAAAAGUAAAuGuuuuAuuuuuuuuuuGuGAuuuAUUUUGGuuGCGuuGuuAuuAuGuAuGuAuuAuuGuGuAu  
 |||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
<sub>11</sub>TTTTATACGGAAGTGAAAGGGAACTAAATGAGACCAACGCAAC**ATATA** 5' pA6 (29-72)  
<sub>11</sub>**TTTTGTATAGAAATGAAGAGAACGGTACTAGGTAAACCAATGCAAATATA** 5' bsA6 (29-75)  
 ||:||||:||||:||||:||||:||||:||||:||||:  
 pA6 (62-102) <sub>11</sub>TAATTAGTGCAGATAGTGATATACATGATGACACATA  
 bsA6 (62-100) <sub>22</sub>TAATTAGTGTAGATAATGATAACATAATAGTAACACATA  
 :||||:||||:  
 pA6 (86-127) <sub>10</sub>TTATAGTAGTATGTA  
 bsA6 (90-129) <sub>10</sub>**TATAGTATATG**  
 100 110 120 130 140 150 160 170 180 190  
 GAuCuAGGuuAuGuuuuAuuGuGuAuuuuAAuUGuuuAAuGuuGAuuuuuuGuuuuGuuG\*UUUGAuuuGuAuuuGuuuGuuGGuuuGu  
 ||||:  
 CTAAC**CATA** 5' pA6 (62-102) pA6 (164-208) <sub>12</sub>TATTATGTGGTAGAT-AGACTGAATATAGATAAGCAACTAAACA  
 C-**AAAATA** 5' bsA6 (62-100) bsA6 (164-208) <sub>14</sub>**TATTAGTAGAT-AGACTGAATATAGACAGATAACCAAACA**  
 :||||:||||:||||:||||:  
 TTGGATCTAGTATAAGATGACACATAAA**TATA** 5' pA6 (86-127) pA6 (192-243) <sub>04</sub>**TATTAAGTG**  
**TTAGATTCAATACGAGATAGCACATAAAATATATA** 5' bsa6 (90-129) bsa6 (190-243) <sub>19</sub>**TTAATTAAAGTG**  
 ||:||||:||||:||||:||||:  
<sub>12</sub>TAATA**GAAGATAGTGTATAGAATTGACAGATTGCAACTAAAAACTACATA** 5' pA6 (113-152)  
<sub>12</sub>**TTTTAATATAGAATGGTCATGAAATTGACGAGTTACAACATAAA-CTATA** 5' bsA6 (105-148)  
 ||:||||:||||:  
<sub>11</sub>**TGATATAGTTAGAAGTTGGAAGATAATGAGACAGAC-AAACTAAACATATA** 5' (138-183)  
<sub>10</sub>**TAATAAGTGGTAGTTAGAGACTGGAAAATAGTAAACAAAC-AAAT-AAATA** 5' bsA6 (139-175)



500            510            520            530            540            550            560            570            580            590  
 AuuAuCAuCCCAuUUUUUAuuGuuGAuGuuuuuuGUuuuuuuuUAuuuuuGuuuuuuuuuuuuAuGGuGuuuuuuGuuAuuGAuuuAuuuuuAuuu  
 |||||:|||||:|||||:  
 TAATAGTAGGGTAGAAGATAACAACTAAACATA 5' (487-526)            pA6 (568-611) <sub>12</sub>TTATTATAGAAGATAGTGAACTGAATAAAATAAG  
**TAATAGTAGGGTAGAAGATAACAACTAAACATA 5' bsA6 (487-526)**            **bsA6 (576-616)** <sub>07</sub>**TACATATAGAATAGTGACTGGATGAAATGAA**  
 :|||||:|||||:|||||:  
 pA6 (521-567) <sub>05</sub>**TAATTGTAAGAGACTGAAAGAAATAAGATAAAAACAAAAAAACAAAAAAACAAAAAA 5'** (589-629) <sub>13</sub>**AATTTAGTGAAAATGAA**  
**bsA6 (520-553)** <sub>04</sub>**TATATAACTGTGAGAGACTAAGAAGAATGAAATAAA-CAAAAAAACAAAAAA 5'** (589-629) <sub>10</sub>**AATTTAGTGAAAATGAA**  
 |||||:|||:#|||:  
**TTTAAA 5' bsA6 (458-500)**            pA6 (557-593) <sub>09</sub>**TATAAATGAGAGT-AAGAGAGAAAATGTCACAAGAAACAATAGCTAAATAA 5'**  
**bsA6 (546-592)** <sub>11</sub>**TAAATGAGAGTGAAAAGAGAGAGATACCGTAGAAGACAATAACTAAATA 5'**  
  
 600            610            620            630            640            650            660            670            680            690  
 AuuuuuGuGuuuuGuuuuuGuuuuAuuAuuuAUGuGuuuuuAuAuUUGuuGGGuuuAUUuGCC\*\*\*GCCAuAuuAC\*\*\*AGuuAuuuAuuuuuuGuAAu  
 |||||:|||||:  
 TAAAAACACAAATCATA 5' pA6 (568-611)            pA6 (680-714) <sub>13</sub>**TAAT**----**TTAGTGAATAGGAGATATTG**  
 TAAAAACACAAACAAATA 5' (576-616)            **bsA6 (671-714)** <sub>14</sub>**TTAATG**---**TTAGTAAGTGGAGAATATTA**  
 |||||:|||||:  
 TAAGAGCATAAGGCAAAGACAAATAATAAATA 5' pA6 (589-629)            pA6 (698-728) <sub>11</sub>**TAATAAGAGATA**GTG  
 TAAGAGCATAAGGCAAAGACAAATAATAAATA 5' bsA6 (589-629)            **bsA6 (699-727)** <sub>11</sub>**TAATAAGAATATGA**  
 :|||||:|||||:  
<sub>12</sub>**TTAAAAGTGAATGATAAAATGTACGAGAATATAGACAACTTAATATATA 5'** pA6 (613-654)  
<sub>11</sub>**TTAAAAGTGAATGATAAAATGTACGAGAATATAGACAACTTAATATA 5'** bsA6 (613-654)  
 |||||:|||:  
 pA6 (640-689) <sub>12</sub>**TTATATAGATAGTTGAGTAGATGG**---**TGGTATGATG**---**TCAATAATATATA** 5'  
**bsA6 (643-667)** <sub>15</sub>**TAAGTAGTCAGTAGATGG**---**CGTTATAGT**---**TCAATAATATATA** 5'  
  
 700            710            720            730            740            750            760            770            780            790  
 AuGAuuuuGCAGuuGAuAAuGG\*\*AuuuuuuGuuGuuuuuGuuGuuuGuuuAGuuuuGuAuuuGAuuuuuGAuAGuuAuuAuAuuGuuGuuGAAuuuG\*  
 |:|||||:  
 TGCTAGAACGTCAACATAAA 5' pA6 (680-714)            pA6 (770-822) **TCTCTTCTTCCCTTTAATAGTATAGTGACAGTTTAGAC-**  
 TACTAGAACGTCAACATAGA 5' bsA6 (671-714)            **bsA6 (773-822)** <sub>09</sub>**TTAATAGTATAGTGACAGTTAGAC-**  
 |||||:|||||:  
 TATTGAGATGTTAGCTATTACC--**TAAAAT**TA pA6 (698-728)  
**TGTTAGAATGTCAATTATTACC--TAAAAT**ATA 5' bsA6 (699-727)  
 :: |||||:|||||:  
<sub>11</sub>**TTT**--**TAAAGAGTGATAGAAATAGCAGATAAGTCAAGACATAAACTAAATA 5'** pA6 (720-767)  
<sub>14</sub>**TTTATT**--**TAGAGAGTAGCAAAGACAGTAAGTAGATCAAACATAAT-ATATA 5'** bsA6 (717-763)  
 :|||||:  
 pA6 (747-789) <sub>11</sub>**TTAAATTAGAGTATAAGTGGAAGCTGTCAATAATATAACAACATAAAA 5'**  
**bsA6 (747-789)** <sub>11</sub>**TTAAATTAGAGTGAAGTTGGAGACTATCGATAATATAACAACATATATA 5'**

800            810            820            830            840  
\*GuuUGuuA\*\*UUGGAGUUAUAGAAUAAGAUCAAAUAAGUUAAUAUA  
  :|||:||||  |:|||||||||  
-TAAGCAAT--AGCCTCAATATCA**GG** 5'  
-TAAGCAAT--AGCCTCAATATCAT**ATA** 5'

Alternate initiating gRNA (procyclic transcriptome only)

750            760            770            780            790            800            810            820            830  
uAGuuuuGuAuuuGAuuuuuGAuAGuuAuuAuAuuGuuG\***uGAAA\***uuG\*\*GuuUGuuA\*\*UUGGAGUUAUAGAAUAAGAUCAAAU  
  :||||:||||:|||||  |:|||  |||  :||||:|||  |||||||||  
pA6 (774-822) \*<sub>14</sub>TAATAGTATGGTGAC-ATTTT-GAC--TAAAGCAGT--AACCTCAATATCATA 5'

B) Cytochrome Oxidase III

0            10            20            30            40            50            60            70            80            90  
 GGUUUAUUGAGGAUUGUUUAAAAuGuuuuuGuuuC\*\*\*\*GuuGuAuAuuuGuuGGuGuuA\*\*\*\*GuGGuGuuuuuGuu  
 |||:|||||||:|:|:|:|||  
 pCO3 (35-70) <sub>11</sub>TATG-TAGTTAAGAAAATGCAGAGATAGAG----CAACATATAATTAAATA 5'  
 bsCO3 (36-70) <sub>12</sub>TATATATGGTAGAAAAGAGATAACAAGAATAGAG----CAACATATAATATATA 5'  
 ||| :|:|||||:|:|:|:|||  
 pCO3 (54-101) <sub>11</sub>AAATAG----TAGTATATAGGCAGTTAGT---CACTACAAAAACAA  
 bsCO3 (51-99) <sub>09</sub>AAAG----TAGTATATAAACAGTTACGAT---CATCACAAAAACAA  
 | :|:|:|:|:|:  
 pCO3 (81-116) <sub>10</sub>TT---TACTATAGAAGTGA  
 bsCO3 (88-115) <sub>07</sub>TCTATAGAAGTAA  
 100        110        120        130        140        150        160        170        180        190  
 uuuuuAuCuuuACCuGCCAuuGuuAuuGuGuAuuGGuuAuuuuGuuuGuuG\*\*\*\*GGAuuuAuuuGuuAuuGUUUG\*\*\*\*GuAGuuuuuuAuuuGuuGA  
 ||| :|:|:|:|||:|:|:|:|||:|:|:|:|||:|:  
 AATATA 5'            pCO3 (141-185) <sub>12</sub>TAATTTAGTAGATAAT----CTTAGATGAGCAGATAACAAAC---CATTATATA 5'  
 TATA 5'            bsCO3 (141-185) <sub>14</sub>TAATTTAGTAGATAAT----CTTAGATGAGCAGATAACAAAC---CATTATATA 5'  
 ||| :|:|:|:|:|:|#|:|:|:|:|:|:  
 AAGAGTGGAAATGGACGATAACATAATATATA 5' (81-116) (163-203) <sub>10</sub>TAGATATAGGGATAGTAGAT---CGTCAAGAGATAACAAACT  
 AGAGATAGAGATGGATAGGTAGCAATA 5' (88-115)            bsCO3 (168-195) <sub>12</sub>TATAGTGAAT---TATCAGAGAATAGACTACT  
<sub>09</sub>TAGATGGATGGTGACAATAACATATATA 5' bsCO3 (108-132)  
 ||| :|:|:|:|:|:|:|:|:|:  
 (117-156) <sub>13</sub>TATATAGTGTAGTGATACATAGCCAATGAGACAAACAC---CTATATATA 5'            pCO3 (195-247) <sub>09</sub>TAATT  
 (117-155) <sub>16</sub>TGATAGTAGTAGTGACACGTGATCAATAAGACAAACAC---C-ATATACA 5'            bsCO3 (195-244) <sub>13</sub>TAATT  
 |||:  
 200        210        220        230        240        250        260        270        280        290  
 uuGuG\*\*\*\*GuuuuuAuuuuuuuuuuuuGuuGGuuuuuGuAuuuGuuuGuuGuuGuuAuuGuuAGAuuuGuuuuuGuGAuuuuuuuACGuGGuuuAuuuGAuuu  
 |||:  
 AATAA---AA 5' (163-203)            pCO3 (258-299) <sub>10</sub>TATAATTTAGATAGAGCATTGAAAGATGTCAAATAAACTAAA  
 AACAC---CAAAATATATA 5' (168-195)            bsCO3 (259-300) <sub>13</sub>TAATTTAGATAAGATATTGAGAAGTGCACTAGATAAACTAAA  
 ||:::    :|:|:|:|:|:|:|:|:  
 AGTGT---TAGAATGAGAGAAAGAACGACTAAAAACATAAAATAATA 5' (195-247)  
 AGTGT---CAAGATGGAAAAGAAGTAGCCGAGAACATAAACAATATA 5' (195-244)  
 ::|:|:|:|:|:|:|:|:|:  
 pCO3 (229-274) <sub>10</sub>TTAGAGATATAGATAGACAATGATAGTGACAATCTAAACAAACATATA 5' (293-320) <sub>10</sub>TAATTAGA  
 bsCO3 (236-279) <sub>11</sub>TTATAAATGAATGACAATGATGATGCTAGACAAGACACTAAAATATA 5'            <sub>11</sub>TAATTGAA



500            510            520            530            540            550            560            570            580            590  
 uuAGAuuuAuuuAAuuuGuuGAuAAAuACAUuuuAUUUGuuUGuuAGuGGuuuAuuuGuuAAuuuuuuuGuuuuGuGUUUUUGGuuuAGGuuuuuuuGuu  
**AATCTAAGTAAATTAAACAACATAATAA 5' (483-522)**  
 :|||:|||:|||:|||:|||:|||:|||:|||:  
 GATTTAAGTAGATTAGGTAACGTATGTAAATAATATA 5' (491-539)  
**GATGTAGGTAATTGAGCAGCTGTTATGTAAATATA 5' (504-535)**  
 :|||:|||:|||:|||:|||:|||:  
 pCO3 (528-565) **12TATAGTAAAGATAGATAGACAATCACTAAAGTGAACAATTAAAATATA 5'**  
 bsCO3 (525-563) **13TAATATGTAGAGTAGATAAACAGTCACTAGATGAACAATTAAATATA 5'**  
 :|||:|||:  
 pCO3 (548-592) **11TTTAAATGAGTGATTAGAGAGACAAGATACAAGAACTAAATCCAAATATATA 5'**  
 bsCO3 (551-593) **13TAATAGATGATTGAAGAGACAGAGATACAAGAGCCAATCCAAATAACA 5'**  
  
 600            610            620            630            640            650            660            670            680            690  
 G\*\*UUGuuGuuuuGuAuuAuGAuuGAGuuuGuuGuuG\*\*\*GuuuuuuGuuuuuuGuGAAACCAGuuAUGAGA\*\*GUUUGCuuGuuAuuuAuuACAuA  
 |   ||:|||:|||:|||:  
 C--AATAACGAAGTATAACTAACTCAAGATATA 5' (585-629)            pCO3 (669-717) **09TTTT--TAAATGTGACGGTAGATGATGTAAT**  
**T--AATAGTAGAGCATAATACTAACTCAATATATA 5' (585-628)**            bsCO3 (669-715) **08TTTT--TAAATGTAGTAGTAAATGATGATGATG**  
 :|||:|||:  
<sup>13</sup>TATAATAGAATATGATGTTAGTTCAAGTAGCAAAC---CAAAAA#CAAAACATA 5' pCO3 (604-647)  
<sup>13</sup>TATAATAGAATATGATATTAGTTGAACAGCGAGC---CATAAAACAAAACATA 5' bsCO3 (604-643)  
 :|||:  
 pCO3 (635-669) **14TAAT---TAGAGAGTAGAGATATTTGGTCAGTACATT--CAAACGTAACATATA 5'**  
 bsCO3 (635-669) **11TAAT---TAGAGAGTAGAGATATTTGGTCAGTACATT--CAAACGTAACATATA 5'**  
 :|||:  
 pCO3 (659-691) **12TAACATAGATACTTGGTTGATACTTT--TAAGCGTAACAATAATTATA 5'**  
 bsCO3 (653-682) **15TAGATAGTAGTGTGTTGGTCATATTCT--CAAGTGTATA 5'**



900            910            920            930            940            950            960            970            980            990  
 GCGAuuuGuuuAuuuuGAuGuuuuAuGuGuuAuGuGuuGuGuAAuuuAuuGGuGuuuuUUUAGUUGuuGAuuA\*GuuAAuuuGuAuuGGUAGUU  
 |||||:|||||:||||| | :|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:  
 CGCTAGACAAATAAGACTA**AATATA** 5' (880-918)        pCO3 (963-1003)    <sub>08</sub>TATATTGGAATTAATGGCTAGT-CAATTGAATATAACCATCAA  
**TGCTAACAGATAAGACTACAAAATATATA** 5' (880-929)      **bsCO3 (965-1003)**    <sub>17</sub>**TGTAAATTAATAGTTAGT-CAATTGAATATA**GCCATCAA  
 ::|||||:|||||:|||||:|||||:|||||:|||||:  
<sub>13</sub>**TGAAATAGAATTGTAAGATGCATAATACATAGACACACATAT**AATATA** 5' pCO3 (907-947)  
<sub>13</sub>**TATATAGTTAGAAGATACATGATACATAAGTACACACATTAAAAGATA** 5'<sub>bsCO3 (920-952)</sub>  
 ::|||||:|||||:|||||:|||||:|||||:  
<sub>12</sub>**TTAAATGTACATGTTAGAGTAACTATAGAAAAGTCACAACTAAATATA** 5'<sub>bsCO3 (940-977)</sub>    <sub>14</sub>**TAATTAGTGTATATTGAGATAGTCACAGAAGAATCAACAACTAAATATA** 5'  
 ::|||||:|||||:|||||:  
<sub>14</sub>**TAATTAGTGTATTAGAGTAACTGCAAGAGAAATCGACAACAAAT-CAATATA** 5'<sub>bsCO3 (951-981)</sub>    <sub>10</sub>**TAATGTACA-TATAGTGACTATAAGAAGATCAACAACTAAAT-CATA****

1000  
 UGUAGGAAG  
 |||||  
 ACATATA 5'  
**ACATATATA** 5'

### C) C-rich region 3

0 10 20 30 40 50 60 70 80 90  
 AGAAAUAAAUAUGUGUAUGAUAAAAuGuuuGA\*\*\*\*UUGGuuuGGuuuuGuuGuuuuUUUAuuGuuuGuuuGuACAuuuuuuuuuGuuuuuu  
 |||:|:||| |||:|:|||||:  
 pCR3 (33-62) <sub>09</sub>TAAAGTGAGATTATAGACT---AACGAGCCAAACAACATGTATA 5'  
 bsCR3 (34-62) <sub>09</sub>**TAGTGTGAT**-**GTAAACT**---**AACAAGTCAAACATATA** 5'  
 | |::|:|:|:|:|:|:|:|:|:|:|:  
 pCR3 (41-88) <sub>14</sub>TAT---AGTAGATTAAAGTGACAAGAGAGTAACAGGCAAACATGTAAA-TATATA 5'  
 bsCR3 (41-89) <sub>08</sub>TTT---AGTAGATTAAAGTGACAAGGGAATAACGAGCAAACATGTAAA**GATATA** 5'  
 | :|:|:|:|:|:|:|:  
 pCR3 (78-118) <sub>16</sub>TATCATAGTATGTGGAGAGAGTAAAGA  
 bsCR3 (78-118) <sub>04</sub>**TATCATAGTAT**GTGGAGAAAGTAAAGA  
 pCR3 (105-140) <sub>13</sub>TATAGTTAT  
 bsCR3 (105-140) <sub>12</sub>**T**TAGTTAT  
 100 110 120 130 140 150 160 170 180 190  
 AuuuGuuuGuG\*\***A\*\***UUUGuuuuuAuGuuuGuuA\*UUUAGuuuuuGuuuuuuAuuGGAuuuuuGuuuuuuAuuuAAuAuGGGuuuAuUGuuGuGuuuA  
 |||:|:||| | |||:|:|:|:|:|:|:|:|:  
 TAAACAAACAC---T--AATATATA 5' (78-118) pCR3 (154-196) <sub>13</sub>TTAATTTAGAAGTAGAGAGGTGAATTATGCTCAAATAACAACATATATA 5'  
 TAAACAAACAC---T--AATATATA 5' (78-118) (162-200) <sub>15</sub>**TACTATAGATAGAAGATAGATTATGCTCAGATGACAACACAA**ATATA  
 ::|:||| | | :|:|:|:|:|:|:|:|:|:|:  
 AGATGGAGCAC---T--GAACGAGAATACAAACAAT-AGATA 5' (105-140) pCR3 (190-230) <sub>06</sub>**TTAATGTAGAT**  
 AGATAGAGCAC---T--GAACGAGAATACAAACAAT-AGATA 5' (105-140) bsCR3 (192-232) <sub>11</sub>**TATAGAT**  
 | :|:|:|:|:|:|:|:|:|:|:  
<sub>13</sub>TATAGAATATGAGTAAT-AGATCGAAGACAGAGAATAACCTAAAGATA 5' pCR3 (122-166)  
<sub>09</sub>**TATAGAGTAAT-GAATCA**AAGACAAGAGATAATCTAAAAACAATATA 5' bsCR3 (129-167)  
 200 210 220 230 240 250 260 270 280 290  
 uuuuuuuuuuuuuAuuuuAuCAuuuGAuAuGuGuAuCA\*AAuuGuuAuuGuuAuuuAG\*UUCGuUUA\*UAuuGuuAuuuUUAuAuuuAuuuAAGUAUGC  
 :|:|:|:|:|:|:|:|:  
 GAAAGAAGAGAATAGAATGGTGAACTATACAACATATA 5' (190-230) (293-308) <sub>12</sub>TAATTAGT-**GAAATGATA**GTGATTAGAGTCATACG  
 AAGAGAGAGAATGGAATAGTAGACTATACACAATAGATA 5' (192-232) (268-312) <sub>05</sub>**AAATAGTA**ATGGAGATATTGAGTGAATTCGTACG  
 |||:|:|:|:|:|#|:|:|:|#|:|:|:  
<sub>12</sub>AATATATATAGT-TTCATGATATGTAAAGGTC-AAGCAAAT-ATAACAAATAATATA 5' pCR3 (226-277)  
<sub>09</sub>**AAACAGT**-TTGATGATAGGTAAAGGTC-AAGCAAAC-ATAACAAATATA 5' bsCR3 (234-265)  
<sub>14</sub>**TT-TCTATGATA**GGTAGGTAGATC-AAGCAAAT-ATAACAATAAGATA 5' bsCR3 (241-279)  
 300 310  
 AAAAAUUUUUGU POLYA  
 |||:|:  
 TTTATTAAAATATATA 5' pCR3 (293-308)  
 TTTATTAAAAATATA 5' bsCR3 (268-312)

D) C-rich region 4. Resequencing of the mRNA indicated that there were 2 errors in the original sequence (yellow highlights).



E) Cytochrome b

0            10            20            30            40            50            60            70            80            90  
GUUAAGAAUAAUGGUUAUAAAUUUUAUAAAAuAuGuuuCGuuGuAGAuuuuuAuuAuuuuuuuuAuuAuuuAGAAuuuGuGuuGUCUUUUAAUGUCAG  
:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
<sub>12</sub>TGATAGGTGTCGTATAAGTAGTATTAGGGATAATAA<sub>AAAAAAA</sub> 5' pCYb (32-64)  
<sub>06</sub>**TATAGGTGTCGTATAAGGTAGTATTTGAGGATAATAA<sub>AAAAAAA</sub> 5'** bsCYb (32-64)  
|||||:||||:||||:||||:||||:||||:||||:||||:||||:  
pCYb (53-91) <sub>11</sub>TTAATAAGGGAAATAATGAGTCTTAAGTGTGACAGA<sub>AAAAAAA</sub> 5'  
bsCYb (54-91) <sub>05</sub>TAT**CAATAGGAGGGTAATGAGTCTT****GATGTAACAGA**<sub>AAAAAA</sub> 5'

F) Maxicircle unidentified reading frame II

0            10            20            30            40            50            60            70            80            90  
UUUUAUAUAGAAAGGUUAUAAUCUAUAAUGAuuuGuuuGGuuGuuuuA\*\*\*\*AuuuAGuuuuAuuuUUGuGCUUUGAUUGuAGUCGUGUUUUUGA  
:|||||:||||:||||:||||:||||:||||:  
pMURF2 (30-79) <sub>11</sub>TTTAAAATTGTTAGGTAATGAGAT---TAAATTAAAATAAAACACGAAAGATA 5'  
bsMURF2 (30-79) <sub>08</sub>TTTAAAATTGTA**AGT**TAATGAGAT---TAAATTAAAATAAAACACGAAAGATA 5'

G) NADH Dehydrogenase subunit 3

0 10 20 30 40 50 60 70 80 90  
 UCAAAAAAUCCUCGCCUUUUUACUUUAGUUUGUUAUCAuuAuuuuuAuAuuuGuuuuUG\*A\*UAuuGuGGuuuA\*\*UUAuuuuAuuuAuAGGuuuuuuuu  
 |||:|||:|||:|||:|||:|||:  
 pND3 (33-76) 12TATAGTAGTAAAGATGTGGTAAAAGC-T-ATAGTACCAAAT--ATTATATA 5' (99-143) 11TA  
 bsND3 (30-73) 12TATAATAGTGATAAAAGATGTGAATAAAGAC-T-ATAACACCAAAT--TATA 5' (98-141) 12TAA  
 |||:|||:  
 pND3 (63-113) 11TTAATATTGAAT--AGTGAGATGAGTGTCTAAAAGAA  
 bsND3 (63-108) 09TTAATATTAGAT--AGTGAGATGAATATTCAAAGAAAA  
 100 110 120 130 140 150 160 170 180 190  
 uAuGuuuuuuAuGuuuuuuAuGuuuuuuGuuGCAuuuuuuuGuuCGuuGuuGuuuGuGGuuuCGuGuGuGGuUUGuAuGAAuAuGAAuUCACGuuuG\*GUGuuuu  
 |||:|||:  
 ATACAAAGAACATACA 5' (63-113) pND3 (158-205) 15TAAGTGTATTAGATGTGCTGTGTTGAGTGTAAAT-TACAAAA  
 ATACAAAAA-TACATA 5' (63-108) bsND3 (158-205) 13TAAGTGTATTAGATATGCTGTGTTGAGTGTAAAT-TACAAAA  
 |||:|||:  
 ATATGGAAAATATGAAAGATAGTGTGAAGAAACTAACAAAAGTATA 5' pND3 (99-143) pND3 (190-229) 13TAC-TATAGAA  
 ATATAGAGAACATAGAAATGGCGTAAGAGACTAACAAAAGAATA 5' bsND3 (98-141) bsND3 (190-233) 13TAC-TATAGAA  
 |||:|||:  
 18TATAATTGATGGAAGTAGCAGTAGATACTAGAACGACACTAAC-TATATA 5' pND3 (130-170)  
 11TAATTAGTAAGAGTGACAATAAACACTAGAGGCACATCAAACATATATA 5' bsND3 (130-174)  
 200 210 220 230 240 250 260 270 280 290  
 AuACAuuGGGuuuuAUGGuuuuGuuAGuUGuUUGGuuuuuuGuAuuGuuAAAuuCCAUUAuuGuGuUUUGuuGuuuGuuuuuGUGAuA\*GuGuuGuuuuAu  
 |||:  
 TATGTATATA 5' (158-205) pND3 (253-299) 14TCTTTAATAGATATAAGATAGTGAGCAAGAGTTATTAT-CACAACAAAATAGTATA 5'  
 TATGTATATA 5' (185-205) bsND3 (253-299) 06TTTAATAGATATAAGAGTGACAGATAGAAACATTAT-CACAACAAAATAATATA5'  
 |||:  
 TATGTAGTCTAGATATAAGACAATCAACAATATATA 5' (190-229) pND3 (284-329) 06TAAT-TGTAATAAGATAG  
 TATGTAGCTTAGATACAGAGTAGTCACAAACAATATA 5' (190-233) pND3 (285-328) 10TT-TATAATAAAATAG  
 |||:  
 pND3 (223-263) 11TTTAGTAAGTAGGAGATATAGCAATTAGGGTAATAACATATATT 5'  
 bsND3 (222-263) 16TATTAATAGATAGAGAGTGAGCAATTAGGGTAATAGACATATAAA 5'



#### H) NADH Dehydrogenase subunit 7

0 10 20 30 40 50 60 70 80 90  
 UGAUACAAAAAAACAUGACUACAUGAUAuCAuuuuAuGuuAuuuuuGGuAGuuuuuuuACAuGuuGuAuCGuuuuACAuG\*GUCCACAGCAuCCC  
 :||||::|:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
 pND7 (36-69) <sub>14</sub>TATTATAGT-GAATACGGTGAGAGTTATCAGAGAAATGTAATAATATA5 (108-137) TTAGATTTTAGAG  
 bsND7 (28-71) <sub>12</sub>TATTATAGTAAGATGCAATGAAAGCCGTCAAGAGAAATGTAACATATAAAA 5'  
 ||:||||:|:||:||:||:||:||:||:||:||:||:||:||:||:  
 pND7 (59-91) <sub>13</sub>TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA5'  
 bsND7 (58-91) <sub>12</sub>TAAAGTGTAAATATAGCGAAGTGTAAAT-CAGGTGACATAATATATA5'  
 100 110 120 130 140 150 160 170 180 190  
 G\*\*\*CAGCACAuG\*\*GuGuuuuAuGuuGuuuAuGuGuuGA\*AuuuAuGuuuA\*\*UAUUGAuUGuAuGuuAuA\*\*\*G\*GuuAUUUGCAUCGUG  
 | #||#||||| |:::|||||||||:|||:  
 C---TCATGTAC--CGTAAATACAACAAATAATATA 5' pND7 (108-137)  
 :||||:||:||:||:||:||:||:||:||:||:||:||:||:  
<sub>14</sub>TTAATAAGTGTATGAAGATGCCATT-TAGATAGCAAAT--ATAACTACATA 5' pND7 (124-170)  
<sub>16</sub>TATATAGTAAATGACATGGAAGTGTACT-TAAATAACAAAT--ATATATA 5' bsND7 (121-166)  
 |||:||:||:||:||:||:||:||:||:||:||:  
 pND7 (152-190) <sub>12</sub>TAATAGTGAGT--ATAGTTGACATGGTAT---C-TAATAATACGTAGCATTTAA5'  
 bsND7 (151-199) <sub>12</sub>TAAAGTGTAGAT--GTGATTGATATGATGT---C-CAATAA-ATGTAGCATAAA5'  
 200 210 220 230 240 250 260 270 280 290  
 GUACAGAAAAGUUUAUGUGAAUAAAAGUGUAGAACAAUGUCUUCGGuAUUUCGACAGGUAGAuGuuGuA\*GuGuuuGuuGuAAuGAGCAuuuGuuGu  
 :||||:||:||:||:||:||:||:||:  
 pND7 (246-269) <sub>14</sub>AAAATAAGGAAATCTATGAGGCTGTTCAATACACAACTATA 5'  
 bsND7 (246-269) <sub>12</sub>AAAATAAGGAAATCTATGGGGCTGTTCAATACACAACTATA 5'  
 ||:||:||:||:||:||:||:||:||:||:||:  
 pND7 (261-311) <sub>08</sub>TTTTAATATAGT-TACAAGTGACATTATTCGTGAATAACA  
 bsND7 (261-293) <sub>10</sub>TTTTAATGTAGT-TATAAGTGACATTGCTCGTGA-CACA  
 ||:  
 pND7 (297-338) <sub>13</sub>TACATA  
 bsND7 (292-324) <sub>13</sub>TGAAATAGTG

300            310            320            330            340            350            360            370            380            390  
 CuuuA\*\*\*UGuuuuGAGuAuAuGuuGCGAuGuuGuuuGuCGuuACGuuGuGCAuuAuGCGuuAuuGuA\*\*\*GAuuuAC\*\*\*CCGuAGuuuuA  
 |||||        |||||  
 GAAAT---ACAATATA 5' (261-311)            pND7 (352-398)        07TATAATGTGCAGATGATTAACGT---CTTAAATG---GACATCAAAAG-ATATA5'  
**GAA-T---ATA-TATA 5' (261-293)**            **bsND7 (353-402)**        13AAAATGTGCAGATAATTAATGT---CTTAGATG---GGTATCAAAATTATATATA5'  
 ||:||        ::||||:||:|||||:|||||:|||:  
 GAGAT---GTAAAGCTTATATGCAACGCTACAACAAATATA 5' (297-338)            pND7 (390-424)        16TAATTG---AGTATTGAGAT  
**GAAAT---ATAAGACTCATATACGA-GCTACAACAA-TATATA 5' (292-324)**            **bsND7 (391-424)**        14TAATG---GATATTAAGAT  
 :||||:||:||:||:||:||:||:||:||:||:||:||:  
 pND7 (327-373)        13TATTATAGTAAGTAGCAGTGTGACGTGTAATATGCAAATAATTAATATA 5'  
**bsND7 (327-365)**        12TAATTGTAGTAAGTAGCAATGTAACCGTAGATATGCAAATAACATA 5'  
  
 400            410            420            430            440            450            460            470            480            490  
 AuGGuuuGuuGuGuAuAuCAuGuAuGGuuuuGG\*AuuuAGGuuGuuuGuCUCCGuuG\*UUuGAuCAuuuGAGGAA\*\*\*CG\*UGACAAuGuGACAu  
 ||:||:||:||:||:|||||:  
 TATTAGATAGCGCATATAGTACATAAAATTATA 5' (390-424)            pND7 (486-530)        11TTTTAATTGTTGTAA  
**TATCAAGTGACATATAGTACATAACATTAAA 5' (391-424)**            **bsND7 (486-530)**        11TTTTAATTGTTGTAA  
 :||||:||:||:||:||:||:||:||:||:  
 13TTATATAGTATATGTCAGAGCT-TAGATCTAGTAAACAGAGGAATATATA 5' pND7 (412-452)  
 13TATAGTGTATGTTGGAACATC-TAAGTTCGATAGACAGAGGCAAT-A-TATA 5' bsND7 (414-458)  
 :||: ||:||:||:||:||:||:||:||:  
 pND7 (453-485)        13TATAAT-AGTGTAGTGTAGTTCTT---GC-ATTGATTAACTACTGTAA  
**bsND7 (453-485)**        13TATAAT-AGTGTAGTGTAGTTCTT---GC-ATTGATTAACTACAAATAA  
  
 500            510            520            530            540            550            560            570            580            590  
 uuuuGAuuuAuG\*\*UUGuGGuuGuCGuAuGCAuuuGGCUUUCauGGuuuuAuuA\*GGuAUUCUUGAUGAuuuuGuuuuuGuuGAuuuuuuGuuG  
 ||  
**AATA 5' (452-501)**            pND7 (564-615)        12TTTATTAAGATAGAGATTAAGCGGTTAGAAGATAAC  
**ATA 5' (453-485)**            **bsND7 (564-615)**        10TTTATTAAGATAGAGATTAAGCAGTTAGAAGATAAC  
 ||:||:||:||:||:||:||:||:||:  
 GAGACTGAGTAT--AGCATTAACAGCATACGAATATA 5' (486-530)            pND7 (584-630)        15TATAGTTAAAGAGTGAC  
**GAGACTGAGTAT--AGCATTAACAGCATACGAATATA 5' (486-530)**  
 ||: ||:||:||:||:||:||:||:||:  
 14TAATTGTATAT--AGCATTGACGGTATGTGTGAGTCGAAAGTACTAAATATA 5' (508-548)            pND7 (596-642)        18TAAATTTGAT  
 11TTATAT--AGTGCTAGTAGTATATGTAGGTTGAGAGTACCAAAATAATATA 5' bsND7 (508-553)        bsND7 (596-640)        11TTAAT  
 ||:||:||:||:||:||:||:||:  
 pND7 (526-569)        12TAATATGTAAATTGAGAGTATCAGAGTAAT-CTATAAGAACTACTATATATA 5'  
**bsND7 (526-569)**        13TAATATGTAAATTGAGAGTATCAGAGTAAT-CTATAAGAACTACTATATATA 5'  
 ||:||:||:||:||:||:||:  
 pND7 (540-576)        05TACTGATAGTATTGAGATAGT-CTATGAGAGTACTAAAACAAATAATAATA 5'  
**bsND7 (540-571)**        14TAATCGTTAGTGTCAAGATAGT-CTATAAGAACTACTAAATAATATA 5'







1200      1210      1220      1230      1240      1250      1260      1270      1280      1290  
 GCGuGGuuuuuuAuuGCAuGAuuuAGuuGC\*\*\*C\*GuuuuAGGuAAuAuuGAuGuuGuuuuGGAuCCGUAGAUCGuuA\*GuuuuAuAuGuG\*\*A\*\*\*\*  
 |||||:||||:||||||| | :|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
 CGCACTAGAAGATAACGTAAATATATA 5' (1181-1218)      pND7 (1269-1320) <sub>12</sub>TAATTAGTAGT-CAGAATATGTGC--T----  
 TGTATCAGAAATGACGTACTAAAT 5' (1183-1224)      bsND7 (1269-1320) <sub>17</sub>TAATTAGTAGT-CAGAATATGTGC--T----  
 |||||:||||:||||||| | :|:|:|:|:|:|:|:|:|:|:|:  
<sub>14</sub>TAATAATGTGTTAGATCAATG--G-CAAGATTCAATTATAACTACAAACATATA 5' pND7 (1210-1257)  
<sub>09</sub>TAATGTGTTAGATCAAT--G-CAAGATTCAATTATAACTACAAACATATA 5' bsND7 (1233-1257)  
 ||:|:|:|:|:|:|:|:|:|:|:  
 pND7 (1251-1282) <sub>12</sub>TAAGTGGACAGAATTGGTTCTAGCAAT-CAA-TATATATA 5'  
 bsND7 (1240-1270) <sub>16</sub>TTTATTGTAGTTATAGTGAAGACTTAGGCATA--GCAAT-CAAATATA 5'  
  
 1300      1310      1320      1330  
 \*GGUUAUUGuAGGAUUGUUAAAUUGAAUAAAAA  
 |-:|:|||||:|||||||:  
 -CTAGTAACATCCTAACAAATATA 5' (1269-1320)  
 -TTAGTAACATCCTAACAAATATA 5' (1269-1320)

I) NADH Dehydrogenase subunit 8

0            10            20            30            40            50            60            70            80            90  
 CAAUUUAAUAAAAGUUUUGGUUGAUUAuuAuuuuuuuAuuuuuuuGuAuGuuuuuuuuGAuuuuuuuGuuuuuuuuUUUUUGuuuGuuuuuu  
 |||:|:|||:|:|||:|:|||:|:|||:|:|||:  
 pND8 (29-68) <sub>08</sub>TATTATATAGTGAAGAATAGAAAGATAAAGACATACAAAAAA 5' (87-136) <sub>14</sub>TAAATGAGTAAGAA  
 bsND8 (28-56) <sub>04</sub>TAGGGAGATAGTAAAAGAGTAGGAGGTAGGATAAAA 5'      bsND8 (86-139) <sub>12</sub>TAAATGAGTAAGAG  
 :|:|:|:|:|:|:|:|:  
 pND8 (55-98) <sub>11</sub>TTATATAGAGAGAAGTAAAGAACAAAGAAGAAAACAAACAAAATATA5'  
 bsND8 (54-97) <sub>16</sub>TATATATAGAAAGAGACTGAGAACAAAAGAACAAACAAA-TATA 5'  
 100        110        120        130        140        150        160        170        180        190  
 AuAuGuGUuuuGuuuGuuGuGuuA\*\*\*\*CuAUUU\*GuuuA\*\*\*CCCAuuGAGuuAACCAuuGuuAGuuuAuuGGuuCGuGGUAACCAuuuuuuuGCGUUUU  
 |||:|:|:|:|:|:|:|:  
 TATGCATGAGACGAGTAACACAAT---GATAAA-TA 5' (87-136) (161-187) <sub>11</sub>TTAATTAGATAGTCAGTATCATTGGTATAAAACGCAAAT  
 TATATGCAAGATAGGCAGTACAAT---GATAAA-CAAATATA 5' (86-139) <sub>15</sub>TGTAATTAAGTAGTTAGGTATCATTGGTAAAGACGCAAA  
 :|:|:|:|:|:  
 (111-153) <sub>13</sub>TTAAGTAGTGTAAT---GATGAA-TAGGT---GGGTAACTCAAATGGTAAATATATA5' p(186-230) <sub>13</sub>TTAAAGAGTGTAAA  
<sub>12</sub>TATTCAGATAGTATAGT---GATAGA-CAGAT---GGGTGACTTA-TTGGTAACATATA5' pND8 (187-228) <sub>11</sub>TAAAGAGCGTAAA  
 :|:|:  
 pND8 (117-170) <sub>11</sub>TTTATAAT---GATGAG-TAAGT---GAGTGATTAGTGGTAACAAATCAAATAAAC 5'  
 200        210        220        230        240        250        260        270        280        290  
 uAUU\*\*\*GGuGuGGuuuAGAGCGuuGuAuuGCuuGuCGuuuAuGuGAuuuAuuuGCCuA\*\*\*\*GuuuAGCAuuGGAuG\*\*\*UUCGuGuGGGuGGAGu  
 AATATA 5' pND8 (161-187)  
 TGATATA 5' bsND8 (160-199)  
 |||:|:|:|:|:|:  
 ATAG---CTATACTAACGTCTCGCAACATAACATA 5' (186-230)  
 GTGA---CCATATAAATCTCGCAACATATATTATA 5' (187-228)  
 :|:|:|:|:|:  
<sub>06</sub>TAATAGCTATAGTAAGTCTTGAGTATAATGGACAGCAAATACAA 5' (213-244)  
<sub>12</sub>TAGATGTTGCAGATTAGTACCAAATACATAAA 5' bsND8 (219-245)  
 :|:|:|:|:|:  
 pND8 (237-267) <sub>15</sub>TTAAATGTATTGAGTTAGATGGGAT---TAATATGTAACCTAC 5'  
 bsND8 (246-271) <sub>11</sub>TAATGTAGTGAGTAAATGGGAT---TAGATTGATACCTAC---AAGCATA 5'  
 |||:|:|:|:|:  
 pND8 (240-288) <sub>10</sub>TATATATTGAGTTAGATGGGAT---CAAGTCATAGCCTAC---AAGTAT 5'  
 bsND8 (259-285) <sub>05</sub>TATAATGTTAGTATATTGAGTTAGATGGTAT---TAGATCGTAGCCTAC---AAGAATATA 5'





J) NADH Dehydrogenase subunit 9

0	10	20	30	40	50	60	70	80	90
---	----	----	----	----	----	----	----	----	----

UUAAUAUCAACUAAUUUUUUUUUAACAUuAuAUGGuGuAuAuUUUUuGuuuAuuuCGuuuAuGuuuuGuuuAAuuUUAuuuuA\*\*UUGuuuGu  
 :||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
 pND9 (33-71) <sub>14</sub>**TATTA**TAGTATGTATGAAGATACGAGTAAAGCAAATACAAATATA 5'  
 bsND9 (25-72) <sub>11</sub>TTTGTAATATAGTGTATGTGAGAATATAAATAAGCAAATACAAAATATA 5'  
 :||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
 pND9 (60-105) <sub>11</sub>TATTAGTGATATAAGAGTGAATTGAAATAAGAT--AACAAACA  
 bsND9 (60-101) <sub>11</sub>TAAATGTAGTGAATATGGAGATAGATTAAGATAAAAT--AACAAACA  
 :||| :||:||:||:  
 pND9 (87-124) <sub>11</sub>**TTAAT**--AGTGAATA  
 bsND9 (87-124) <sub>10</sub>TATAAT--AGTGAATA

100	110	120	130	140	150	160	170	180	190
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

GuuGuAGAuGGuGuuUUGuuuGuuuuGuuGAuuGuAGuuuuuuGuuuuuuuuAuuGuuuuGuuAGuuuuuuuuuuGuuuuAUUGuAuGuuuuuAuuuuuuuAA  
 ||||:|||  
 CAATATA 5' (60-105) :||||:||||:||||:||||:||||:||||:||||:  
**TA-TATA 5' (60-101)** pND9 (176-216) <sub>14</sub>TAATAGTGTAGAGATAGGGAATT  
 :||||:||||:||||:||||:  
 TAGTATCTGTACAGAGATAAACAAA-CAACTATA 5' (87-124)  
**TGACATCTACTACAGAGCAGACAAA-CACTATAAA 5' (87-124)** bsND9 (176-216) <sub>14</sub>TAATAGTATATAGAGATAAGAAATT  
 :||||:||||:||||:||||:  
 pND9 (117-160) <sub>13</sub>TTAAATAGAGTAATTGACATCGGAAGATAAACAAAATA-TATA 5'  
 bsND9 (117-162) <sub>12</sub>TTAAATAGAATAGTTAGTGTCAAAGAGTAAAGAAATAACAAAACAATATA 5'  
 :||||:||||:||||:||||:  
 pND9 (149-193) <sub>15</sub>**TTGATAGT**AGATAATCAAAGGAAGATAAACATACAAAATAATATA 5'  
 bsND9 (147-187) <sub>15</sub>**TAGA**ATAGTGAGATAATCAAAGAGAAGCAGAATAACATACAATATA 5'  
  

200	210	220	230	240	250	260	270	280	290
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

uuuGuGAuuuuuuGuuuuuuAuAuuGUUGuGAuUUGuuAuGuuGAuuuuuuGuGGuuuuuGuuuuuuGuCGuuuuAuGuuGuuGUuAuAuuuuAuuuuGuuuG  
 :||||:||||:||||:  
 AGACACTGAAAACAAAATATAACATA 5' (176-216) pND9 (272-312) <sub>06</sub>TATACAGTGTATGTGAAATGAGACGAAC  
**AGATGCTAAAACAAAATATAACATA 5' (176-216)** bsND9 (273-314) <sub>12</sub>**TTT**ATAGTAGTATATAAGATGGAACGAAT  
 :||||:||||:||||:  
 TAATATTGAAGATAGAGATATAGTACTAGACAATAACTAAATATA 5' pND9 (201-242) pND9 (303-339) <sub>13</sub>**TAA-**  
**11TAATTAAGAGTAGGGATATGGTACACTAGATAGTAACTAACTAAAAAAA 5' bsND9 (204-249)** bsND9 (305-339) <sub>13</sub>**TAAT**  
 :||||:||||:||||:  
 pND9 (239-279) <sub>12</sub>TTTAATTAGAGATACTGGAAATAGAAGCAGCAGAATACAACA**TATTA**AA 5'  
 bsND9 (239-286) <sub>11</sub>TTTAATTGAAAGTACTAGAGATAAGAGTAGCAAGATAACAAAATATA 5'



600        610        620        630        640        650        660        670  
 GCAuACC\*\*AuuUUUAuuuG\*CuuAuuuuAuuuA\*\*\*AuA\*\*\*UCACCGuUGUAUUUCUAAAUUCCACUUCC  
 |||||  
 CGTATATATA 5' (568-604)  
**TATATA 5' (569-600)**  
 :||||| ||||#|||:||| |||||  
 TGTATGG--TTAAT**TATAGAC-GTAATATA 5' (582-612)**  
**CGTGTGG--TTATTAATAGAC-GTAA 5' (582-611)**  
**CGTGTGG--TTGAAGATAAAC-GTAA 5' (597-625)**  
 ||||:||:|||: |||||:|||:|||     ||| #||||:|||||||  
<sub>12</sub>**TTAAAGAGTAAAT-GTAATGAAGTGAAT---TAT---GTGGTAACATTAAGAATATATA 5' pND9 (609-644)**  
<sub>12</sub>**TTTAAAGTGAAT-GTAATAAGATAGAT---TG---ATAG 5' bsND9 (609-640)**  
 :||||: ||||:||:|||:||     ||| |||||:|||:|||  
<sub>12</sub>**TGAAAT-GTAGTGAGATAAGT---TAT---AGTGGTGACATTAGGAAATATATA 5' pND9 (615-659)**  
<sub>12</sub>**TTAAAGTTAGT-GTAATAAGATAAGT---TAT---AGTGGCACATTAAGTATATA 5' bsND9 (618-658)**

K) Ribosomal Protein S12

0	10	20	30	40	50	60	70	80	90
---	----	----	----	----	----	----	----	----	----

CUAAUACACUUUUGAUAAACAAACUAAAGUAAAGuGuGuGA\*UUUUUGUAUG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGu  
 |||:||||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:  
 pRPS12 (35-76) <sub>12</sub>TATTTAGAGTGGAGAGACGTACATT-GAACACATGC-CAACAAATA (96-121) <sub>12</sub>TATTATAGTA  
 bsRPS12 (38-78) <sub>18</sub>TAGTGAAGAGAGTGATATGCT-AAAGACATAC-CAACAATATATA (96-121) <sub>11</sub>TAATAGTA  
 |||:|||:|||:|||:|||:|||:|||:|||:|||:  
 pRPS12 (43-78) <sub>14</sub>TATATAGTTAGAAGATGCATGTACT-AGAAGTATAC-CAACAACATATA 5'  
 bsRPS12 (43-78) <sub>12</sub>TATATAGTTAGAAGATGCATGTACT-AGAAGTATAC-CAACAACATATA 5'  
 ::|::: :|::: :|::: :|::: :|::: :|::: :|::: :|:::  
 pRPS12 (63-109) <sub>12</sub>TATAGTATGT-TAATGATAGATG-TGAGACAGAATAAACAA  
 bsRPS12 (66-99) <sub>07</sub>TAGTAGTGT-CAATAGTAAATG-CAAACAAATAATA5'  
 ::|::: :|::: :|::: :|::: :|::: :|:::  
 pRPS12 (74-106) <sub>12</sub>TAATATGTCATAGTAGATG-CAAGATAAGATAAACAA  
 bsRPS12 (73-115) <sub>16</sub>TCTTTATAGTAAATG-TAGAGCAGAGATAGACA  
  

100	110	120	130	140	150	160	170	180	190
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

uuuAuGuuAuuAuAuGAGuCCG\*\*CGAuuGCCAGuuCCGGuAACCGACGuGuAuuGuAuGC\*\*\*C\*\*\*\*GuAuuuuAuuUAuuAuAAuuuuGuuuGGGuGu  
 |||:|||:  
 AAATATAATATA 5' (63-109)  
 |||:|||:  
 pRPS12 (169-208) <sub>12</sub>TATATAGAGTGAATATGTTAGAATGAGCCTATA  
 bsRPS12 (156-207) <sub>12</sub>TTATATG--G---TATAAGATAGATGTGTTAGAATAGACTTACA  
 AAATATA**TATA** 5' (74-106)  
 GAATACAATAATATATA 5' (73-115)  
 ::|:::  
 AAATGCAATGATATATTAGGC--ACTAA 5' (96-121)  
 GAATGTAGTGATATATT**CAGGT**--AGCTAACGTGTCAAATATA 5' (96-121)  
 ::|:::  
 AAATGCAATGATATATTAGGC--GCTAACGGATTAAGATATA 5' (96-131)  
 ::|:::  
<sub>10</sub>TATTCAGT--GTTAGTGGATTGAGGCTATTGGTTGCACATAACATTCA 5' (119-158)  
 ::|:::#||:||#||:|||:  
 TTTTAATGTGTTAGGATCATTGGCTGCATATGATATACG--G---CAATATA 5' pRPS12 (139-170)  
<sub>11</sub>TTTAGT**AGATTGAGGCTATTGGTTGCACATAACATTCATA** 5' bsRPS12 (133-158)

200            210            220            230            240            250            260            270            280            290  
 uGCuuGuuuuuuuuGuuGuuuuAuuGGuuuAGuuAuG\*\*UCauuAuuuAuuAuAGA\*\*\*GGGUGGUuGGuuuGuuGAuuuACCC\*\*\*G\*\*\*\*GuG\*UAA  
 ||||||||| :  
 ACGCAACAATAAATA 5' (169-208) pRPS12 (267-322) 07TTTAAAGTGACTAGAT**AGG**---T---CAT-ATT  
**ATGCAACATATA** 5' (156-207) **bsRPS12 (288-322)** 12TTTAAAGTA**ACTAGATGG**A---T---CAT-ATT  
**bsRPS12 (269-308)** 10TTAGTACAAGAGCAGTTAAATGGG---C---TAC-ATT  
 ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||: ||||:  
 14TATAT--AGTAGTGAATGATGTCT---TTTACCGTCAAACAACTAGAATATA 5' pRPS12 (234-280)  
 15TATAT--AGTAGTGA**G**TGAT**A**TCT---TTTACCGTCAAACAACTAGAATATA 5' **bsRPS12 (234-280)**  
 ||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
 ATGTAGTAGAGAAGATGACGAAATAGTCAAATCAAT-C--AGTAATATATA 5' (194-235)  
**ATGTAGCAGAGAAGATGACGAAATAGTCAAATCAAT-C--AGTAATATATA** 5' (198-235)  
 :||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
 14TATAATAGAGAGAGTAACGGAGTAATCGAGTCATGC--AGTAATA**T**ATATA 5' pRPS12 (203-246)  
 07TATAATAGAGAGAGTAAT**GGG**ATAG**T**AAATCAATAC--AGTAAT**TAA**ATA 5' **bsRPS12 (203-245)**  
  
 300            310            320            330            340            350  
 AGuAuuAuACA\*CG\*\*UAuuGuAAGGuuAGA\*UUUAGGuuAUAAAGAUAUGUUUUUU [AUUA] POLYA  
 ||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
 TTATAGTATGT-GC--ATAACATATA 5' (267-322)  
**TCATAATATATA** 5' (269-308)  
**TTATAATGTGT-GC--ATAACATATA** 5' **bsRPS12 (288-322)**  
 || |: ||||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:  
 16TAAGT-GT--ATAATGTTCAATCT-AAGTCTATATTCTATACAATATAAA 5' pRPS12 (309-349)  
 24TAAGT-GT--GTAATGTTCAATCT-**AA**ATCT-TAGTCTATACAAATAAA 5' **bsRPS12 (309-336)**

**APPENDIX C. All gRNA major classes pulled for ATPase 6 in the EATRO 164 procyclic (shaded gray) and bloodstream (white) transcriptomes.**

Populations of gRNAs are bordered boxes. A) ATPase 6; B) Cytochrome Oxidase III; C) C-Rich Region 3; D) C- Rich Region 4; E) Cytochrome b; F) Maxicircle Unidentified Reading Frame II (Murf II); G) NADH Dehydrogenase Subunit 3; H) NADH Dehydrogenase Subunit 7; I) NADH Dehydrogenase Subunit 8; J) NADH Dehydrogenase Subunit 9; K) Ribosomal Protein S12.

**A) ATPase 6**

5'	3'	Reads	ATPase 6 gRNA Sequences
31	75	2,044	AT ATAAACGTAACTGAAATGAATCACGAGAGAAGATAAAGATATAT AT <sub>12</sub>
29	72	1,630	ATATAC AACGCAACCAGAGTAAATCATGAAGGGAAAGTGAAAGGCATATT T <sub>11</sub>
31	75	143	AT ATAAACGTAACTGAAATGAATCGCAGAGAAGATAAAGATATAT AT <sub>21</sub>
31	75	489	AT ATAAACGTAACCAAATGGATCATGGAAGAGAAGTAAAGATATGT AT <sub>09</sub>
29	75	395	AT ATAAACGTAACCAAATGGATCATGGAAGAGAAGTAAAGATATGTT T <sub>11</sub>
35	75	224	AT ATAAACGTAACCAAATGGATCATGGAAGAGAAGTAAAGAT T <sub>13</sub>
29	62	203	ATATACAACGCAACC AGATAAAATCATGAAGAGAAGTGAAGGTATATT T <sub>09</sub>
33	75	144	AT ATAAACGTAACCAAATGGATCATGGAAGAGAAGTAAAGATAT T <sub>12</sub> *
31	62	95	AACGCAACC AGATAAAATCATGAAGAGAAGTGAAGGTATAT AT <sub>14</sub>
39	75	82	AT ATAAACGTAACCAAATGGATCATGGAAGAGAAGTAA T <sub>16</sub>
37	75	40	AT ATAAACGTAACCAAATGGATCATGGAAGAGAAGTAAAG TCAACTTAAT <sub>08</sub> *
62	102	1,435	ATACA ATCATACACAGTAGTACATATATAGTGATAGACGTGATTAA T <sub>11</sub>
62	100	154	ATAAAAA CATAACAAATGATATATACATAGTAATAGATGTGATTAA T <sub>23</sub>
62	102	92	ATATAA ATCATACACGATAATATATGCCCTAGTAAACAGATGTGATTAA T <sub>18</sub>
74	102	10	ATATAT ATCATACACGATAATGCATATGTAGTAAAC T <sub>13</sub>
64	100	9	ATAAAAA CATAACAAATGATATATACATAGTAATAGATGTGATT T <sub>14</sub> *
86	127	2,158	ATAT AAATACACAGTAGAATATGATCTAGGTTATGTATGTGATAT T <sub>11</sub>
90	129	177	ATATA AAAATACACGATAGAGCATAACTTAGATGTTATGTATA T <sub>20</sub>
84	118	85	ATAATAATACAC ATAGAACATGACCTAGATTGTACATAGTGATATAT T <sub>12</sub>
82	118	38	ATATAATAATACAC ATAGAACATGACCTAGATTGTACATAGTGATATAT T <sub>13</sub>
86	118	28	ATATAATAATACAC ATAGAACATGACCTAGATTGTACATAGTGATAT T <sub>15</sub>
83	118	22	ATAATAATACAC ATAGAACATGACCTAGATTGTACATAGTGATATATA AT <sub>18</sub>
91	129	21	ATATA TAAAATACACGATAGAGCATAACTTAGATTGTATATGAT T <sub>18</sub>
113	152	743	ATAC ATCAAAAATCAACGTTAGACAGTTAAGATATGTGATAGAA GATAAT <sub>12</sub>
105	152	54	AC ATCAAAAATCGACATTGAGATAATTGAGGTATGTGATAGAGTATAATT T <sub>11</sub>
105	148	128	ATATC AAAATCAACATTGAGCAGTTAAAGTACGTGCTAAGATATAATT T <sub>12</sub>
116	152	84	AC ATCAAAAATCAACATTGAGCAATTGAGGTACATGATA TGATATAAT <sub>09</sub>
104	148	16	ATATC AAAATCAACATTGAGCAGTTAAAGTACGTGGTAAGATATAATT T <sub>07</sub> *
138	183	430	AT ATACAAATCAAACAGACAGAGTAATAGAAGGTTGAAGATTGATAT AGT <sub>11</sub>
144	177	210	ATATC ATCAAACAAACAGAATAATAGAGAATCAGAGGT GAATGTTAAGT <sub>15</sub>
139	175	3,468	ATAAA TAAACAAACAAAATGATAAAAGGTAGAGATTGATG GTGAATAAT <sub>08</sub> *
132	177	205	ATAT ATCAAACAAACAAAGTAATAGAAAGTCAGAGATTGATGTTAATAA T <sub>10</sub>
164	208	54	ATAT ATAAACACAAATCAACGAATAGATATAAGTCAGATAGATGG TGTATTAT <sub>12</sub> *
176	210	36	AT AAACAAACACAAATCAGTAGACCGAGTACAAGT GAGATGGACGTATAGAT <sub>07</sub>
165	208	24	ATAAT ACAAACACAAATGATAGACCAATACGAGTTAGATGGACG TAT <sub>06</sub>
158	208	14	ATAT ACAAACACAAACTGACGAATAGATACAGATTAAGTGAATGAAATAAT T <sub>11</sub>
164	208	218	AT ATAAGCACAAACCAATAGACAGATAAGTCAGATAGATGA TTAT <sub>14</sub>
192	243	172	ATATAAATTAAACAACATAGATTACAGTAGATAGAAGTAAATGTGAATTA T <sub>04</sub>
218	248	147	ATC AGACTATGTGAGTTAGATGACGTGAATTATA CTGTATAT <sub>12</sub>
190	243	52	ATATAAATTAAACAACATGAATGATGATAAAGTAAATGTGAATTAAT T <sub>19</sub>
189	243	10	ATATAAAATTAAACAACATGAACTATGATGATAAAGTAAATGTGAATTAATG TACT*
207	243	7	ATATAAATTAAACAACATGAACTATGATGATAAAGGT T <sub>10</sub>
193	246	5	*ATTGTATAATTAAACAACATGAACTATGATGATAAAGTAAATGTGAATT T <sub>03</sub>
224	269	864	ACATAA TAATACAATAATACGAGATTAGACTATGTGAATTAAATGATATGA T <sub>11</sub>
226	269	808	ACATAA TAATACAATAATACGAGATTAGACTATGTGAATTAAATGATAT T <sub>13</sub>
221	262	149	ATATAT ATAATACAAAATGAACTGTATAAGTAGACAAATGTGAATT TT
219	262	105	ATATAT ATAATACAAAATTGAACTGTATAAGTTAGACAAATGTGAATTAT T <sub>14</sub>
218	262	57	ATATAT ATAATACAAAATTGAACTGTATAAGTTAGACAAATGTGAATTATA T <sub>21</sub>
221	267	32	AC ATACAATAATACAGAATTAAACTGTGTAAGTTAGATAGTGTAAATT T <sub>10</sub> *

5'	3'	Reads	ATPase 6 gRNA Sequences cont.
248	292	25,157	ATATA AAATACAAATTGAGTAGGTACTACAATGATATGAGATTA T <sub>13</sub>
253	298	5,405	ATATAT ACAACAAATATAGATTCAAGTAAGTGATGTAGTAATATGA T <sub>11</sub>
255	292	244	ATATA AAATACAAATTGAGTAGGTACTACAATGATAT TATTATTAAT T <sub>15</sub>
252	292	134	ATATA AAATACAAATTGAGTAGGTACTACAATGATATAGA TTATTAAT T <sub>07</sub>
262	304	4,881	ATAT ATACAAAACAACAGATATAGATTCGGATAGGTATATGA GATCT T <sub>13</sub>
259	304	395	AT ATATAAAACAACAGATATGAATTCAAGTGAGTGATACAGTA TAT T <sub>15</sub> *
254	298	290	ATATAT AACAAACAAATACGAATTGAGTAGGTAGTATGATGATA T <sub>15</sub>
266	313	586	AAAAAA AAAAAAACAATACAAGATGACAGGTATAAGTTGGATGAGTAAT T <sub>12</sub>
281	313	14	AAA AAAAAAACAATACAGAACATACTAGGTATAGATT AGATATGTGAT T <sub>09</sub> *
284	310	13	ATAT AGAACAAATACAAAATAACGAGTACAG T <sub>08</sub>
283	310	2	TAT AGAACAAATACAAAATAACGAGTACAGG ATAAGTGATAT T <sub>08</sub>
281	312	2	AT AGAAAACAATACAGAACATACTAGGTATAGATT AGATATGTGAT T <sub>19</sub>
291	329	133	AAATATAT AAATGCAATATACGATAGAGAAATGATATAAGATGATAA T <sub>16</sub> *
293	329	12	AAATATAT AAATGCAATATACGATAGAGAAATGATATAAGATGAT T <sub>14</sub>
301	345	647	ATAT AAACAAACAAAAGTAGAAGTGCAGTATATGATAGAAAATGATGT CAAAT T <sub>11</sub>
300	346	263	ATAT AAACAAACAGAAATAGAAATCCAATATACGATAAGAAAATGGTATA T <sub>12</sub>
301	335	125	ATAT ACAAAACAT AAATAAAAGTGCAGTATATGATAAGAGATAATAT T <sub>11</sub>
301	345	33	ATAT AAACAAACAAAAGTAAAGTGCAGTGTATGATAGAAAATGATGT CAAAT T <sub>15</sub>
331	375	24,736	ATAT AATTATTAACAAAGAGAAAGTCACGTAAAAGGTAGAATGAAGATA T <sub>12</sub>
332	378	8,776	AT ATAAATTATTAACAGAAAGAGATCATGTAGAAAGTAGAGATAAT T <sub>14</sub>
331	371	3,561	ATATAA ATAAACAAAAGAAATCACGTAGAACAGAGATAATAGAGATA T <sub>13</sub>
331	375	712	ATAT AATTATTAACAAAGAGAAAGTCACGTAAAAGTAGAATGAAGATA TTAT T <sub>05</sub>
332	378	387	AT ATAAATTATTAACAGAAAGAGTCATATAGAAAATAAGTAGAGAAAT T <sub>12</sub>
331	374	302	ATAT ATTATTAACAAAGAGAAATCATATAAGAGACAGAATGAGATA T <sub>15</sub>
332	378	144	AT ATAAATTATTAACAGAAAGAGATCATGTAGAAAGTAGAGATAAAAAT TTT
331	371	1,547	ATATAA ATTAACAAACAAAAGAAATCACGTAGAACAGAGATAATAGAGATA T <sub>15</sub> *
360	407	181	ATATAT ACATCCATAAAATTATCATCAGTTAATAGATTGTTAAATGAAAAA TTTT
349	389	41	ATATAA ATCACCACAAATAAGTATTGATGAGAGAAAGTTATATA T <sub>12</sub>
352	401	705	ATATAA ATAAAACATCACTAACTAATGGATTGTTAAGTAGAAGAGAATCAT T <sub>11</sub>
360	407	429	ATATAT ACATCCATAAAATTATCATCGGTTAATAGATTGTTAAATGAAAAA T <sub>11</sub>
354	401	160	ATATA ATAAAACATCACTAACTAATGGATTGTTAAGTAGAAGAGAATC CT T <sub>07</sub>
362	407	143	ATATAT ACATCCATAAAATTATCATCGGTTAATAGATTGTTAAATGAA T <sub>12</sub>
349	389	104	ATATAA ATCACCACAAATAAGTATTGATGAGAGAAAGTTATATA T <sub>09</sub>
361	407	74	ATATAT ACATCCATAAAATTATCATCGGTTAATAGATTGTTAAATGAAA TTTTTTTTTTTTTTT
387	435	1,428	ATATAT AACACAACAAAGAACGAATGAGAGAAAGTATCTGAGATTATT T <sub>14</sub> *
387	435	1,049	ATATAT AACACAACAAAGAGACGAATAGAAAAGATATCTGTGAAATTATT T <sub>13</sub>
387	437	934	ATATAT AAAACACAAATAGAAAACGGATAAGAGAGATAATTCTAGAGTTATT T <sub>13</sub>
387	435	664	ATATAT AACACAACAAAGAGCAATAGAAAAGATATCTGTGAAATTATT T <sub>12</sub>
390	435	635	ATATAT AACACAACAAAGAGACGAATAGAAAAGATATCTGTGAAATT T <sub>11</sub>
398	435	28	ATATAT AACACAACAAAGAGACGAATAGAAAAGATATCTGTGA T <sub>05</sub>
397	435	23	ATATAT AACACAACAAAGAGACGAATAGAAAAGATATCTGTGAA T <sub>13</sub>
424	464	25,624	ATAT ATGACACAAACGAGGGAAAGATACTCTAAAGGACACAGTGAAA T <sub>12</sub>
427	467	2,307	ATAT ATAACGACACAAATAGAGAAAGATGCTCTGAGAGATGTATA T <sub>13</sub>
421	460	1,864	ATAAAT TACAACAAAGAAAGATACTCTAGAACAGCACAGTGAGAAAT T <sub>16</sub>
424	457	368	AAATTAACGACA AACAAAGAGAAATACTCTGAGAAATATGATGAAA T <sub>12</sub>
424	464	6,879	ATAT ATGACACAAACGAGGGAAAGATACTCTAAAGGTACAGCGAAA T <sub>13</sub> *
427	468	1,781	ATAT AACAACGATACGACAGAGAAAGATATTCTAAGAGATATGACA T <sub>13</sub> *
424	457	232	AAATTAACGACA AACAAAGAGAAATACTCTGAGAAATATGATGAAA T <sub>12</sub>
5'	3'	Reads	ATPase 6 gRNA Sequences cont.
455	491	368	ATATATAATTAC AAACAAACGAGAGATGTCGGTAAATAATGATATAAT T <sub>11</sub>
455	497	22	ATAT ATTACAAAACAGACGTAAAGATGTCGATGAATGGTGGTATAAT T <sub>14</sub>
452	477	54	AATT ACGTCGATAGATAACGATAATGAG ATTAATTTT
476	500	14	AAATT TAAATTACAAGACAAACGTAGAAGC T <sub>24</sub>
458	500	1	AAATT TAAATTACAAGACAAACGTAGAAGC GTCGATAGATAATGATAT T <sub>15</sub>
487	526	8,723	ATACAA ATCAACAAATAGAAGATGGGTGATAATAGATTGTGAGATA T <sub>17</sub>
487	528	1	ATAC ACATCAACAAATAGAAGATGGGTGATAATAGATTGTGAGATA T <sub>27</sub>
487	526	635	ATACAA ATCAACAAATAGAAGATGGGTGATAATAGATTGTGAGATA T <sub>16</sub>
521	567	232	AA AAAA AAAAAAAAACAAAAATAGAATAAGAAAGTCAGAGAATGTTAAT T <sub>05</sub>
520	553	85	AAAAA AAAAAA AC AAAATAAAGTAAGAAGAATCAGAGAGTGTCAATA TTTTTT
521	553	11	AAAAA AAAAAA AC AAAATAAAGTAAGAAGAATCAGAGAGTGTCAAT T <sub>12</sub>

5'	3'	Reads	ATPase 6 gRNA Sequences cont.
557	593	69,619	AATAAATCGATAACAAAGAACACTGTAAAAGAGAGAA TGAGAGTAAATAT <sub>09</sub>
549	593	2,587	AC AATAAATCAATAACAGAGAAATCATAGAGAGGAAAGATAGAAAT T <sub>12</sub> *
549	592	181	ATAT ATAAATCAATGACAAGAACGACTGTAGAAAAAGAGAGTGAAAAT T <sub>13</sub>
546	592	313	AT ATAAATCAATAACAGAAGATGCCATAGAGAGGAAAGTGAGAGTAAA T <sub>11</sub>
546	593	76	AT AATAAATCAATAACAAAGAACATTGTAAAAGAGAAAAGTGAGAATAAA T <sub>13</sub>
549	592	30	ATAT ATAAATCAATGACAAGAACGACTGTAGAAAAAGAGAGTGAAAAT T <sub>13</sub>
568	611	670	ATACT AAACACAAAAATGAATAAAAATAGTCAGTGTAGAGATATTAT T <sub>12</sub>
576	616	115	AT AAACAAAACACAAAAATAAGTAAGTAGTCAGTGTAGATAAGA TATACAT <sub>07</sub> *
589	629	854	AT AAATAATAACAGAACGGAATACGAGAATAAGTAAAGTGA TTTAAT <sub>13</sub>
589	629	3,292	AT AAATAATAACAGAACGGAATACGAGAATAAGTAAAGTGA TTTAAT <sub>10</sub>
613	654	618	ATATAT AATCCAACAGATATAAGAGCATGTAAAATAGTAAGTGAAAAT T <sub>12</sub>
613	657	183	AT ATAAATCCAACAAGTATAAGAACATATAGAATAGTAGGTGAAAAT T <sub>12</sub>
613	654	459	ATAT AATCCAACAGATATAAGAGCATGTAAAATAGTAAGTGAAAAT T <sub>11</sub>
613	657	399	ATAT ATAAATCCAACAAGTATGAAGACACGTAAAATAGTAATGAAAAT T <sub>14</sub> *
640	689	39,063	ATAT ATAAATAACTGTAGTATGGGGTAGATGAGGTTGATAGATATA T <sub>12</sub>
647	689	678	ATAT ATAAATAACTGTAGTATGGGGTAGATGAGGTTGAT T <sub>11</sub>
640	689	131	ATAT ATAAATAACTGTAGTATGGGGTAGATGAGGTTGATAGATATA T <sub>11</sub>
654	689	234	ATAT ATAAATAACTGTAGTATGGGGTAGATGA TTTGATAGATATAT <sub>12</sub>
643	667	4,401	ACATATATAATAACTGTGATATTGGGGTAGATGGATCTGATGAAT T <sub>14</sub> †
640	668	250	ATATAATAATAACTATAATAAGGTGTAAGTGAGGTTCAGTGAATATA T <sub>14</sub> †
640	662	165	AAAAGTGTATATGGAGGTAAGTGAAATTGATAGATGA TTAATTT†
680	714	7,581	AAAATA CAACTGCAAGATCGTGTATAGAGGATAAGTGATT TAAT <sub>13</sub>
680	719	1,291	ATATATA ATTATCAACTGTGAGATTATATTACAAGGAATAAGTGATT T <sub>13</sub>
672	716	119	ACACA ATCAACTGCAGAATTATATTACAGAGAGTGAGTAATTGTAAT AAT <sub>12</sub>
680	714	105	AAATA CAACTGCAAGATCGTGTAGAGGATAAGTGATT TAAT <sub>11</sub>
671	714	319	AGATA CAACTGCAAGATCATATTATAAGGGTGAATGATTGTAAT T <sub>14</sub> T*
685	714	309	ATATA TAACTGCAGAATCATATTATAAGGGATGAA CGATTGT <sub>13</sub>
698	728	740	ATT AAAATCCATTATCGATTGTAGAGTTATGT GATAGAGAATAAT <sub>11</sub>
686	728	12	ATATATT AAAATCCATTATCGATTGTAGAGTTATGTTAGAGAATAA TAT <sub>21</sub>
699	727	88	AAATCCATTATCAGTTGGAGATTGTA CTATAAAGAATAAT <sub>14</sub>
699	727	24	ATATAT AAATCCATTATTAACGTAAAGATTGTA GTATAGT <sub>17</sub>
720	767	4,588	AT AAATCAAATACAGAACTGAATAGACGATAAAAGATAGTGAGAAATTT T <sub>11</sub>
715	755	2,272	ATATATAT AAACTAAACAAATAGCAGAGACAGTGAGAGATTGTTAT AAT <sub>13</sub>
728	765	920	AT ATCAAATACAAAATGAGCAGATGACAGAGATAGTAA TGATTAT <sub>12</sub>
720	767	165	AT AAATCAAATACAGAACTAGATGAAACATAGAGATAGTGAGAAATTT T <sub>12</sub>
717	763	1,613	ATATA TAAATACAAAATGAGATGACAGAACGATGAGAGATTATT T <sub>14</sub> *
718	763	488	ATA TAAATACAAAATGAGATGACAGAACGATGAGAGATTAT AAT <sub>17</sub>
720	763	452	ATATA TAAATACAAAATGAGATGACAGAACGATGAGAGATT T <sub>12</sub>
720	767	269	AT AAATCAAATACAGAACTAGATGACGATAGAGATAGTGAGAAATTC TTT*
720	767	177	AT AAATCAAATACAGAACTAGATGACGATAGAGATAGTGAGAAATTT T <sub>17</sub>
747	789	13	ATAAT ACAACAATATAACTGTGCAAGGTTGAATATGAGATTAAAT T <sub>11</sub>
747	789	587	ATATAT ACAACAATATAATAGCTATCAGAGGTTGAATGTGAGATTAAAT T <sub>11</sub>
745	789	91	ATATAT ACAACAATATAATAGCTATCAGAGGTTGAATGTGAGATTAAATGA T <sub>11</sub>
790	822	8,663	ATA CTATAACTCCAATGACGAAATCGATTTA CAGTGATATGATAATT T <sub>12</sub>
770	822	1	GGA CTATAACTCCGATAACGAATCAGATTTGACAGTGATATGATAATTATT*
773	822	428	ATATA CTATAACTCCGATAACGAATCAGATTTGACAGTGATATGATAATT T <sub>09</sub>
774	822	223	ATATA CTATAACTCCGATAACGAATCAGATTTGACAGTGATATGATAAT
777	822	195	ATATA CTATAACTCCGATAACGAATCAGATTTGACAGTGATATGAT T <sub>14</sub>

### B) Cytochrome Oxidase III

5'	3'	Reads	COIII gRNA Sequences
35	70	1,179	ATAATT AATATACAACGAGATAGAGACGTAAAAGAAT TGATGTAT <sub>12</sub>
36	73	826	AT ATAATATACAACGAGATGAAGGCATAGAGAAA AGATGGTATATAAT <sub>14</sub>
29	70	112	ATATAC AATATACAACCGAATGAGAATATAAGAAAAGTGTGATA TTAT <sub>11</sub>
36	70	14,200	ATATAT AATATACAACGAGATAAGAACATAGAGAAA AGATGGTATATAAT <sub>13</sub>
54	101	1,386	ATAT AAAACAAAAACATCACTGATATTGACGGATATATGATGA TAAAT <sub>12</sub>
51	99	721	ATATAT AACAAAAACACTACTAGCGTTGACAGATATATGAAAT T <sub>12</sub>
51	99	364	ATAT AACAAAAACACTACTAGCATTGACAAATATATGATGAAAT T <sub>13</sub> *
52	101	229	AT AAAACAAAAACACTGCTAATATCGACGAATATATGATGGAA AAT <sub>14</sub> *
50	92	122	ATATAAAAT ACACCACTGATATCAACGAGTATATGAGATA T <sub>14</sub>
49	95	78	ATATAT AAAACACCACGTGACATCGATAAGTATATGAGTGA TTAAT <sub>15</sub>
81	112	834	GTA GAGTGAAGATAGAGAAAATAAAGATATCGTT T <sub>13</sub>
81	116	550	ATATATAATAACAATA GCAGGTAAGGTGAGAAAGTGAAGATATCATT T <sub>10</sub>
81	131	1	TACATAATAACAGTGGCGGGTAGAGATAGAAGAATAAAGATACTATT T <sub>08</sub>
88	115	443	AACGATGGA TAGGTAGAGATAGAGAAAATGAAGATATT TTAT <sub>05</sub> *
88	115	117	AATAACAATGGA TAGGTAGAGATAAAGAAAATGAAGATATC T <sub>06</sub>
88	115	96	AATAACGATGGA TAGGTAGAGATAGAGAAAATGAAGATATC T <sub>07</sub>
81	112	64	ATACATAACAGTGGCAGA GTGAAGATAGAGAAAATAAAGATATCATT T <sub>11</sub> *
108	132	6	AT ATATACAATAACAGTGGTAGGTAGA T <sub>09</sub>
117	156	104	ATATATA TCCAACAAACAGAGTAACCGATACATAGTGTAGTG ATAT <sub>13</sub>
117	155	465	ACATATA CCAACAAACAGAATAACTAGTCACAGTGTGATG ATAGT <sub>16</sub>
118	150	36	ATATA CAACAAACAAAATAATCGATGCACAGTGTAGT AGTAGT <sub>13</sub>
141	185	126,513	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATAGATGA TTTAAT <sub>13</sub>
141	185	761	ATAT ATTACCAAACAATAGATGAGTAGATTCTAATAGATGA TTTAAT <sub>11</sub> TAAGTTTT*
134	188	542	ATAT AAAACTACCAAACAGTAAATAGATAAGTTCTAATAAGTGTGAGATAATT T <sub>11</sub>
146	185	350	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATA TATGATTTAAT <sub>12</sub>
142	185	199	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATAGATA TTTAATTAT <sub>05</sub>
141	185	183	ATAT ATTACCAAACAATAGACAGTAGATTCTAATAGATGA TTTAAT <sub>13</sub>
145	185	172	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATAGAT CTGATTTAAT <sub>13</sub>
143	185	116	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATAGAT CATTAAATTAT <sub>05</sub>
131	188	109	ATAT AAAACTACCAAACAGTAAATAGATAAGTTCTAATAAGTGTGAGATAATTAT TGTTAT <sub>16</sub>
147	185	106	ATAT ATTACCAAACAATAGACGAGTAGATTCTAAT CGATGATTTAATTAAT <sub>17</sub>
141	185	99	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATAATGA TTTAATTAAT <sub>08</sub>
141	185	62,901	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATAGATGA TTTAAT <sub>14</sub>
141	187	810	ATAT AAAACTACCAAATAATGAACAGATAATTCTAGTGAGTGA TTTAAT <sub>13</sub> *
143	185	685	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATAGAT T <sub>13</sub>
141	185	497	AT ACTACCAAACGATAAGCAGATAAGTCTCAGTGATG TGTAAT <sub>15</sub>
134	188	332	ATAT AAAACTACCAAACAGTAAATAGATAAGTTCTAATAAGTGTGAGATAATT T <sub>10</sub>
141	188	308	ATAT AAAACTACCAAACGATAGACGAATAAGTTCTGATAAGTGA TATAT <sub>12</sub>
142	185	140	ATAT ATTACCAAACAATAGACGAGTAGATTCTAATAGAT T <sub>15</sub>
163	203	2,487	AAA ATAATCAACAAATAGAGAACTGCTAGATGATAGGTGA TATAGAT <sub>13</sub>
163	211	861	ATATT AACCCACAAATCAATAAGTAAGAGACTACTAGATGATGATAA T <sub>14</sub>
168	195	232	ATATATAAAACCACAAATCAT CAGATAAGAGACTATTAGTGATA T <sub>12</sub> *†
185	216	138	ATATAT AATAAAACCACAAATTAGCAAGTAAGAAG GTATCAGATGATAATTAT <sub>06</sub> †
204	247	5,008	ATATAT ACAAAACAAATACAGAGATCGACGAGAAAGAAAGTGTGAGATT TAT <sub>12</sub>
195	247	688	ATAAATAAAATACAAAATCAGCAAGAAGAGAGTAAGATTGTGATTAAT T <sub>08</sub>
199	243	462	ATAT ACAAAATACAAAATCGATAGAAAAGAAAAGTGTGAGATCATGATT TAT <sub>12</sub>
195	247	457	ATAAATAAAATACAAAATCGACAGAGAGAAAAGTAGGATTGTGATTAAT T <sub>12</sub>
195	244	83,016	ATAT ACAAAATACAGCCGATGAAGAAAAGGTGAACTGTGATTAAT T <sub>12</sub>
199	243	9,810	ATATAT ACAAAATACAGAAACTGACGAAAGAGAGAAATGAAGTTATGAT CT <sub>17</sub>
202	247	2,126	ATATAT ACAAAACAAATATGAGAACTAACAAAGAGAGAAAGTGTGAGATTAT T <sub>12</sub> *
201	247	1,676	ATATAT ACAAAACAAATATGAGAACTAACAAAGAGAGAAAGTGTGAGATTATA T <sub>17</sub>
204	247	1,379	ATATAT ACAAAACAAATATGAGAACTAACAAAGAGAGAAAGTGTGAGATT T <sub>13</sub>
199	244	363	ATATAT ACAAAATACAGAAGCCAACGAGAGAAGGAATAAGATTGTAAT T <sub>10</sub>
202	243	269	ATATAT ACAAAATACAGAAACTGACGAAAGAGAGAAATGAAGTTAT T <sub>14</sub>
204	244	149	ATAT ACAAAATACAGAGCCGATGAAGAAAAGGTGAGACT TTGATTAAT <sub>12</sub>
195	243	130	ATAT ATAAATACAAAACTAACGAAAGAAAAGATGGAACGTGTGGTTAAT T <sub>12</sub>
204	243	126	ATAT ACAAAATACAGAAACTGACGAAAGAGAGAAATGAAGTT T <sub>19</sub>
202	244	102	ATAT ACAAAATACAGAGCCGATGAAGAAAAGGTGAGACTGT TATTAT <sub>14</sub>

5'	3'	Reads	COIII gRNA Sequences cont.
229	274	289	ATA TACAAAACAATCTAACAGTGTAGTAACAGATAGATATAGAGATT T <sub>10</sub>
236	279	575	ATAT AAAATCACAGAACAGATCTGTAGTAGAACAAGTAATAAGTAATAT T <sub>11</sub>
229	270	559	ATATAT AAAACAATCTAACAGATGACGGATAGATATAGAGATT TAT <sub>16</sub>
238	268	120	ATATAACCACAAT ACAGATCTGACAGTAATGTGTAGGTTAAAT T <sub>05</sub>
238	279	37	ATAT AAAATCACAGAACAGATCTGTAGTAGAACAAGTAATAAGTAATAT TTCT <sub>14</sub>
234	264	33	ACAAAACAT ATCTAGCAGAACAGTGACGAATAGATACAA T <sub>07</sub>
234	279	29	ATAT AAAATCACAGAACAGATCTGTAGTAGAACAGTAATAAGTAATATAA TTTT
258	299	18,740	ATATAT AAATCAAATAACTATGTAGAAAGTTACGAGATAGATTTAATA T <sub>10</sub>
265	308	4,452	AAA AAAACACAAAAATCAAGTGAACATATGTAGAGGATTGTAAGATAA T <sub>11</sub>
265	310	4,418	ATAAAAACACAAAAATCAAGTGAACATATGTAGAGGATTGTAAGATAA T <sub>13</sub>
258	300	814	ATATAT AAAATCAAATAAATTACGTAGAGGTTACAGAATAAGTTAAT T <sub>10</sub>
267	306	435	ATATAT AACACAAAAATCAGATAGACTATGTAGAAAGATTGTGAAAT T <sub>11</sub>
261	299	181	ATATAT AAATCAAATAACTATGTAGAAAGTTACGAGATAGATT T <sub>08</sub>
258	300	347	ATATAT AAAATCAAATAGATCACGTAGAGGTTAGAATAGATTTAAT T <sub>14</sub>
261	307	328	ATAC AAAACACAGAACATCAGATAGATCACGTAGAGGTTAGAATAGATAA T <sub>08</sub>
262	300	244	ATATATATT AAAATCAGATAAGCCACGTAGAGGATTGTAAAGTGAATT AT <sub>12</sub>
261	300	195	ATATATATT AAAATCAGATAAGCCACGTAGAGGATTGTAAAGTGAATT T <sub>09</sub>
258	289	182	AATCACGTAAAGATCGTAGAATGAGTTAAT T <sub>13</sub>
263	307	71	ATAC AACACAGAACATCAGATAGATCACGTAGAGGTTAAAGATAAAT AT <sub>08</sub>
257	300	39	ATATAT AAAATCAAATAGATCACGTAGAGGTTAGAATAGATTTAATA T <sub>12</sub>
293	320	1,024	AATACTGT ATATGTGTAGTAAGATATAGAGATTAA T <sub>10</sub>
284	321	21	ATATAAA GATACAAACGTAATAAGGCATAGAAGTTAAGTGAATTAT TGT <sub>12</sub>
293	335	1	AAAAAAACAATACTGGATATGTGTAGTAAGATATAGAGATTAA TAACT <sub>06</sub>
291	332	1,923	ATATAT AAACAATACTGGTACGATGTAAAGATGTGAAAGTTAAT T <sub>14</sub>
293	321	151	AATACTGA GATACGACGTGATAAGATATAGAAGTTAA T <sub>11</sub>
323	365	33,179	ATAT ATAAAACAAACTCGCTATGTAAAGAACGTAAAGTGTAAAAAGTGTATT AT <sub>12</sub>
330	365	856	ATAT ATAAAACAAACTCGCTATGTAAAGAACGTAAAGTGTAAAAAGTGTATT T <sub>09</sub>
323	365	301	ATAT ATAAAACAAACTCGCTATGTAAAGAACGTAAAGTGTAAAAAGCGTATT AT <sub>14</sub>
323	365	4,275	ATATAT ATAAAACAAACTCACTGTGTAAAGATTGTAGAAAGTTGTATT AT <sub>24</sub>
323	360	139	ATACAT ATAAACTCACTGCATAAGAACATAGAGAGTGTATT AT <sub>11</sub> *
345	391	522	ATATAT ATAATACAACAAGGAGCGTCATAAGTAAAGTGAATTGTTTATAT T <sub>12</sub>
345	391	203	ATATT ATAATACAACAGAAAATGTCTAAAGTGTGAGATGAATTGTTTATAT T <sub>08</sub>
345	389	490	ATATAA AATACAACAAGAGACGTCGTAATAGAGTAAATTGTTTATAT TTT
349	389	461	ATATATAA AATACAACAAGAGACGTCGTAATAGAGTAAATTGTTT T <sub>06</sub>
347	389	365	ATATATAA AATACAACAAGAGACGTCGTAATAGAGTAAATTGTTTAT T <sub>11</sub>
343	389	262	ATATATAA AATACAACAAGAGACGTCGTAATAGAGTAAATTGTTTATAA T <sub>15</sub>
357	391	106	ATATT ATAATACAACAGAAAATGTCTAAAGTGTGAGATGA TTCGTTTAT <sub>06</sub>
353	389	79	ATATAA AATACAACAAGAGACGTCGTAATAGAGTAAATT TTT
352	388	48	ACAAAT ATACAACAAAAGATGCCGTAGATAAGATAGATTG GTATATTT
354	390	41	ATATAT TAATACAACAGAACAGCTATAAGTGTGAGATAGATT GATTATAT <sub>27</sub>
362	406	101	ATATAT ATAAACATAAAATCAGATAGTACAATGAAGAGTGTATAGATAA T <sub>09</sub>
362	406	139	ATAT ATAAACATAAAATCAGATAATACGTGAAGAGTGTATAGATAA T <sub>06</sub>
376	418	9,942	ACA TATAACACAAAAATAGACATAGACTGAATGATGCAGTGAA T <sub>13</sub>
378	422	1,634	ATAT AACTTACAACACAGAGATAGACATAGATCAGATAATGTGATAA T <sub>13</sub>
384	418	185	ACA TATAACACAAAAATAGACATAGACTGAATGATGCA TTGAAAT <sub>11</sub>
376	418	169	ACA TATAACACAAAAATAGACATAGACTGAATGATGCAAGTAAAT TTTTAACT <sub>06</sub> *
378	422	375	ATAT AACTTACAACACAGAACATAAGCAGATAGTGTGATAA TTAAT <sub>17</sub>
397	426	15	AAAACGAT AGCAAATTCTGACGTGAAAATAGATGTAA T <sub>11</sub>
397	436	10	AAAA AAAAACGAAAGCAGATTACCGTACAGAGATAGATAG T <sub>10</sub>
411	449	20,772	ATATATA ATATAAGTAAATGAGAGACGAGGGTAGACTTGTGATAC TAT <sub>12</sub>
409	438	1,420	ATAATAAGGTAT ACAGAGAACGGAAGCAGACTTATGTATATAA T <sub>12</sub>
413	449	224	ATATATA ATATAAGTAAATGAGAGACGAGGGTAGACTTGTGAT T <sub>10</sub>
410	449	125	ATATA ATATAAGTAAATGAGAGACGAGGGTAGACTTGTGATATA T <sub>09</sub>
413	452	545	ATATAT AACATATAAGTAAATAGAAGATGGAAGCGAATTGTCGAC T <sub>16</sub>
418	453	60	ATATATATAACAAAC AGACATATAAGTAAAGAGATGAAGGTAAATT T <sub>09</sub>
413	452	20	ATATAT AACATATAAGTAAATAGAAGATGGAAGCGAATTGTCGAC TCT <sub>05</sub>
418	467	451	ATAC TATAATAACAAACAAAATGTGTAAAGGTAGATAAGAAGTGAAGGTAAATT ATATTTT
437	469	238	A TACATAATAACAAATGAGATATATAAGGTGAAT CGAAAGTGAATAT <sub>12</sub>
427	460	218	ATATTAT AACAAACAAAACGTATAAGGTAAAGTGAAGGAAATGGG TGTAATAT <sub>12</sub>
443	474	37	A ATAATCACATAATAATGATAGAACGTATAAG ATAGATGAAAT T <sub>10</sub>
5'	3'	Reads	COIII gRNA Sequences cont.

461	497	66,677	ATACAT AATACCAATAGAACAGAACATTGTAGTCATGTGATA TTCAT <sub>14</sub>
461	497	5,941	ATACAT AATACCAATAGAACAGAACATTGTAGTCATGTGATA TTCAT <sub>13</sub>
456	499	936	ATAT AAAATACCAATAAGAACAGAACATTATAGTTGATGATGATAA AT <sub>12</sub>
462	497	371	ATACAT AATACCAATAGAACAGAACATTGTAGTCATGTGAT T <sub>11</sub>
461	497	139	ATACAT AATACCAATAGAACAGAACATTGTAGTCATGTGATA TTCATAT <sub>11</sub>
457	499	6,621	ATATT AAAATACCAATAGAACAGACTGTGATTATATGATGAATA T <sub>14</sub> *
458	499	3,602	ATATT AAAATACCAATAGAACAGACTGTGATTATATGATGAAT T <sub>14</sub>
454	499	295	ATAT AAAATACCAATAAGAACAGAACATTATAGTTACATGTGATAATA T <sub>10</sub>
462	499	185	ATATT AAAATACCAATAGAACAGACTGTGATTATATGAT T <sub>15</sub>
460	501	185	AAA AAAAAATACCAGTAGAACAGAACATAATCATGTGATAA TTTT
461	499	107	ATATT AAAATACCAATAGAACAGACTGTGATTATATGAT T <sub>14</sub>
483	522	288	AAATA ATCAACAAATTAAATGAATCTAAAGGTATCAGTGAAAA T <sub>14</sub>
491	539	145,031	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAATATT T <sub>10</sub>
488	539	4,441	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAATATTAA TTCT <sub>08</sub>
490	539	1,630	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAATATT T <sub>15</sub>
498	539	1,118	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAATTTAGAGAAT TTTT
495	539	1,052	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAAT T <sub>14</sub> *
491	539	849	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAATATT T <sub>10</sub>
487	539	847	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAATATTAA T <sub>14</sub>
486	539	653	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAATATTAG TTTT
504	535	599	ATAT ATAAAATGTATTGTGACGAGTTAATGGAT GTAGAAGAT <sub>12</sub>
491	539	482	ATA TAAATAAAATGTATTGTCAATGGATTAGATAAATTAGAGAATATT T <sub>11</sub>
502	539	448	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATT TAGAATAT <sub>12</sub>
491	539	353	ATA TAAATAAAATGTATTGTCAATGGATTAGATGAATTAGAGAATATT T <sub>09</sub>
504	535	1,604	ATAT ATAAAATGTATTGTGACGAGTTAATGGAT GTAGAAGAT <sub>15</sub>
528	565	188	ATATAT AAAATTAACAAGTGAATCACTAACAGATAGATAGATG ATAT <sub>12</sub> *
528	564	181	ATATT AAATTAACAGATAAGCCACTGACAAATAGATAGAGTG ATAT <sub>12</sub>
524	564	77	ATATAT AAATTAACAAATAGACTACTAAATAAGTGAAGATGTATT AATTATATATTTT*
525	563	3,592	ATATAT AATTAACAAGTAGATCACTGACAAATAGATGAGATGTAT AAT <sub>13</sub> *
528	565	640	ATATAT AAAATTAACAGATAGATCATTACGAGTAGATAAGTG ATAT <sub>11</sub>
523	563	336	ATATAT AATTAACAAGTAGATCACTGACAAATAGATGAGATGTATT T <sub>08</sub>
526	564	420	ATATT AAATTAACAAATAAACTATTAAATGGATGAGTGAGATGT ATTAT <sub>15</sub>
527	563	144	ATATAT AATTAACAAGTAGATCACTGACAAATAGATGAGATGT T <sub>16</sub> *
548	592	99,540	ATATAT AAACCTAAATCAAGAACATAGAACAGAGAGATTGTGAGTAAATT T <sub>12</sub>
555	594	27,720	ACAT AAAACCTAAACTGAGAACATAGAACAGAACATTAGTGA TTAAAT <sub>12</sub>
547	592	1,442	ATATAT AAACCTAAATCAAGAACATAGAACAGAGAGATTGTGAGTAAATT T <sub>11</sub>
551	592	599	ATATAT AAACCTAAATCAAGAACATAGAACAGAGAGATTGTGAGTAA T <sub>13</sub>
554	592	504	ATATAT AAACCTAAATCAAGAACATAGAACAGAGAGATTGTGAG AAAAT <sub>13</sub>
558	592	309	ATATAT AAACCTAAATCAAGAACATAGAACAGAGAGATTAG GGAGTAAAT <sub>14</sub>
548	592	238	ATATAT AAACCTAAATCAAGAACATAGAACAGAGAGATTAGCGAGTAAATT T <sub>09</sub>
550	592	176	ATATAT AAACCTAAATCAAGAACATAGAACAGAGAGATTGTGAGTAA AT <sub>11</sub>
545	592	150	ATATAT AAACCTAAATCAAGAACATAGAACAGAGAGATTGTGAGTAAATT T <sub>10</sub>
558	594	147	ACAT AAAACCTAAACTGAGAACATAGAACAGAGAACATTAG GGATTAAT <sub>11</sub>
558	593	128	ACATT AAAACCTAAACTGAGAACATAGAACAGAGAGATTAA GGATTAAT <sub>13</sub>
548	592	114	ATATAT AAACCTAAATCAAGAACATAAAACAGAGAGATTGTGAGTAAATT T <sub>10</sub>
551	593	84,835	ACAAT AAAACCTAAACCGAGAACATAGAGCAGAGAACATTAGATAA T <sub>13</sub>
555	593	6,897	ACATT AAAACCTAAACTGAGAACATAGAGCAGAGAACATTAGATAA TTAAAT <sub>14</sub>
552	593	4,207	ACAAT AAAACCTAAACCGAGAACATAGAGCAGAGAACATTAGATAA T <sub>13</sub>
556	593	1,803	ATATAT AAAACCTAAATCAGAGACGCAGAACATTAGAGAGATTGTGAGATA TAT <sub>10</sub> *
553	593	1,137	ACAAT AAAACCTAAACCGAGAACATAGAGCAGAGAACATTAGATAA T <sub>12</sub>
551	596	722	A AAAAACCTAAACCGAGAACATAGAGCAGAGAACATTAGATAA T <sub>15</sub>
557	593	486	ACATT AAAACCTAAACTGAGAACATAGAACAGAGAGATTAA T <sub>12</sub>
551	603	456	ACAAAAAAACCTAAACCGAGAACATAGAGCAGAGAACATTAGATAA T <sub>05</sub>
555	593	206	ATAT AAAACCTAAACCAAAGATATGAGAACAGAGAGATTGTGA TATGT <sub>13</sub>

5'	3'	Reads	COIII gRNA Sequences cont.
585	629	85,916	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGAGATTAA T <sub>12</sub>

587	629	9,106	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGAGATT T <sub>11</sub>
586	629	1,991	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGAGATT T <sub>12</sub>
585	631	1,447	ATAT ATAAAACCTCAATCATAGTATAAGATGACGACAATGAGAAGATTAA T <sub>13</sub>
585	629	276	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGAAGATTAA T <sub>12</sub>
592	629	269	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGA TTTAAT <sub>14</sub>
591	629	263	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGA TTTAAT <sub>07</sub>
588	629	221	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGAGATT AAT <sub>13</sub>
594	629	150	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGAGATT TAGATTAA T <sub>06</sub>
580	622	145	ATATAT ATCATAATACAAGGCAATGACGACGAGAAGATTAGATTAA T <sub>11</sub>
590	629	131	ATATA GAACTCAATCATAATATGAAGCAATAACAATGAAGAGA ATTAATTAA T <sub>06</sub>
603	634	17,313	AT ATAACAAACTTAATCGTAATATGAAACACGA GAATGAGAAAAT <sub>13</sub>
585	628	8,519	ATATAT AACTCAATCATAATACGAGATGATAATGACGAAGAGATTAA T <sub>13</sub>
580	624	3,564	ATATA TAATCATAATACAGAACGATGGCAGTGAGAGATTAGTTAA T <sub>12</sub>
587	630	2,850	ATATATAT CAAACTCAATTGTAGTAGCAGAGACAATAATGATGAGAAGATT T <sub>10</sub>
587	628	628	ATATAT AACTCAATCATAATACGAGATGATAATGACGAAGAGATT T <sub>08</sub>
585	630	582	ATATATAT CAAACTCAATTGTAGTAGCAGAGACAATAATGATGAGAAGATTAA T <sub>11</sub>
586	628	231	ATATAT AACTCAATCATAATACGAGATGATAATGACGAAGAGATTAA T <sub>11</sub>
587	631	135	ATATAT ACAAACTCAGTATAGTATAAGACAGTGATAATGAGAGAATT T <sub>12</sub>
591	630	112	ATATATAT CAAACTCAATTGTAGTAGCAGAGACAATAATGATGAGAAG T <sub>13</sub>
604	647	10,315	ATACAAAAAC AAAAACCAAACGATGAACTTGTAGTATAAGATAATA T <sub>13</sub>
605	647	7,606	ATACAAAAAC AAAAACCAAACGATGAACTTGTAGTATAAGATAAT T <sub>13</sub>
604	643	250	ATACAAAAACAAAAT ACCGAGCGACAGATTGATTGTAGTATAAGATAATA T <sub>11</sub>
607	647	637	ATACAAAAAC AAAAACCAAACGATGAACTTGTAGTATAAGATAATA T <sub>21</sub>
604	643	173	ATACAAAAACAAAAT ACCGAGCGACAAGTTGATTATAGTATAAGATAATA T <sub>13</sub> *
605	643	150	ATACAAAAACAAAAT ACCGAGCGACAAGTTGATTATAGTATAAGATAAT T <sub>15</sub> *
635	669	482	ATATACAATGCAAACCTTA CATGACTGGTTTATAGAGATGAGAGATTAA T <sub>14</sub>
635	669	268	ATATACAATGCAAACCTTA CATGACTGGTTTATAGAGATGAGAGATTAA T <sub>15</sub>
635	676	60	ATATACAATGT ACTCTCATATAATTGTTCATAGAGATAGAAGATTAA T <sub>15</sub>
637	669	42	ATATACAATGCAAACCTTA CATGACTGGTTTATAGAGATGAGAGATT T <sub>08</sub> *
633	679	25	ATATA CAAACTCTTATAACTGGTTTACGAGAATGAGAAATTAAAT T <sub>15</sub>
659	691	649	ATAT TAAATAACAATGCGAATTTCATAGTTGGTT CATAGATACAAT <sub>12</sub>
653	682	39	AT ATGTGAACTCTTATAACTGGTTTGTAG TGATGATAGAT <sub>15</sub>
669	717	758	ATAAT AGAACACCACAGCTTAATGTAAGTAGATGGCAGTGAAATT T <sub>10</sub>
669	722	374	ATAT ATACAAAAACACCGTAATTGATGAGTAGATAGTAGTGAAATT T <sub>07</sub>
669	706	354	ATATAGAACCAAAACAG TGCAATTAGTGTAGATGATAGTGAAATT T <sub>07</sub> *
669	715	160	ATATAGAACCAAT AACATCGCAGCTTAGTGTAGTAATAGTGAAATT T <sub>07</sub>
684	722	1,329	ATATAT AACCAAAACACTCGCAGTTGATGATAAAGTGAC TGTGAAAT <sub>11</sub>
689	726	1,181	ATAT ATAGAACCAAGAACACCATAGTTGATGATGATAG TGATAGTGAAAT <sub>14</sub> *
669	715	911	ATATAGAACCAT AACACCATGATTGATGAGTAGTAATGATGATGAAATT T <sub>09</sub>
669	722	738	ATAT AACCAAAACACTGTAATTGATGAGTAGATAGTAGTGAAATT T <sub>08</sub> *
695	730	129	ATAT TAAAATAGAAGACTAAAGACACTGTAATTGATGAGTAGAATGATGAAATT T <sub>10</sub> *
675	715	107	ATATAGAACCAT AACACCATGATTGATGAGTAGAAATGATGATGAAAT AT <sub>08</sub>
677	715	102	ATATAGAACCAT AACACCATGATTGATGAGTAGAAATGATGATGAA T <sub>11</sub>
675	717	90	ATAAATT AAAACACCATAATTGATGATGATAAGTAATGATGAAAT AT <sub>17</sub>

5'	3'	Reads	COIII gRNA Sequences cont.
706	753	31,331	ATATAT AAATGTAATAGATCTGATGAAAGTGAGGTAGAATTGAGAATATT T <sub>10</sub>
707	753	8,037	ATATAT AAATGTAATAGATCTGATGAAAGTGAGGTAGAATTGAGAATAT AT <sub>14</sub> *

			COIII gRNA Sequences cont.
5'	3'	Reads	COIII gRNA Sequences cont.
880	918	1,822	ATATAA ATCAGAATAACAGATCGCAATAGAGAGAAATTAGTTAA TAT <sub>14</sub>
882	927	600	AA ATATAAAACATCAAGATAATTGGATTGTGATAGAGAAAATT T <sub>11</sub>
880	929	222	ATATATAAAACATCAGAATAGACAAATCGTAATAGAGAGAAATTAGTTAA T <sub>14</sub>

880	921	209	ATAT AATATCAAAATAAACAGATCGTAGTAAAAGAAGTTAGATTAA T <sub>12</sub>
882	929	159	ATATATAAAACATCAGAATAGACAAATCGTAATAGAGAACTTAAGTT T <sub>13*</sub>
881	929	28	ATATATAAAACATCAGAATAGACAAATCGTAATAGAGAACTTAAGTTA TTTT
907	947	787	ATATAA TATACACACAGATAACATAATACGTAGAATGTTAAGATAAGT T <sub>16</sub>
909	950	126	ATATA ATTACACACACAGATACTGTGATATATAGAATGTTAAGGTAA TATAAT <sub>10</sub>
913	946	558	ATAC ATACACACAAATATATAACATATAGAGCATTGAG TTAGATAAT <sub>15</sub>
905	944	3,472	ATATAA ACACACAAATATATGGCATATAGAGCATTGAAGTAGATAA T <sub>13*</sub>
905	944	1,321	ATATAA ACACACAAATATATGACATATAGAGCATTGAAGTAGATAA T <sub>15</sub>
920	952	162	ATAG AAAATTACACACATGAATACTACATAGTACATAGAA GATTGATATAT <sub>13</sub>
935	977	2,920	ATATA AATCAACAACGTGAAAAGATATCAATGAGATTGTACATGTAAAT T <sub>15</sub>
940	977	522	ATATA AATCAACAACTAAGAAGACACTGTGATAGAGTTATATGTG ATTAAT <sub>14*</sub>
935	977	354	ATATA AATCAACAACGTGAAAAGATATCAATGAGATTGTACATGTAAAT T <sub>08</sub>
939	977	28	ATATA AATCAACAACGTGAAAAGATATCAATGAGATTGTACATGT T <sub>09</sub>
942	983	1,160	ATAT AACTAATCAACAGCTAAGAGAACGTCAATGAGATTATGTG ATTAAT <sub>14</sub>
951	981	127	AT ACTAATCAACAACTAGAAGAATATCAGTGA TATACATGTAAAT <sub>10</sub>
951	981	25	AT ACTAATCGACAACTAGAGGGACATCAGTGA TTTATACGTAT <sub>15</sub>
963	1003	111	ATA TACAAACTACCAATAAGTTAACTGATCGGTAAATTAAGG TTATAT <sub>15</sub>
965	1003	68	ATATA TACAAACTACCGATATAAGTTAACTGATTGATAATTAA TGTTCT <sub>14</sub>

### C) C-Rich Region 3

5'	3'	Reads	CR3 gRNA Sequences
34	62	18	ATATGT ACAACAAAACCGAGCAATCAGATAT AGAGTGAAAT <sub>09</sub>
34	62	3	ATATAT ACAACAAAACTGAACAATCAAATGT AGTGTGAT <sub>09</sub>
36	64	1	ATCGCAAGTCGT GGACACAAAACTGAACAATCAAAT T <sub>11</sub>
34	62	1	ATATAT ACGACAAAATGAACAATCAAATGT AGTGTGAT <sub>08*</sub>
41	88	140,541	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAGATGAT AT <sub>14</sub>
40	88	34,162	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAGATGATT T <sub>12</sub>
48	87	3,434	ATATAATT AAATGTACAGACAAATGATAGAGAGACGATGAGATTAAGT TATAT <sub>12</sub>
51	86	3,016	AAAATT AATGTACAAAATAACGATAGAGAGACAGTGAAT TGAT <sub>13</sub>
51	88	2,553	ATAT AAAATGTACAGACGGAGCAGTGAAGAGAACAGTGAAT TACAT <sub>11</sub>
47	88	1,270	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAGATG T <sub>11</sub>
48	88	902	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAGAT T <sub>14</sub>
41	88	693	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTAAAATTAGATGAT AT <sub>12</sub>
52	88	648	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATT T <sub>09</sub>
51	88	632	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATT TATGATAATTATTT
46	88	468	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAGATGA AATAT <sub>14</sub>
40	77	213	ATATAT AAAATGTACATACGAAACGATAAAAGGGCAGTGAATTAGATAATT TT
41	88	212	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAGATGAT ATCT <sub>07</sub>
50	88	182	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAG TGTAAAT <sub>11</sub>
41	88	181	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAGATGAT ATCT <sub>16</sub>
49	88	177	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAATTAGA AGATATAT <sub>11</sub>
55	88	170	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGA T <sub>19</sub>
41	88	155	ATATAT AAAATGTACAAACGAACAATGAGAGAACAGTGAATTAGATGAT AT <sub>13</sub>
40	88	151	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTAAAATTAGATGATT T <sub>12</sub>
41	88	135	ATATAT AAAATGTACAAATGGACAATGAGAGAACAGTGAATTAGATGAT AATAT <sub>07</sub>
56	88	133	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGA T <sub>12</sub>
41	88	106	ATATAT AAAATGTACAAACAGACAATGAGAGAACAGTGAATTAGATGAT AATAT <sub>08</sub>
40	89	598	ATAT AGAAATGTACAAACGGACAATAAGGGAACAGTGAATTAGATGATT TCT <sub>07</sub>
41	89	526	ATAT AGAAATGTACAAACGGACAATAAGGGAACAGTGAATTAGATGAT AAT <sub>10*</sub>
40	77	437	ATATATAAAATGTACAT ACGAACGATAAAAGGGCAGTGAATTAGATAATT T <sub>09*</sub>
47	83	285	ATATAA GTACAAACAAACAGTGAAGAGATAACGAGACTGAGTA TATATTT*
50	88	255	ATAT AAAATGTACAAAATAGCAGTAGAGAGCAGTGAATT T <sub>09</sub>
47	89	237	ATAT AGAAATGTACAAACGGACAATAAGGGAACAGTGAATTAGATG TTTAAT <sub>09</sub>
51	88	152	ATAT AAAATGTACAGACAAGCAGTAGAGAGACAGTGAATTAGATTA TACAGTTT
39	77	112	ATATATAAAATGTACAT ACGAACGATAAAAGGGCAGTGAATTAGATAATT T <sub>12</sub>
51	86	70	AAAATC AATGTACAAAATAAGCGATAAAGAGAACAGTGA T <sub>15</sub>
48	89	53	ATAT AGAAATGTACAAACGGACAATAAGGGAACAGTGAATTAGAT T <sub>12</sub>
47	83	31	ATATAA GTACAAACAGACAGTGAAGAGATAACGAGACTGAGTA TATTT <sub>17</sub>
78	118	573	ATATAT AATCACAAACAAATAGAAAATGAGAGAGGTATGA TACTAT <sub>15</sub>
93	123	1,321	ATATAT AAACAAATCACGAAACGAGTAGAAAAT TGGAAAGATGTAT <sub>12</sub>
78	118	685	ATATAT AATCACAAACAAATAGAAAATGAAAGAGGTATGA TACTATTTT
89	123	251	ATATT AAACAAATCACAAATGAGTAGAGAACGAGA TTGATGTATAT <sub>16</sub>
86	121	243	ATAC ATAAATCACAAACGAATAAGGGCAGAACAGAG TTGTATGATAT <sub>05</sub>
76	124	120	ATATAT AAAACAAATCATAAACGAATAAAGAGTGAAGGAAGGTATATAAT TTT
105	140	15,770	A TAGATAACAAACATAAGAGCAAGTCACGAG GTAGATATTGATAT <sub>14</sub>
105	140	692	A TAGATAACAAACATAAGAGCAAGTCACGAG ATAGATATTGAT <sub>12</sub>
98	142	442	AT ATTAAAATAACAAACATAGAAAATAGATCACAGATAGATGA TTGATGAAT <sub>06</sub>
122	166	27,885	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAAGA TAT <sub>13</sub>
124	166	212	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAA T <sub>12</sub>
125	166	139	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATA T <sub>17</sub>
127	166	114	ATAT ACAAAAATCCAATGAAAAATAAGACTGAGTAGTGTGGATG CAAT <sub>15</sub>
123	162	309	ATATAT AAATCCAATAAGAAATGAAAGCTAGATAGTAGTGTAGTATAAG T <sub>16</sub>
124	166	421	ATATAT ACAAAAATCCAATGAAAAATAAGACTGAGTAGTGTGGATGTA TTATTT <sub>14</sub>
129	167	370	ATAT AACAAAAATCTAATAGAGAACAGAAACTAAGTATGAG ATAT <sub>09</sub>
123	161	348	ATATTAT AATCCAATAGAACAGAGACTAAATGATAGATGTAAA T <sub>09*</sub>
132	168	110	ATATAT AAACAAAATTCAATAGAACAGAGACTGAGTAAT TGATATGT <sub>05</sub>
125	167	53	ATATAT AAACAAAATCCGATAAAAGATGGAAACTAAGTGTAGATATA T <sub>13</sub>
154	196	1,492	ATAT ATACAACAAATAACTCGTATTAAGTAGAGAGATGAAGATTTAAT T <sub>13</sub>
156	199	136	ATAT TAAACACAACGATAGATCTATATTAAGTAGAGATAGAAATT T <sub>16</sub>
162	200	2,176	ATAT ATAAACACAACAGTAGACTCGTATTAGATAGAGATAGA TATCAT <sub>15</sub>
156	200	114	ATAT ATAAACACGACAATAGATTATTAAGATAGAGATAGATTTA T <sub>12</sub>
157	200	98	ATAT ATAAACACGACAATAGATTATTAAGATAGAGATAGATTTA T <sub>07*</sub>

5'	3'	Reads	CR3 gRNA Sequences cont.
190	230	2	ATATACA ACATATCAAGTGGTAAGATAAGAGAAGAAAGTAGATGTAAT T <sub>06</sub>
190	230	1,243	ATATACA ACATATCAAGTGATAAGATAAGAGAAGAAAGTAGATGTAAT T <sub>12</sub>
192	232	279	ATAGATA ACACATATCAGATGATAAGGTAAAGAGAGAGAATAGATATA T <sub>13*</sub>
226	277	19	ATAT AATAACAATATAAACGAACCTGGATGATGATAGTACTTTGATATATA AT <sub>12</sub>
231	265	2	ATATAACAATAC AAATGAGCTAGATAATGGATGATAGTTGATAT T <sub>12</sub>
237	265	1	TATAAAATAACAATAC AAGTGAATTAGATAATGGATGATG TTTCAAT <sub>15</sub>
241	279	936	AT AGAATAACAATATAAACGAACTAGATGATGGATAGTA TCT <sub>17</sub>
234	265	528	ATATAACAATAC AAATGAACTGAGTAATGGATAGTAGTTGA CAAT <sub>09</sub>
293	308	2	ATATAT AAATTATTCGATACT GAGATTAGTGATAGTAAAGTGATTAAT <sub>13</sub>
268	312	125	AT ATAAAAATTATTCGATGCTTAAGTGAGTTAGAGGTAATGATA AAT <sub>07</sub>

D) C-Rich Region 4

5'	3'	Reads	CR4 gRNA Sequences
25	64	121	ATAAT AAAAATGCACAACCTAGAATTGAAAGTAAAGTGATA TATAT <sub>14</sub>
25	64	596	ATAAT AAAAATGCACAACCTAGAATTGAAAGTAAAGTGATA TATAT <sub>14*</sub>
25	64	308	ATATT AAAAATGCACAACCTAGAATTGAAATAAGTGATGGTA TATAT <sub>13</sub>
25	62	83	ATATAATT AAATGCACAGCCAAAGTTAAGGTAGAATAGTGATA TAT <sub>14</sub>
48	103	296	ATA TATATAAAACACAGACATACTAAGTAAGAGAAAAGAGAGGTGTATGATT T <sub>12</sub>
48	103	175	*ATA TACATAAAACACAGACATACTAAGTAAGAAAAAGAGAGATGTATGATT T <sub>12*</sub>
65	98	108	ATATAT TAAAACACAAATACATCAGATAGAAGAGA TTGAGTGTATAAT <sub>17</sub>
51	103	18	ATA TACATAAAACACAGACATACTAAGTAAGAAAAAGAGAGATGTATG T <sub>21</sub>
59	103	11	ATA TACATAAAACACAGACATACTAAGTAAGAAAAAGAGAGAT T <sub>15</sub>
52	103	10	ATTATA TACATAAAACACAGACATACTAAGTAAGAAAAAGAGAGATGTAT T <sub>15</sub>
87	134	7,793	ATATAT AAACAACAATAGAGTATATCATAGACTGTATATGAAGCATAAT T <sub>10</sub>
89	134	426	ATATATAT AAACAACAATAGAGTATATCATAGACTGTATATGAAGCATAA CT <sub>08</sub>
88	134	208	T ATATAT AAACAACAATAGAGTATATCATAGACTGTATATGAAGCATAAA AT <sub>13</sub>
93	139	200	AAACAAAACAACAGTCAAATATACCGTAGATTGTATGTGAAAT TATATAT <sub>06</sub>
90	134	143	ATATATAT AAACAACAATAGAGTATATCATAGACTGTATATGAAGCATA T <sub>13</sub>
93	142	2,052	A AAAAACAAAACAACAGTCAAATATACGTAGATTGTATGTGAAAT TAT <sub>05</sub>
90	138	584	ATAT AACAAAACAACGATAAGATGTATCATGAGCTGTATATGAGATATA T <sub>25</sub>
91	138	139	ATAT AACAAAACAACGATAAGATGTATCATGAGCTGTATATGAGATATA T <sub>13</sub>
127	171	6,971	ATAT ATATACTCACACAAATAGATGACAGAGATAGAAAGTAAGATGATA TAT <sub>15</sub>
128	171	204	ATAT ATATACTCACACAAATAGATGACAGAGATAGAAAGTAAGATGAT T <sub>13</sub>
154	192	3,603	ATATTA AACTATAACAAGGCAGATAGAACGTACCTATATAGATAA T <sub>14</sub>
174	200	6	AT AAATAAACAACTATAATAAGACAAGTG CGATGACT <sub>10</sub>
166	196	5	ATATA AAACAACCTATAACAAAATAGATAGAATGTGC GTAT <sub>06</sub>
171	200	3	AAATAAACAACTATAATGAAGTGTAGAAG GTACGT <sub>12</sub>
169	196	1	ATATA AAACAACCTATAACAAAATAGATAGAATG GCGTAT <sub>06</sub>
171	196	1	A AAACAACCTATAATGAAGCAGATAGAA GTACGTATAT <sub>17</sub>
186	232	31	ATAAAAAAATCACAGCCTAAAATGACGAGAGAAAAGTAAATGGTTATA TAGAT <sub>05</sub>
186	225	22	ATATAT ATCACAACTCTAAAGATAACGAAAAAGAGCAGATAATTGTA TATTTT
186	228	7	ATATAT AAAATCACAACTTAGAATGACGAAAGAGAAATAGATAGTTGTA TATTGT <sub>22</sub>
186	228	4	ATATAT AAAATCACAACTTAGAATGACGAAAGAGAAATAGATAGTT TTAAAAAAAT <sub>16</sub>
186	232	166	ATAAAAAAATCACAGCCTAAAATGACGAGAGAAAAGTAAATGGTTATA T <sub>16</sub>
186	227	152	ATAT AAATCACAACTCTGAAATAGCAGAGAAGAGTAAATGATTATA T <sub>13</sub>
189	232	25	ATAAAAAAATCACAGCCTAAAATGACGAGAGAAAAGTAAATGGTT TTT*
189	227	12	ATATAT AAATCACAACTGAAATAGCAGAGAAGAGTAAATGATT T <sub>09*</sub>
213	261	14,358	ATAT ATAAACTATACAAATTGAAAGCACTGATAGAAGGTTGTGATTTAA T <sub>12</sub>
213	261	264	ATAT ATAAACTATACAGTCAGACTGTGAGAGATCGTGATTTAA T <sub>13</sub>
216	261	114	ATAT ATAAACTATACAAATTGAAAGCACTGATAGAAGGTTGTGATTT TTT
222	261	112	ATAT ATAAACTATACAAATTGAAAGCACTGATAGAAGGTTG ATTTAATAT <sub>08</sub>
213	258	5,062	ATATAT AACTATACAGTCGAGACATCAATGAGAGATTGACTTAA T <sub>13</sub>
213	259	1,370	CCCTCGGCCGCAGCGATATA AAACATATACAATTAGAGCATCAGTAGAAGATTGTGACTTAA T <sub>14</sub>
213	259	1,074	ATATA AAACATATACAATTAGAGCGTCAGTAGAAGATTGTGACTTAA T <sub>15*</sub>
213	258	490	ATATAT AACTATACAAATTGAGATATCAGTGAAGAATTGTGATTTAA T <sub>11</sub>
215	258	340	ATATAT AACTATACAGTCGAGACATCAATGAGAGATTGACTT T <sub>14</sub>
223	258	183	ATATAT AACTATACAGTCGAGACATCAATGAGAGATT T <sub>05</sub>
251	300	23	AAAAAAAAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATTATAT T <sub>11</sub>
242	284	9	AAAAAAA AAAATTATAACGTCTAAAGAAGATAGACTATATGATTTAA T <sub>09</sub>
251	295	77	ATAT AGAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATTATAT TTTT
250	295	64	AT AGAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATTATATA T <sub>06</sub>
253	295	25	ATAT AGAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATTAT T <sub>16</sub>
253	298	22	AAAAAAAAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATTAT TTT
244	295	20	AT AGAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATTATAGTT T <sub>12</sub>
250	297	17	AAAAAAAAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATTATATA T <sub>12</sub>
251	301	16	AAAAAAAAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATTATATA T <sub>05</sub>
255	295	11	ATAT AGAAAATAAACAGAATTGTGACGTCTAAAGAAGATAGATT TTTT
280	320	643	AAA AAAAACACAAGGCAGATAGAGAAAAGAGATAATAAGAT GT <sub>10</sub>
279	320	93	AAAAA AAAAACACAAGGCAGATAGAGAAAAGAGATAATAAGATT T <sub>05</sub>
5'	3'	Reads	CR4 gRNA Sequences cont.
307	351	725	ATAAAAAT ACCAAACAGATCAGATAAAAGCAGTGTATAGAGAATATAAAAT T <sub>11</sub>

311	353	675	ATAT AAACCAAACAAGCTGAATAAGAACAGTGATATAGAAGATATA TAT <sub>14</sub>
306	351	259	AAAAT ACCAAACAGATCAGATAAAGGCAGTGATATAGAGAATATAAATA T <sub>12</sub>
309	351	214	ATAAAAT ACCAAACAGATCAGATAAAGGCAGTGATATAGAGAATATAA T <sub>12</sub>
324	354	140	AT AAAACCAGACAAACTGAATGAAGACAGTAACGTAGAAGATAT T <sub>10</sub>
307	352	485	TTGTTGGTTGATTAAT AACCAAACAAGTCGAGTAGAGACAGTGATATAAAGGTATAAAAT T <sub>09*</sub>
310	353	374	ATAT AAACCAAACAGACCAAGTGAAGATGCCAGTATAAGAGATATGA TATAT <sub>11</sub>
309	352	134	AAAT AACCAAACAAGTCGAGTAGAGACAGTGATATAAAGGTATAAA T <sub>16*</sub>
310	352	42	AAAT AACCAAACAAGTCGAGTAGAGACAGTGATATAAAGGTATAA T <sub>12</sub>
312	352	38	AAAT AACCAAACAAGTCGAGTAGAGACAGTGATATAAAGGTAT T <sub>13</sub>
343	388	863	ATAT ATAACACAAAAACATAACAGAGGTATAGAGAGAAATTGAATGA T <sub>15</sub>
340	390	2,308	ATAT AAATAACACAAAATACGACGAGAAATATAAGAGAGATTGAGTAAATT T <sub>05*</sub>
343	388	255	ATAT ATAACACAAAACATAACAGAGGTATAGAGAGAAATTGAATGA T <sub>15*</sub>
374	417	411	AAAAA AAAAAACACATAGAAAGTGAATCAGAGAATGACATAAGATATA T <sub>11</sub>
377	417	447	A AAAAAACACATAGAAAATAAGTCAGAGACTATATGAGAT TGTTATAAT <sub>05</sub>
374	415	423	ATATAT AAAACACATAAGAGATGAATCAAGAGGTATATAAGATATA T <sub>17</sub>
404	457	4,211	AAAAAAAAAAACAAAGACAAGAAATCACTCAGAATAGAAGATGGTATAA T <sub>14</sub>
405	456	620	AAAAAAAAAAACAAAGACAAGAAATCACTCAGAATAGAAGATGGTATA T <sub>13</sub>
406	457	392	AAAAAAAAAAACAAAGACAAGAAATCACTCAGAATAGAAGATGGTAT T <sub>12</sub>
405	457	406	AAAAAAAAAAACAAAGACAAGAAATCTGCTAAGATAGAGAATGATATA T <sub>15</sub>
409	457	185	AAAAAAAAAAACAAAGACAAGAAATCACTCAGAATAGAAGATGG GATAAAT <sub>08</sub>
404	458	227	AA AAAAAAAAAAAACAAAGACAAGAAATCACTCAGAATAGAAGATGATATAA T <sub>08*</sub>
405	456	186	AAAAAAAAAAACAAAACAAGAAGACTATCTGAGGTAGAAAATGATATA T <sub>14</sub>
411	443	97	ATATAAACAT AAACAAAGAGACCATCGAAATAGAGAAT T <sub>15</sub>
405	457	34	AAAAAAAAAAACAAAGACAAGAAATCACTCAGAATAGAAGATGATATA T <sub>09</sub>
406	458	25	A AAAAAAAAAAAACAAAGACAAGAAATCACTCAGAATAGAAGATGATAT T <sub>09*</sub>
442	489	3,197	ATAACACCCACAGATAGAAATGAACGTTAATGAGAGAGAAAGATGAAA T <sub>13</sub>
440	489	238	ATAACACCCACAGATAGAAATGAACGTTAATGAGAGAGAAAGATGAAAGT T <sub>11</sub>
442	487	730	ATAA AACAAACCACAAGTAGAAATAGACATAATGAGAGAGAAAGATAAAA T <sub>12</sub>
442	486	246	ATAT ATAACCCACAATAGAGACAATGTAAGTAAAGAGAGAAAGTAAA T <sub>19</sub>
443	487	208	ATAA AACAAACCACAAGTAGAAATAGACATAATGAGAGAGAAAGATAAA T <sub>06*</sub>
446	487	149	ATAA AACAAACCACAAGTAGAAATAGACATAATGAGAGAGAAAGAT T <sub>11</sub>
453	497	1,153	AT ACAAAATAACAACAAATCATAAGTAAGAATAGATGTAGATGAGAAA T <sub>11</sub>
453	494	282	ATATAT AAATAACACAAATCACAGATGAAGGTAGATATAAGTGTAGATGAGAAA T <sub>14</sub>
453	497	2,511	AT ACAAAATAACAATGATCACAAATAAGAGTGAAGTGTAGATAGAGAA T <sub>15*</sub>
453	494	306	ATATAT AAATAACACAAATCACAGATGAAGGTAGATATAAGTGTAGATGAGAAA T <sub>13*</sub>
475	519	766	ATATAT ATAACACAAATAATCAGATTAGCAGAGTAATGATAGTTATAAT T <sub>13</sub>
480	524	390	ATACAAATAACAACAAATAGCCAAGTTAATAGAGTGTGATGATTTA AT <sub>14</sub>
478	524	1,913	ATACAAATAACAACAAATGATCAGACTAATAGAGTAATAGTGTATTATA T <sub>14</sub>
481	524	105	ATACAAATAACAACAAATGATCAGACTAATAGAGTAATAGTGTATT T <sub>12</sub>
479	524	101	ATACAAATAACAACAAATGATCAGACTAATAGAGTAATAGTGTATTAT T <sub>15</sub>
504	554	146	AA AAAAAACGCATATAAGTGGCTTACATAGATAATGATGATAATT T <sub>03</sub>
503	547	53	AAAAAAAAAA GCATATAAAATAGATCTATATATGAGTGTAGTGACAAATT T <sub>15</sub>
504	547	50	AAAAAAAAAA GCATATAAAATAGATCTATATATGAGTGTAGTGACAAATT T <sub>11</sub>
507	554	98	AAAAA AAAAACGCATATGAGTAGATCTGTACATAGATAGTGTGATA T <sub>12</sub>
507	556	84	ATAGAAAACGCATATGAGTAGATCTGTACATAGATAGTGTGATA CT <sub>06</sub>
508	554	37	CAACTCGACTGCGTGAA AAAAACGCATATGAGTAGATCTGTACATAGATAGTGTGAT TTT
542	575	511	AAATTTATAAAACCAAT CCGTAATTATCTGAGATGAGAAGTGTATA AT <sub>12</sub>
542	584	997	ATATTTATAAAACCC ATACTGTGATTATCTAAGGTAGAAAGTGTGTA AT <sub>21</sub>
542	596	352	ATA TTTATAAAACCAATATCATAGTTATCTGAGGTAGAGAGTGTATA ATTTAATGTAT T <sub>18</sub>
537	568	62	CATACC TAATTATCTGAAGTAGAAGATGTATGTAAAT T <sub>14</sub>

E) Cytochrome b

5'	3'	Reads	CYb gRNA Sequences
32	59	327	AAATAATAGGGATTATGATGAGATATG CTGTGGATAT <sub>14</sub>
32	60	255	AAAATAATAGGGATTATGATGAGATATG CTGTGGATAT <sub>12</sub>
32	61	190	AAAAATAATAGGGATTATGATGAGATATG CTGTGGATAT <sub>13</sub>
32	64	176	AA AAAAAAATAATAGGGATTATGATGAGATATG CTGTGGATAGT <sub>09</sub>
32	62	65	GT AGAAAATAATAGGGATTATGATGAGATATG CTGTGGATAGT <sub>12</sub>
32	64	35	AA AAAAAAATAATAGGAGTTATGATGGAATATG CTGTGGATATTTT*
32	64	17	AA AAAAAAATAATAGGAGTTATGGTGAGATATG CTGAGAGTAT <sub>05</sub>
53	91	18,713	AAAAAA AAAAGACAGTGTGAATTCCTGAGTAATAAGGAAATAAT T <sub>11</sub>
51	91	10,649	AAAAAA AAAAGACAATATAGATTCCTGGGTGATAAAAGGGATAATAA CT <sub>11</sub>
52	91	908	ATAA GAAAGACAATATAGTTCTGGTAATGGAGAGAATAATA T <sub>16</sub>
54	91	339	AAAAAA AAAAGACAATGTAGATTCCTGAGTAATGGGAGGATAA CTATTTATTTT*
54	91	8,406	AAAA AAAAGACAATGTAGATTCCTGAGTAATGGGAGGATAA CTAT <sub>05</sub> *
54	91	2,265	AAAAAA AAAAGACAATGTAGATTCCTGAGTAATAGGGAGGATAA CTAT <sub>16</sub>
53	91	204	AAAAAA AAAAGACAATGTAGATTCCTGAGTAATGGGAGGATAAT T <sub>07</sub>
56	91	95	AAAAAA AAAAGACAATGTAGATTCCTGAGTAATGGGAGGAT T <sub>05</sub> *

F) Maxicircle Unidentified Reading Frame II

5'	3'	Reads	Murf II gRNA Sequences
30	79	2,605	ATAG AAAGCACAAAATAAAATTAAATTAGAGTAATTGGATGTTAAAATT T <sub>11</sub>
30	79	125	ATAG AAAGCACAAAATAAAATTAAATTAGAGTAATTGAATGTTAAAATT T <sub>08</sub>
34	79	17	ATAG AAAGCACAAAATAAAATTAAATTAGAGTAATTGAATGTTAA CAT <sub>12</sub>
33	79	15	ATAG AAAGCACAAAATAAAATTAAATTAGAGTAATTGAATGTTAAA T <sub>09</sub>

### G) NADH Dehydrogenase 3

5'	3'	Reads	ND3 gRNA Sequences
33	76	313	ATATATT ATAAACCACATGATATCGAAAATGGGTAGAAATGATGATA T <sub>12</sub>
31	79	70	ATAT ATAATAAACACAGTATCAGAGACAGATATAGAAGTGATGATAGT T <sub>13</sub>
30	73	625	ATAT TAAACCACAATATCAGAAATAAGTGAGAAATAGTGATAATA T <sub>12</sub>
31	79	511	ATAT ATAATAAACACAGTATCAGAGACAGATATAGAAGTGATGATAGT T <sub>09</sub>
42	79	397	ATAT ATAATAATCACAGTATCAGAGACAGATATAGAA TGATGATAGT <sub>12</sub>
33	79	26	ATAT ATAATAAACACAGTATCAGAGACAGATATAGAAGTGATGATA T <sub>16</sub>
63	113	1,006	ACATAAGAAAACATAAAGAAAATCTGTGAGTAGAGTGATAAGTTATAAT T <sub>11</sub>
65	113	411	ACATAAGAAAACATAAAGAAAATCTGTGAGTAGAGTGATAAGTTATA T <sub>15</sub>
64	113	84	AT ACATAAGAAAACATAAAAAGAAAATCTGTAGTAGAGTGATAAGTTATAA GT <sub>14</sub>
63	108	59	ATACAT AAAAACATAAAAAGAAAATTATAAGTAGAGTGATAGATTATAAT T <sub>09</sub>
65	108	33	ATACAT AAAAACATAAAAAGAAAATTATAAGTAGAGTGATAGATTATAA TTTT
66	108	27	ATACAT AAAAACATAAAAAGAAAATTATAAGTAGAGTGATAGATTAT T <sub>12</sub>
68	108	14	ATACAT AAAAACATAAAAAGAAAATTATAAGTAGAGTGATAGATT T <sub>10</sub>
99	143	826	AT ATGAAAACAATCAAAGAAGTGTGATAGAAAGTATAAAGGTATAA T <sub>11</sub>
98	141	542	ATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATAAA T <sub>12</sub>
101	141	501	ATATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATA T <sub>09</sub>
98	141	240	ATATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATAA T <sub>13</sub>
100	141	199	ATATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATA T <sub>11</sub>
99	141	188	ATATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATAA TCT <sub>13</sub>
100	141	168	ATATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATA T <sub>18</sub>
99	141	52	ATATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATAA T <sub>10</sub>
98	141	24	ATATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATAAG T <sub>09</sub>
100	141	21	ATATAA GAAAACAATCAGAGAAAATGCGGTAAGAGATAAGAGATATG T <sub>20</sub>
130	170	507	ATATAT CAAATCACAGGAAGATCATAGATGACGATGAAGGTAGTTAA TAT <sub>18</sub>
130	170	413	ATACAT CAAATCACAGGAAATCATAGATGGCAATGAAGATAAGTTAA T <sub>11</sub>
130	174	2,397	ATA TATACAAACTACACGGAGATCAAATAACAGTGAGATGATTAA T <sub>12</sub>
132	174	133	ATA TATACAAACTACACGGAGATCAAATAACAGTGAGATGATT T <sub>15</sub>
158	205	412	T ATGTATAAAACATTAATGTGAGTTGTGTCGTGAGATTATGTGAA T <sub>15</sub>
155	189	236	AAATAAAACACC AACGTGAATTATATTGTATAGATCGTATGAGAAT AT <sub>15</sub>
158	205	1,756	ATAT ATGTATAAAACATTAATGTGAGTTGTGTCGTATAGATTATGTGAA T <sub>14</sub>
190	229	4,864	ATATAT ACAACTAACAGAATATAGATCTGATGTATAAGATATCA T <sub>14</sub>
190	229	123	ATATAT ACAACTAACAGAATATAGATCTGATGTATAAGATATCG T <sub>13</sub>
190	233	227	ATAT ACAAAACAACGTGAGACATAGATTGATGTATAAGATATCA T <sub>13</sub>
198	234	139	ATAT AAACAAACAACATAAAATGAAATCTGATGTGTGA TATACT <sub>08</sub>
191	233	40	ATAT ACAAAACAACGTGAGACATAGATTGATGTATAAGATATC T <sub>11</sub>
223	263	3,225	TTAT ATACAAATAATGGGATTTAACGATATAAGAGGTGAATGATT T <sub>11</sub>
222	263	3,165	ATAT ATACAAATAATGGGATTTAACGATATAAGAGGTGAATGATTA T <sub>15</sub>
222	264	384	AAAAAT AACACAGATAATGGAATTAAATGATATGAGAAATGGATGATTA T <sub>05</sub>
223	264	164	AAAAT AACACAGATAATGGAATTAAATGATATGAGAAATGGATGATT TCT <sub>17</sub>
223	263	106	ATAT ATACAAATAATGGGATTTAACGATATAAGAGGTGAATAATT TTTT
222	263	232	AAAT ATACAGATAATGGGATTTAACGATGTGAGAGATAGATAATT T <sub>16</sub>
223	263	177	AAAAT ATACAGATAATGGGATTTAACGATGTGAGAGATAGATAATT T <sub>05</sub>
220	263	32	AAAT ATACAGATAATGGGATTTAACGATGTGAGAGATAGATAATTAT T <sub>13</sub>
253	299	260	ATAT GATAAAACAACACTATTATGAGAACGAGTGATAGATAATAGATAAT TTCT <sub>14</sub> *
253	299	106	ATAT AATAAAACAACACTATTACAAAGATAGACAGTGAGATATAGATAAT T <sub>13</sub>
253	298	83	ATACAT ATAAACAAACACTATTACAGAACGAGACAGTGAGATATGAGATAAT T <sub>15</sub>
253	299	133	ATAT AATAAAACAACACTATTACAAAGATAGACAGTGAGATATAGATAAT T <sub>07</sub>
252	299	89	AT AATAAAACAACACTATTACAAAGATAGACAGTGAGATATAGATAATG AT <sub>06</sub>
253	299	57	ATAT AATAAAACAACACTATTACGAAGATAGACAGTGAGATATAGATAAT T <sub>12</sub>
252	299	22	AT AATAAAACAACACTATTACGAAGATAGACAGTGAGATATAGATAATG AT <sub>16</sub>
284	329	786	ATAT AAAACCACAAAAATAGAAAGCTATAATAGAGATAGAATAATGTTA A <sub>07</sub>
288	320	435	ATTTTAAGTT AGAGTGAGAAATTGTAGTGAAATAAGATGATA AAAT <sub>11</sub>
285	329	282	ATAT AAAACCACAAAGGTAGAAGATCGTAATAGAGATAGAATAATATT T <sub>09</sub>
300	333	268	AT AACAAAAACACAGAGATGAGAGATTGTAATAAG TATAGTGATAAT T <sub>13</sub>
285	329	168	ATAT AAAACCACAAAAATAGAAAGCTATAATAGAGATAGAATAATGTT T <sub>12</sub>
282	329	64	ATAT AAAACCACAAAAATAGAAAGCTATAATAGAGATAGAATAATGTTATT CT <sub>09</sub>
285	328	3	AT AAACCACAAAGATAGAAGGCCATAATAGAGATAAAATAATT T <sub>10</sub>

5'	3'	Reads	ND3 gRNA Sequences cont.
322	369	48,018	ATAT AATACCACATGAATCTTATATGTACGATGGAAGATGAGAATTAT T <sub>13</sub>
324	369	5,034	ATAT AATACCACATGAATCTTATATGTACGATGGAAGATGAGAATT T <sub>11</sub>
321	369	3,062	ATAT AATACCACATGAATCTTATATGTACGATGGAAGATGAGAATTATG TTTCT <sub>06</sub>
322	369	142	ATAT AATACCACATGAATCTTATATGTACGATGGAAGATGAAAATTAT TCT <sub>10</sub>
322	370	721	AT AAATACCACATGAATCTTATATGTACGATGGAAGATGAGAATTAT T <sub>12</sub>
321	370	291	AAATACCACATGAATCTTATATGTACGATGGAAGATGAGAATTATG CAGT <sub>08</sub>
320	368	149	ATAT ATATCACACAAATTCTATACATAATAGAGAATGAGAGTTACAA T <sub>06</sub>
324	370	43	AT AAATACCACATGAATCTTATATGTACGATGGAAGATGAGAATT T <sub>05</sub> *
347	388	95	ATATGTAAC AATATACGTGATTTAGAGATACTATGTGAATTCTAT T <sub>22</sub>
345	388	82	ATATGTAAC AATATACGTGATCTCAGAGATACTATGTGAATTCTATGT T <sub>07</sub>
355	388	1,307	ATATAAACAAAC AATATATGTGGTTTCGAAAGTGTCATGTGAATT AT <sub>12</sub>
350	388	633	ATAGAT AATATACGTGATCTTAGAGGTACCATGTGAGTCT GAGTTAT <sub>12</sub>
349	388	376	AGTATGCGTGATTTAGAGATATTGTATGAATT T <sub>08</sub>
355	388	104	AGTATGCGTGATTTAGAGATATTGTATGAATT ATAT <sub>15</sub> *
402	438	417	ATATA TATAATACAACAAGGAGCGTCATAAGTAAAGTGAA TTCGTTATAT T <sub>13</sub>
402	438	128	ATAT TATAATACAACAGAAAATGTCTATAAGTGAGATGAA TTCGTTATAT T <sub>09</sub>
402	435	753	ATATATAA AATACAACAAGAGACGTCGTAATAGAGTAA TTCGT <sub>08</sub>
403	438	126	ATAT TATAATACAACAGAAAATGTCTATAAGTGAGATGAA TTCGTTAT T <sub>06</sub>

## H) NADH Dehydrogenase Subunit 7

5'	3'	Reads	ND7 gRNA Sequences
36	69	100,761	ATATA ATAAATGTAAGAGACTATTGAGAGTGGCATAAG TGATATTAT <sub>14</sub>
35	69	765	ATATA ATAAATGTAAGAGACTATTGAGAGTGGCATAAGG GATATTAT <sub>11</sub>
36	71	240	AAAT ATACAAATGTAAGAGACTATTGAGAGTGGCATAAG TGATATAAT <sub>13</sub>
36	69	121	ATATA ATAAATGTAAGAGACTATTGAGAGTGGCATAAG TGATATTAT <sub>13</sub>
38	71	113	ATAT ATACAAATGTAAGAGAGCTATCGAGAGTGCATA TGTGATATTAT <sub>11</sub>
28	71	35,079	AAAT ATACAAATGTAAGAGAGACTGCCGAAAGTAACGTAGAATGATATT AT <sub>12</sub>
27	71	20,487	AAAT ATACAAATGTAAGAGAGACTGCCGAAAGTAACGTAGAATGATATT TT <sub>08</sub>
31	71	546	ATAAAT ATACAAATGTAAGAGAGACTGCCGAAAGTAACGTAGAATGAT TTTT
34	71	402	AAAT ATACAAATGTAAGAGAGACTGCCGAAAGTAACGTAGAAT T <sub>15</sub>
36	69	354	ATATA ATAAATGTAAGAGAGACTATTGAGAGTGGCATAAG TGATATTAT <sub>12</sub>
24	71	327	ATAAAT ATACAAATGTAAGAGAGACTGCCGAAAGTAACGTAGAATGATTTGTT TTTT
24	71	298	AAAT ATACAAATGTAAGAGAGACTGCCGAAAGTAACGTAGAATGATTTATT T <sub>11</sub>
28	71	194	AAAT ATACAAATGTAAGAGAGACTGCCAAAAGTAACGTAGAATGATATT ATCT <sub>09</sub>
59	91	4,376	ATATATAGATGCA GTGGGTGAAATGTAAGATGATATAGATGTGAA TGT <sub>13</sub>
58	91	1,242	ATATATAAACATA GTGGACTAAATGTAAGCGATAAAATGTGAAA T <sub>12</sub>
108	137	13	ATATAATAAACACATAAAAGTGCATGTACT#CGAGAT TTTAG
95	132	1	AT ATAAACACATAAAACTATGTGATGTAGGAT CTGTGAATTAAT <sub>09</sub>
124	170	3,944	ATAC ATCAATATAAACGATAGATTACCGTAGAAGTATAGTGAATAAT T <sub>14</sub>
124	168	3,640	ATATA CAATATAAACAGTAGATTCACTGCAGAAGTATGATAGATAAT T <sub>11</sub>
139	170	1,507	AT ATCAATATAAACATAAGTCGTATAGA TTTACAGTAGATAAT <sub>12</sub>
121	166	396	ATACAT ATATAAACATGAATTCACTGTGAAGAGATACGATAGATGATA T <sub>14</sub>
139	170	384	AT ATCAATATAAACATGAGTTCTGTATAGA TTTACAGTAGATAAT <sub>12</sub>
121	166	178	ATAT ATATAAACATAAAATTCATCGTAAGGTACAGTAAATGATA T <sub>16</sub>
122	166	26	ATAT ATATAAACATAAAATTCATCGTAAGGTACAGTAAATGATA T <sub>15</sub>
152	190	988	AAATTACGATGCAT AATAATCTATGGTACAGTTGATATGAGTGATAA T <sub>12</sub>
147	199	563	ATATATA CACGATCGAGATAATCTATAGTATGATTGATATAAGTGATAAATTT T <sub>09</sub>
150	199	200	ATATATA CACGATCGAGATAATCTATAGTATGATTGATATAAGTGATAAAT AT <sub>12</sub>
151	199	2,089	AAA TACGATGTAATAACCTGTTAGTATAGTTAGTGTAGATGATA T <sub>12</sub> *
147	179	295	ATATATAAACATAACG TGTAAACAGTCATATAGATGATAAATTT T <sub>09</sub>
152	190	143	AAATTACGATGCAT AATAATCTATGGTACAGTTGATATGAGTGATAA T <sub>13</sub>
246	269	2,470	ATATCAAC ACATAATCTGACTTGTGGAGTAT CTAAGGAATAAAT <sub>14</sub>
246	269	837	ATATCAAC ACATAATCTGACTTGTGGGTAT CTAAGGAATAAAT <sub>12</sub> *
261	311	3,259	ATAT AACATAAACATAAGTCTTATTACAGTGAACATTGATATAATTT T <sub>08</sub>
261	293	2,874	ATATATAAACAC GATGCTCATTATGATAGATACTGATGTAATTT T <sub>10</sub>
261	310	146	ATAT ACGTAACAGACAATAGGTGTTATTGCACTAGATATTGATGTAATTT T <sub>10</sub>
260	293	133	ATATATAAACAC GATGCTCATTATGGTAGATACTGATGTAATTAA T <sub>07</sub>
261	293	1,105	ATATATATAAACAC AGTGCTCGTTACAGTGAATATTGATGTAATTT T <sub>11</sub> *
297	338	888	AT ATAAATAACATCGAACGTATATTGAAATGTTAGAGATA CAT <sub>13</sub>
295	338	765	ATATA ACAACAAACATCGTAATATGTCGAGTAGAGATAAT TAAATAT <sub>13</sub>
292	334	91	ATACT ACAACATCGCGATATACTTGAATGTAAGGTGATAAA GT <sub>11</sub>
292	324	394	ATATATAAACACATCG AGCATATACTCAGAATATAAAGGTGATAAA GT <sub>12</sub>
292	324	178	ATATATAAACACATCG AGTATATACTCAGAATATAAAGATGATAAA GT <sub>09</sub>
293	324	90	ATCGG AGCATATACTTGAGATAAAAGATGATAA T <sub>11</sub>
327	373	2,128	A TATAATTAATAAACGTATAATGTGAGTGTGACGTGATGATATT AT <sub>13</sub>
327	365	1,225	ATACAT ATAAACGTATAGGTGCTATGTAACGTGATGATGTT AAT <sub>11</sub>
327	365	20,110	ATACAT ATAAACGTATAGTGCCTATGTAACGTGATGATGTT AAT <sub>12</sub>
327	373	997	A TATAATTAATAAACGTATAATGTGAGTGTGACGTGATGATGATATT T <sub>13</sub>
330	378	517	AAA TTACAATTAAATAGACGTATAAGTCATAGTGTAGTGTAGGATAAT T <sub>17</sub>
329	378	111	AAA TTACAATTAAATAGACGTATAAGTCATAGTGTAGTGTAGGATAATG AT <sub>13</sub>
352	398	283	ATATA GAAAATCACAGGTAAATTCTGCAATTAGTAGACGTGTAAT AT <sub>07</sub>
353	402	3,043	ATATA TATTAACACTATGGTAGATTCTGTAATTATAGACGTGTAAT AT <sub>13</sub> *
352	385	602	ATATATAACACTACGA GTAGATTCTATGATTGATGAAACGTGTAAT T <sub>11</sub> *
354	402	556	ATATA TATTAACACTATGGTAGATTCTGTAATTATAGACGTGTAAT T <sub>12</sub>
352	402	173	ATATA TATTAACACTATGGTAGATTCTGTAATTATAGACGTGTAAT TTCTTTT
390	424	2,283	ATATTAA ATACATGATATACCGCATAGATTATTAGAGTTATG AGTTAAT <sub>16</sub>
391	427	200	AAATT ACCATACATGATATACAGTGAACATTAGAATTAT AGGTAAATGT <sub>06</sub>
390	424	162	ATATTAA ATACATGATATACCGCATAGACTATTAAAGTTATG AGTTAAT <sub>12</sub>
391	424	298	AAATTACA ATACATGATATACAGTGAACATTAGAATTAT AGGTAAAT <sub>14</sub>

5' 3' Reads

ND7 gRNA Sequences cont.

412	452	15,183	ATATATAA GGAGACAAATGATCTAGATTGAGACTGTATATGATATAT T <sub>13</sub>
414	451	223	ATATATAACAAC GAGACAGATAATCTAGATTTGAAGTTATATGTATAT T <sub>12</sub>
416	452	199	ATATATAA GGAGACAAATGATCTAGATTGAGACTGTATATGAT T <sub>13</sub>
416	464	5	AT ATTATAATAACGGAGATGAGCAATTAGATTGAGCTTCAAGAGTTATATGTAT T <sub>13</sub>
407	450	3,126	ATATA AGACAAACAATCTAAATCTGAGACTGTATATGATATGTATAAT T <sub>10</sub>
414	458	2,332	ATAT ATAACGGAGACAGATACTGAACTAAGGTTATGTGATAT T <sub>12</sub>
410	450	132	ATATA AGACAAACAATCTAAATCTGAGACTGTATATGATATGTAT T <sub>13</sub>
453	485	3	ATAAAATGTCATCAATTA GTTACGTTCTTGAGTTGAGTTGATAAT AT <sub>13</sub>
453	485	207	ATAAAATAACATCAATTA GTTACGTTCTTGAGTTGAGTTGATAAT AT <sub>13</sub>
486	530	2,523	ATATAA GCATACGACAATTACGATATGAGTCAGAGAATGTTGTTAATTT T <sub>11</sub>
499	526	2,002	ATATATAA ATACGACAATCATAACGTGAATCAGAGA CTGTGATTAAT T <sub>11</sub>
486	530	1,136	ATATAA GCATACGACAATTACGATATGAGTCAGAGAATGTTGTTAATTT T <sub>11</sub>
486	528	344	ATATA ATACGACAATCACGATATAGATTTAAAGATGTTGTTAATTT T <sub>10</sub> *
508	548	1,129	ATAT AAATCATGAAAGCTGAGTGTGTATGGCAGTTACGATATA TGTTAAT T <sub>14</sub>
508	544	116	ATATAAAA CATGGAAGCTAAGTGTGTATGATGATTATGATATA TGATTAAT T <sub>06</sub>
508	553	1	ATA TAATAAAACCATGAGAGTTGGATGTATATGATGATCGTGTAT T <sub>17</sub>
526	569	3,050	ATATAT ATCATCAAGAATATCTAATGAGACTATGAGAGTTAAATGTATA AT <sub>12</sub>
526	569	4,850	ATATAT ATCATCAAGAATATCTAATGAGACTATGAGAGTTAAATGTATA AT <sub>13</sub>
527	569	351	ATATAT ATCATCAAGAATATCTAATGAGACTATGAGAGTTAAATGTAT T <sub>06</sub>
540	576	544	ATAATAAT AAACAAAATCATTGAGAGTATCTGATAGAGTTATGA TAGTCAT T <sub>05</sub>
531	574	291	AAAT ACAAAATCATCAGGGATACTGGTAAGATTGTGAAAGTTAAGT T <sub>16</sub>
540	571	8,794	ATATAAT AAATCATCAAGAATATCTGATAGAACTGTGA TTGCTAAT T <sub>14</sub> *
540	571	1,694	ATACAAT AAATCATCAAGAGTATCTAATAGAACTGTGA TTGCTAAT T <sub>14</sub>
564	615	1,896	AT ATATTATCAACACAAATAGAAGATTGGCGAAATTAGAGATAGAATTATT T <sub>12</sub>
567	615	360	AT ATATTATCAACACAAATAGAAGATTGGCGAAATTAGAGATAGAATT T <sub>09</sub>
562	603	267	ATATT ACAACACAAAGAAATCAATGAAGTCAGAGATAAAAGTTATTAA T <sub>15</sub>
564	596	173	TAGATATCAACACAT CAGAGAATCAATGAAACTAGAGATAGAGTTATT T <sub>10</sub>
567	611	1,316	ATACA TATCAACACGACAAGAGATCAGTGAATTAGAGTAAAGTT T <sub>13</sub> *
564	615	1,043	AT ATATTATCAACACAAATAGAAGATTGACGAAATTAGAGATAGAATTATT T <sub>10</sub>
564	611	754	ATACA TATCAACACGACAAGAGATCAGTGAATTAGAGTAAAGTTATT T <sub>09</sub> *
567	615	703	AT ATATTATCAACACAAATAGAAGATTGACGAAATTAGAGATAGAATT T <sub>15</sub>
563	611	188	ATACA TATCAACACGACAAGAGATCAGTGAATTAGAGTAAAGTTATCA T <sub>06</sub>
564	611	146	ATACA TATCAACACGACAAGAGATCAGTGAATTAGAGTAAAGTTATC T <sub>12</sub>
568	603	131	ATATT ACAACACAAAGAAATCAATGAAGTCAGAGATAAAAGT AT <sub>08</sub>
584	630	831	ATAAT TAACAAACAAATATGATATTATCAGTGACAGTGAGAAATTGATA T <sub>15</sub>
596	642	20	ATA TATAACAATCCATAGCAGATAGACGTGATATTATTGATGATAGT TTAAAT T <sub>18</sub>
596	640	512	AAAT TAACAATCCATAACAAGTGGCGTGATATTGTCAATGATAAT T <sub>11</sub>
596	648	159	ATAAAATCATAACAATCTATAATAGACGAGCGTGATATTGTCAATGATAAT T <sub>10</sub>
596	642	32	ATA TATAACAATCCATAGCAGATAGACGTGATATTATTGATGATAAT T <sub>12</sub>
629	670	1,258	ATAT GATAAACGATTACCTACAGATAATGAGTCATAGTGTATTTATA T <sub>13</sub>
632	674	463	ATAT ATAAAATAACGATTACCTGTGAATGATAGATTATGATGATTT T <sub>11</sub>
630	674	419	ATAT ATAAAATAACGATTACCTGTGAATGATAGATTATGATGATTTAT T <sub>26</sub> *
630	669	299	ATACAT ATAAACGATTACTCATAGATAGACGAGTCATAGTGTATTTAT T <sub>08</sub> *
630	671	101	ATT ATTAAACGATTACTTACGAGTGACAGATTGTGATGATTTAT T <sub>14</sub>
629	671	65	ATATT ATTAAACGATTACTTACGAGTGACAGATTGTGATGATTTATA T <sub>15</sub>
656	699	59	ATATAT ATGACAAAATACGTAAGTGCAGATAAAAGTAAGTGTATTT T <sub>12</sub>
666	701	2,991	ATAT AGATGACAAACCAGTAGACGTGAATAAGATAG TGATTGATCAT T <sub>12</sub> *
679	727	2	AAAAAA AAAAACTAAATCATATAAATTAAAGGGTGTGAACCTGTGAAAT T <sub>13</sub>
679	711	1	ATAACTCAGAAAGTGTAGATCGTGTAAAT T <sub>15</sub>
679	725	1,728	AAATTAATCATATAAAGTGGAAAGATGCCAGATTGTGTAAT TGTAGT T <sub>16</sub>
679	727	373	AAAAAA AAAAACTAAGTCATATAAGTCAGAAAGTGTAGATCGTGTAAAT TTT
711	758	3,229	ATATAT ACGAGACAAAATACCTAGATTAGAGATTGAGTTATA AT <sub>13</sub>
710	741	22	AA TTGACTATTAGAGATTGAGTTATAT TTT
709	741	2	AAAACGAGACAAGATACCAA TTGACTATTAGAGATTGAGTTATATA T <sub>14</sub>
725	764	55,757	ATATAA TAACGAACGAGGCAGAGTATCATTAGACTATTAGA TTCAAT T <sub>14</sub>
731	759	2,867	ATATATAAT GACGAGATAAGACATCACTTAGACTGT AGAGAT T <sub>14</sub>
722	765	8,175	ATATA TTAACGAACGAGATAGAACATTGCTTGAGTTATTGAGAAT AT <sub>14</sub> *
719	765	552	ATATA TTAACGAACGAGATAGAACATTGCTTGAGTTATTGAGAAT T <sub>12</sub>
735	765	115	ATA TTAACGAACGAGATAGAACATTGCTTGAGTT TT

5' 3' Reads

ND7 gRNA Sequences cont.

756	794	48,263	AT ACTAAATAAACGACGATCTTATACTGTATCTGATGAATG TGATATTAAT <sub>14</sub>
757	794	318	AT ACTAAATAAACGACGATCTTATACTGTATCTGATGAAT ATGATATTAATT <sub>TT</sub>
751	794	174	AT ACTAAATAAACGACGATCTTATACTGTATCTGATGAATGGGATA TTAAAT <sub>14</sub>
755	794	167	AT ACTAAATAAACGACGATCTTATACTGTATCTGATGAATGA TATTAATTGT <sub>07</sub>
756	794	2,829	AT ACTAAATAAACGACGATCTTATACTGTATCTGATGAATG TGATAT <sub>14*</sub>
772	806	367	ATATA TCATAACAACAGATAAACGATGATCTCACA GTATAGTTAAT <sub>12</sub>
782	816	77	ATAACTTATAACAACAGATAAACGAGG GAATCTTATATTGTAGTTAAT <sub>16</sub>
778	830	25,793	ATAT ATAAAACATAAAATATGACTTGAGCACTTAAGTGAATGATGAT GAT AT <sub>12</sub>
777	822	5,729	ATATAT TAAAATACAACCTTATGATGACTAAGTGAATGATGATT CAAT <sub>10</sub>
775	830	2,644	ATAT ATAAAACATAAAATATGACTTGAGCACTTAAGTGAATGATGATT TT T <sub>08</sub>
781	833	2,373	AAATA ATAACAAAACATGAGATATAACTTGAGTAGATGAATGAT T <sub>11</sub>
780	833	2,079	AAATA ATAACAAAACATGAGATATAACTTGAGTAGATGAATGATA T <sub>13</sub>
778	833	903	AAATA ATAACAAAACATGAGATATAACTTGAGTAGATGAATGATAGT AT <sub>13</sub>
778	830	450	ATAT ATAAAACATAAAATATGACTTGAGCACTTAGTGAATGATGAT ATTAAAT <sub>13</sub>
778	830	331	ATAT ATAAAACATAAAATATGACTTGAGCACTTAGTGAATGATGAT ATTAAAT <sub>14</sub>
778	830	134	ATAT ATAAAACATAAAATATGACTTGAGCACTTAGTGAATGATGAT AT <sub>07</sub>
792	845	98	ATAT AAAACAATAATCATGATGAGATATAAAGTCAGTTGTGATAATT T <sub>10</sub>
790	839	52	ATAT ATAATCATAACAAGATGTAGAGTACGATTTATGATTAAT T <sub>12</sub>
790	839	343	ATAT ATAATCATAACAAGATATAGAGTACGATTTATGATTAAT T <sub>11</sub>
792	844	190	ATAT AAAACAATAATCATGATGGGATATAAAGTCAGTTGTGATAATT T <sub>10</sub>
790	839	99	ATAT ATAATCATAACAGAGCATAGAATAACAGTTATAGTATTAAT T <sub>TG</sub> T <sub>05*</sub>
790	844	89	ATAT AAAACAATAATCATGATGGGATATAAAGTCAGTTGTGATAATTAA TAT <sub>05</sub>
843	877	26,437	ATAT ATAAACGGTCAAATGTTACTTATAAGATAGAG TTGATAAT <sub>14</sub>
834	877	1,279	AT ATAAACGGTCAAATGTTACTTGTGCTTAGGTAGAGGTGATAATT T <sub>13</sub>
829	865	188	ATATATAACGATC AATGCATTATCTATGAAGCAAGAACAGTGATTATAAT T <sub>15</sub>
831	857	121	AATGTGTG ACCTATAAAAGTGAGAATAATGATTATA T <sub>12</sub>
832	865	23	ATATATAACGATC AATGCATTATCTATGAAGCAAGAACAGTGATTAT TTTT
832	857	20	AATGTGTG ACCTATAAAAGTGAGAATAATGATTATA T <sub>11</sub>
834	879	1	ATAT AAATAAACGATTAGATGATCACTTATAGAATAAAATGATGATT T <sub>15</sub>
863	902	2,491	ATATA ATATGCCATATCAGATAGACGTAGAATAATGATTTAAT T <sub>10</sub>
866	902	126	ATATAATAA ATACGCATATTAGATGAGTGAGTAAAGTAAGTGATTA T <sub>10</sub>
872	916	3	ATATAA AATCAACAAATTGATGATATCGAGTAAATGAGAAATAAT T <sub>09</sub>
872	919	3,227	A TACAAATCAACAAATTGCGCTATATTAAGTGGGTGAGAAGTGAAT T <sub>17*</sub>
867	919	528	A TACAAATCAACAAATTGCGCTATATTAAGTGGGTGAGAAGTGAATGATT T <sub>15</sub>
872	919	295	A TACAAATCAACAAATTGCGCTATATTAAGTGGGTGAGAAGTGAAT TAAAT <sub>19</sub>
871	919	142	A TACAAATCAACAAATTGCGCTATATTAAGTGGGTGAGAAGTGAATG T <sub>15</sub>
901	939	1,365	ATATAT ACTAATAAAAGGCATTGCTTACAGATTGATAGATTAT T <sub>12</sub>
899	931	292	ATATATTACCAACAT AGAGACATTGCTTATGAGTTAACAGATTATAT T <sub>12</sub>
901	939	249	ATATAT ACTAACAAAAAGATATTGCTTATAGATCAGTGAATTAT T <sub>11</sub>
907	947	24	ATAAAGAACCAACAGAAGAATATTGCTGTAAGTTAATGA TCTTGTAT <sub>10</sub>
907	951	64	CCAAGGCA AAAGATAAAGAACACAGAACAGAACATTGCTGAGTTAATGA TCT <sub>17</sub>
932	978	44,852	ATATAT TCAAAACAAACAGATAGAACCGGAGACGAGAACATTGATAA T <sub>13</sub>
934	988	11,706	ATATAAATAATCAAACAGACGAATGAAACTAGAGATAGAGAAATTAAAT T <sub>12</sub>
931	978	3,513	ATATAT TCAAAACAAACAGATAGAACCGAGACGAGAACATTGATGA T <sub>12</sub>
933	986	2,467	ATAAATAATCAAACAGACGAATGAAACTAGAGATAGAGAAATTAAATA T <sub>12</sub>
940	978	182	ATATAT TCAAAACAAACAGATAGAACCGGAGACGAGAAC CTTGATAAT <sub>05</sub>
936	988	175	ATATAAATAATCAAACAGACGAATGAAACTAGAGATAGAGAAATT T <sub>TAT</sub> <sub>15</sub>
941	978	137	ATAT TCAAAACAAACAGATAGAACCGGAGACGAGAAC TATTGATAATTAAAT T <sub>14</sub>
937	988	97	ATATAAATAATCAAACAGACGAATGAAACTAGAGATAGAGAAATT TAT <sub>13</sub>
934	983	18,961	ATATAT AATAATCAAACAGACGAATGAAACTAGAGATAGAGAAATTAT T <sub>12*</sub>
933	983	1,729	ATATAT AATAATCAAACAGACGAATGAAACTAGAGATAGAGAAATTAAATA T <sub>05</sub>
936	983	472	ATATAT AATAATCAAACAGACGAATGAAACTAGAGATAGAGAAATT T <sub>14</sub>
937	983	397	ATATAT AATAATCAAACAGACGAATGAAACTAGAGATAGAGAAATT T <sub>11</sub>

5'	3'	Reads	ND7 gRNA Sequences cont.
959	1000	173,548	ATAAA GGTAAATCACAGTGTAGATAGTCGGATAGATAGATGAAA AT <sub>14</sub>

952	1000	748	ATAAA GGTAATATCACAGTGTAGATAGTCGGATAGATAGATGAAATT T <sub>09</sub>
960	1000	472	ATAAA GGTAATATCACAGTGTAGATAGTCGGATAGATAGATGAA T <sub>13</sub>
966	1000	412	ATAAA GGTAATATCACAGTGTAGATAGTCGGATAGATA TATGAAAAT <sub>13</sub>
964	1000	372	ATAAA GGTAATATCACAGTGTAGATAGTCGGATAGATA AAAAAT <sub>14</sub>
963	1000	371	ATAAA GGTAATATCACAGTGTAGATAGTCGGATAGATA T <sub>15</sub>
961	1000	265	ATAAA GGTAATATCACAGTGTAGATAGTCGGATAGATA T <sub>12</sub>
962	1000	234	ATAAA GGTAATATCACAGTGTAGATAGTCGGATAGATA TAAATTAT <sub>09</sub>
959	1000	3,928	ATATAAA GGTAATATCACAGTGTAGATAGTCGGATAATAGATGAAA AT <sub>11</sub>
983	1017	101	ATATATAAATAACATAT TAATGGCTTGATGGTGTATTGTAGTATAA T <sub>12</sub>
1001	1032	3	ATATA TATAAAATAACGTGTAGTGTCTCGAT AGTAATTGATAAT <sub>08</sub>
1000	1043	1	ATAT ACACCACAAACTGTAGAATGACGTGTAGTGGTTTCAGTG ATAAT <sub>15</sub>
1015	1043	12	AT ACATCACAACATAAAGTAATATGATAA GGGTTTCGAT <sub>15</sub>
1013	1043	208	ATAT ACACCACAACTGTAGAATGACGTGTAGTG AT <sub>18</sub>
1032	1067	6,797	ATATATACAAGC AATGATGTACTCGTAAATAGTGACACTGTGAATT T <sub>12</sub>
1030	1067	2,406	ATATATACAAGC AATGATGTACTCGTAAATAGTGACACTGTGAATTAT T <sub>12</sub>
1032	1078	751	A AATATAAGCAAATGATGTATTGGTGGAGCAGTGATATCGTAAATT T <sub>10</sub>
1057	1085	44	ATTAGAAA GGTGTCGGTATAGTAGATGATATAT AT <sub>10</sub>
1055	1085	40	CTCATTAGAAA GGTGTCGGTATAGTAGATGATATATT TTT
1089	1121	1	AT ATATGATAACAAACAATACTTACTTT GAGGTGTTTGAGTAAAT <sub>14</sub>
1087	1113	412	ATA AATAACAAACAATATTGGCTTTG AGATATTGATATAAGTAAAT <sub>12</sub> †
1099	1143	81,287	ATATAA GAGAACATAAAACTAGCATAGAGGCATAGTAACAAGTGATAT AT <sub>14</sub>
1094	1143	6,615	ATATAA GAGAACATAAAACTAGCATAGAGGCATAGTAACAAGTGATATTT T <sub>10</sub>
1099	1142	5,793	ATATAAA AGAACATAAAACTGACACGAGGGTATAGTGATGAATGATAT AT <sub>13</sub>
1099	1131	705	ATATAAGAACACAAA CAGCATAGAGGCATAGTAACAAGTGATAT AT <sub>14</sub>
1101	1143	550	ATATAA GAGAACATAAAACTAGCATAGAGGCATAGTAACAAGTGAT T <sub>13</sub>
1105	1143	520	TATAA GAGAACATAAAACTAGCATAGAGGCATAGTAACAAG GGATATAT <sub>13</sub>
1106	1143	318	ATATAA GAGAACATAAAACTAGCATAGAGGCATAGTAACAA T <sub>11</sub>
1108	1143	303	ATATAA GAGAACATAAAACTAGCATAGAGGCATAGTAAC T <sub>11</sub>
1094	1142	289	ATATAAA AGAACATAAAACTGACACGAGGGTATAGTGATGAATGATATTT T <sub>10</sub>
1099	1143	259	ATATAA GAGAACATAAAACTAGCATAGAGGCATAGTAACAAGCGATAT ACT <sub>12</sub>
1107	1143	177	ATATAA GAGAACATAAAACTAGCATAGAGGCATAGTAACA T <sub>12</sub>
1099	1143	163	ATATAA GAGAACATAAAACTAGCATAGAACGATAGTAACAAGTGATAT AT <sub>05</sub>
1099	1143	162	ATATAA GAGAACATAAAACTAGCATAGAGGCATAGTAACAAATGATAT AT <sub>05</sub>
1099	1143	128	ATATAA GAGAACATAAAACTAACATAGAGGCATAGTAACAAGTGATAT AT <sub>14</sub>
1099	1143	98	ATATAA GAGAACATAAAACTAGCATAGAGACATAGTAACAAGTGATAT AT <sub>14</sub>
1099	1143	356	ATATAA GAGAACATAAAATTGACATGGAAGCATAGTAATAAGTGATAT AT <sub>12</sub> *
1095	1143	212	ATATAA GAGAACATAAAATCAGTGCAGGATAGTAGTGAGTGATATT AAT <sub>09</sub> *
1108	1145	569	ATATAAACGTAT ACGAGAGCATAGATCAGTGTGAGAATGTTAGTAAT T <sub>14</sub>
1094	1148	295	ATATAAAA TAAACGAGAAATATAAAACTGTGTAGAGATATAGTGATAAGTAATATT T <sub>08</sub>
1107	1157	82	AT ATGTAAATGTAAACGAGAAATATGATTGATGTAGAGATATGATAA T <sub>13</sub>
1107	1146	90,559	ATATTAACACGT AACGAGAATGTGAACACTGACATAGAGATATGATAA T <sub>14</sub>
1108	1146	10,003	ATATTAACACGT AACGAGAATGTGAACACTGACATAGAGATATGATAAT T <sub>14</sub>
1111	1146	2,224	ATATTAACACGT AACGAGAATGTGAACACTGACATAGAGATATGAT T <sub>10</sub>
1114	1146	2,220	ATATTAACACGT AACGAGAATGTGAACACTGACATAGAGATAT T <sub>14</sub>
1110	1146	1,887	ATATTAACACGT AACGAGAATGTGAACACTGACATAGAGATATGATA T <sub>16</sub>
1108	1154	978	ATATA TAAACGTAAATGAGAATATGAAATCAGTGTGAAAATGATAAT T <sub>07</sub> *
1107	1146	400	ATTTAACACGT AACGAGAATGTGAACACTGACATAGAGATATGATA T <sub>13</sub>
1120	1146	266	ATATTAACACGT AACGAGAATGTGAACACTGACATAGAGAT T <sub>13</sub>
1113	1146	243	ATATTAACACGT AACGAGAATGTGAACACTGACATAGAGATATG TTTT
1108	1152	232	ATAT AACGTAAACGAGAGCATAAATTGATGTGAAGATGTGATAAT TTCTTT
1128	1167	3,079	ATATAT AATCCGTACAATGCGAACGTAGACGAGAAATATGAGTTAAC T <sub>14</sub>
1136	1182	1,834	ATAT ATAAATATGCAAGAAATCTGTATGATGTAGATGTGAATGAGAAATAT T <sub>10</sub>
1138	1182	866	ATAT ATAAATATGCAAGAAATCTGTATGATGTAGATGTGAATGAGAAAT T <sub>11</sub>
1128	1167	337	ATATAT AATCCGTACAATGCGAACGTAGACGAGAAATATGAGTTAAC TTT
1131	1167	203	ATATAT AATCCGTACAATGCGAACGTAGACGAGAAATATGAGTT T <sub>19</sub>
1150	1197	213	AT ATAAACATCCAATAGACGAGATGTGAGAGATTTGTATGATGTAAAT T <sub>10</sub>
1150	1195	1,077	AAATATCCAATAAACAGATATGTAGAAGGTCCGTATAATGTGAAT ATAT <sub>13</sub>

5'	3'	Reads	ND7 gRNA Sequences cont.
1181	1218	4,103	ATATATAA ATGCAATAGAAGATCACGCAAATAGATATCTGATAAT T <sub>13</sub>
1183	1224	1,399	ATA TAAATCATGCACTAAACAGATATGTAGATGGATATTTCAGTGA TAAT <sub>14</sub>

1183	1223	1,017	AAATA AAATCATGCAGTAGAGAACCGTGTAAAGTGAGTATCTGATAA TTAT <sub>11</sub>
1183	1224	138	ATA TAAATCATGCAGTGAAGAGCTACGTAAATGGATATTCACTGA TAAT <sub>11</sub>
1183	1224	2,561	ATA TAAATCATGCAGTAAAAGACTATGTAGATGAATATTCACTGA TAAT <sub>14*</sub>
1210	1257	123	ATAT ACAACATCAATATTACTTAGAACCGTAACTAGATTGTGTAAATAA T <sub>14</sub>
1233	1257	167	ATAT ACAACATCAATATTACTTAGAACG ATAACTAGATTGTGTAAAT <sub>10†</sub>
1233	1257	83	ATAT ATAACATCAATATTACCTAGAGTG TCAGTTAGATTATGTGATAAT <sub>14†</sub>
1240	1268	95	AAACTAACGATATT CGGATCTGAGAGTAACATTGATATTATT T <sub>07</sub>
1251	1282	63	ATATATAT AACTAACGATCTATGGGTTAAAGACAGTGT GAAT <sub>12</sub>
1242	1283	6	AC AACTAACGATTTACGGATTTAGAGACAGTGTAAATGTTAT AT <sub>13</sub>
1242	1270	296	A TACGGATTTCAGAAGTGTATTGATGTTAT AT <sub>15</sub>
1240	1268	216	AAACTAACGATATT CGGATCTGAGAGTAACATTGATATTATT T
1240	1270	54	ATATAAACTAACGA TACGGATTTCAGAAGTGTATTGATGTTAT T <sub>16</sub>
1239	1268	25	ATT CGGATCTGAGAGTAACATTGATATTATT T <sub>17</sub>
1240	1268	23	ATT CGGATCTGAGAGTAACATTGATATTGTTT TT
1241	1268	12	GATATT CGGATCTGAGAGTAACATTGATATTATT AT <sub>15</sub>
1269	1320	6	ATA TAAACAATCCTACAATGATCTCGTGTATAAGACTGATGATTAA AT <sub>12</sub>
1269	1320	1,074	ATA TAAACAATCCTACAATGATTCGTGTATAAGACTGATGATTAA AT <sub>17</sub>

## I) NADH Dehydrogenase 8

5'	3'	Reads	ND8 gRNA Sequences
34	58	2	GTGGG ATATGAAAGTAAGAGAATAAAAAAA ATTAAAT <sub>13</sub>
29	68	1	AA AAAAAAAAACATAACAGAAAATAGAAAGATAAGAAAGTGATA TATTAT <sub>08</sub>
28	56	9	AAA ATAGGAGTAGGAGGATGAGAAAATGATAG AGGGATTT*
55	98	2,577	ATAT AAAACAAACAAAAAGAAGAACAGAAATTGAAGAGAGATATAT T <sub>13</sub>
55	98	274	ATAT AAAACAAACAAAAAGAAGAACAGCGAGAAATTGAAGAGAGATATAT T <sub>12</sub>
54	97	236	ATAT AAACAAACAGAAAAGAAAACAAGAGTCAGAGAAAAGATATATA T <sub>17</sub>
55	97	100	ATAT AAACAAACAGAAAAGAAAACAAGAGTCAGAGAAAAGATATAT T <sub>10</sub>
57	98	98	ATAT AAAACAAACAAAAAGAAGAACAGCGAGAAATTGAAGAGAGATAT T <sub>05</sub>
57	97	65	ATAT AAACAAACAGAAAAGAAAACAAGAGTCAGAGAAAAGATAT T <sub>05</sub>
54	97	57	ATAT AAACAGACAGAAAAGAAAACAAGAGTCAGAGAAAAGATATATA TGTAAATTATTT*
59	97	29	ATAT AAACAAACAGAAAAGAAAACAAGAGTCAGAGAAAAGAT T <sub>09</sub>
87	136	6,039	ATAAATAGTAACACAATGAGCAGAGTACGTATAAGAATGAGTAAA T <sub>14</sub>
84	133	633	A AAATAGTAATACAACAGACAGACATATATAGAAAATAAGTGAGAAA T <sub>16</sub>
87	133	593	AAATAGTAACACAATGAGCAGAGTACGTATAAGAATGAGTAAA T <sub>15</sub>
87	135	394	TAAATAGTAACACAATGAGCAGAGTACGTATAAGAATGAGTAAA T <sub>12</sub>
84	131	380	AT ATAGTAATACAACAAACGAGATACTGTATAAGAATAGATGAGAAA T <sub>13</sub>
87	136	292	GTAATAGTAACACAATGAGCAGAGTACGTATAAGAATGAGTAAA T <sub>12</sub>
96	136	191	ATAAATAGTAACACAATAGACGAGATACTGTAGAA TAAGTGATTAAAT <sub>11</sub>
86	139	11,689	ATA TAAACAAATAGTAACATGACGGATAGAACGTATATGAGAATGAGTAAA T <sub>12</sub> *
85	139	1,938	ATA TAAACAAATAGTAACATGACGGATAGAACGTATATGAGAATGAGTAAAAG T <sub>13</sub>
86	139	1,171	A TAAACAAATAGTAATATGACGAATGAAGCGTATATGAGAATAAGTAAA TTTT*
87	139	816	ATA TAAACAAATAGTAACATGACGGATAGAACGTATATGAGAATGAGTAAA TTTAT <sub>14</sub>
88	132	612	ATAT AATAGTAACACAACGAATAGAACATGTATAGAGATGAATGA TAT <sub>17</sub>
84	131	598	ATAT ATAGTAACACAATAATGAGACATATATGAGAATGAGAATGAGAAA TTTT*
92	133	267	ATATA GAATAGTAACACAGCAGATAAGATACTATAGAGATAA TGACAGT <sub>07</sub>
92	132	267	ATATAT AATAGTAACACAGCAGATAAGATACTATAGAGATAA TGACAGT <sub>14</sub>
86	138	221	AA AAACAAATAGTAATATGACGAATGAAGCGTATATGAGAATAAGTAAA T <sub>06</sub>
85	138	196	AA AAACAAATAGTAATATGACGAATGAAGCGTATATGAGAATAAGTAAA T <sub>10</sub>
111	153	70,659	ATATATAATGGTA AACTCAATGGGTGATAAGTAGTAATGTGATGAAT T <sub>13</sub>
111	153	479	ATATATAATGGTA AACTCAATGGGTGATAAATAGTAATGTGATGAAT TTT
111	153	245	ATATATAATGGTA AACTCAATGGGTGATAAGTAGTAATGTGATAAAT T <sub>12</sub>
111	153	147	ATATATAATGGTA AACTCAATGGGTGATAAGTAGTAATATGATGAAT TGTAATAT <sub>15</sub>
115	153	127	ATATATAATGGTA AACTCAATGGGTGATAAGTAGTAATGTGAT T <sub>14</sub>
113	153	121	ATATATAATGGTA AACTCAATGGGTGATAAGTAGTAATGTGATGA T <sub>15</sub>
111	152	804	A TATACAATGGTT ATTCACTGGTAGACAGATAGTGTATGATAGAC TTAT <sub>12</sub>
117	170	4	ACAT ATAAACTAACAAATGGTTGATTTAGTAGTGAGTAGTAATATT T <sub>11</sub>
158	186	169,806	ATAT GAACGCAAAGATGGATTACCACGAGTTAGTAATTGATGAT AT <sub>14</sub>
161	187	124,070	ATATAAT AAACGCAAATATGGTTACTATGAACGTGATAGATTAAT T <sub>12</sub>
161	187	799	AAACGCAAATATGGTTACTATGAACGTGATAGATTAAT T <sub>11</sub>
161	198	66	ATATAAT AAACGCAAATGGTTACTATGAACGTGATAGATTAAT T <sub>14</sub>
161	207	2,257	ATATA TAATAAGACGCAAAGATGGTTACTGTGAATTGATGAGTTAAT T <sub>12</sub> *
161	199	966	ATATATAGT AAAACGAGAAAATGGTTACTATGGATTGATGAATTAAT T <sub>16</sub>
160	199	804	ATATAGT AAAACGAGAAAATGGTTACTATGGATTGATGAATTAATG T <sub>15</sub> *
186	230	24,763	ATA CAATACAACGCTCTGAATCATATCGATAAAAGTGTGAGAAAT T <sub>13</sub>
186	229	182	ATA CAATACAACGCTCTGAATCATATCGATAAAAGTGTGAGAAAT TTAAT <sub>12</sub>
186	230	109	ATA CAATACAACGCTCTGAATCATATCGATAAAAGCGTGAGAAAT T <sub>11</sub>
186	220	204	ATA CAATACAACACTCTGAATCATATCGATAAAAGTGTGAGAAAT TTAAT <sub>10</sub>
194	230	160	ATA CAATACAACGCTCTGAATCATATCGATAAAAGTG GGAGAAAT <sub>11</sub>
195	230	149	ATA CAATACAACGCTCTGAATCATATCGATAAAAGT TGAGAAATTAAAT <sub>11</sub>
196	230	199	ATA CAATACAACGCTCTGAATCATATCGATAAAAG AGTGGAAATTAAAT <sub>05</sub>
186	228	161	T ATACAACGCTCTAAATTATACCACTGTGAAAATGCGAGAAAT T <sub>14</sub>
187	228	52	ATATTAT ATACAACGCTCTAAATTATACCACTGTGAAAATGCGAGAAAT AT <sub>11</sub> *
189	229	31	ATATA AATACAACGCTCTAGATCATATCACTGTGAGAGTGTGAAA T <sub>13</sub>
213	244	173	A ACATAAACGACAGGTAATATGATGTTCTGAAT GATATCGATAAT <sub>06</sub>
209	239	19	ACAT AACGACAAGTGTATAACGTTTAAGTACA GTGATAAT <sub>19</sub>
219	245	10	AAATA TACATAAACGATGAGTAATATGACGT GTAGAT <sub>13</sub>
219	254	1	AAATTAAATCACATAATGACGAGCGATACTGTGCT GTAGATGATACT <sub>06</sub>

5'	3'	Reads	ND8 gRNA Sequences cont.
240	288	3	TATGAACATCCGATACTGAACTAGGGTAGATTGAGTTATATA T <sub>10</sub>
237	267	2	CATCCAATGTAT AATTAGGGTAGATTGAGTTATGTAATAT T <sub>15</sub>
246	271	340	ATACGAACATCCATA GTTAGATTAAGGTAAATTGAGT GATGTAAT <sub>12</sub> †
259	285	325	ATATAA GAACATCCGATGCTAGATTA TGTAGATTGAGTTATGTAATAT <sub>05</sub> †
270	291	268	ATATA TAACACGAACATCTGATGT ATAAAAGTAAGTAAATTGAT <sub>16</sub> †
276	318	51,820	ATATAC AACGATGATCACTGAGATTTACCTAATATGGATGTT AAT <sub>13</sub> *
276	319	28,492	ATATAT AAACGATGACTACTAGAACATTCTACTCAATGTGAATGTT AAT <sub>14</sub>
277	317	1,829	ATATAT ACGATGACTACCAAAATTCTATCTGATATGAATGT GATAATAT <sub>11</sub>
279	318	594	ATATAC AACGATGATCACTGAGATTTACCTAATATCGAT T <sub>17</sub>
275	318	384	ATATAC AACGATGATCACTGAGATTTACCTAATATGGATGTT T <sub>09</sub>
279	319	334	ATATAT AAACGATGACTACTAGAACATTCTACTCAATGTGAAT T <sub>15</sub>
276	314	215	GATGACTACTAGAACATTCTACTCAATGTGAATGTT AATTATGATAT <sub>14</sub>
276	311	186	ATATACAAT ATGATCACTGAGATTTACCTAATATGGATGTT AAT <sub>11</sub>
284	321	140	ATATATA CAAACGATGATTACCGAGATTTCTATTAAATATGA TTGTCTAATTT
287	319	139	ATATAT AAACGATGACTACTAGAACATTCTACTCAATGTGAAT T <sub>12</sub>
275	319	109	ATATAT AAACGATGACTACTAGAACATTCTACTCAATGTGAATGTT T <sub>12</sub>
289	318	4,761	ATAATAT AACGATGACTATCAGAACATTCTACTCGA GATGAATGT <sub>13</sub>
276	318	1,264	ATATAT AACGATGACTACTGAGACTCTATCTGACGTGAATGTT AAT <sub>12</sub>
275	318	175	ATATAT AACGATGACTACTGAGACTCTATCTGACGTGAATGTT T <sub>07</sub> *
310	353	13,915	ATATTA GATGATAACTCAGTGTAGATTGATCTGTAGATGAT AATAT <sub>13</sub>
301	344	636	ATAT ATAACCTAACATGTAGATCGATTGATGATGATTGCAA T <sub>11</sub>
310	350	312	ATGATAACTCAGTGTAGATTGATCTGTAGATGAT AATATATAATGT <sub>05</sub>
316	344	266	ATATACA ATAACCTAACATGTGAATTAAATCTGTGAGAC TAT <sub>14</sub>
301	344	4	AT ATAACCTAACATGTAGATCGATCCGTAGATGATGATTGCTAA T <sub>11</sub>
331	364	1,709	ATAT ATAACACAGGTGATAATTGATGTGA TGATATCTGTAAT <sub>16</sub>
325	358	137	ATAA ATAACGACGATGATTCACTGATGAAATTGGT ATGTGAAATGATGT <sub>12</sub>
325	355	21	ACGACGATAGCTGATGTAGGTTAGT GTGTGAGATAATGATAT <sub>17</sub>
338	382	1,964	ATATAA ATGCATACAAAGAACATAGTAAGTACAGTGTGATAATTT T <sub>09</sub>
350	386	1,269	ATATA AAACATGTATACAAAAATTACAGTAATGCGACGA AAAATAAT <sub>13</sub>
339	382	1,073	ATATAA ATGCATACAAAGAACATAGTAAGTACAGTGTGATAATT AT <sub>14</sub>
341	373	13,719	ATATATAATGCATAC AAGGCCATGATGAGATACAATGATGATAA AT <sub>11</sub>
338	385	2,361	ATGTAC AACATGCATACAGAGACTATAATAGATACAGTGTGATAATTT TTATTT
342	385	159	ATGTAC AACATGCATACAGAGACTATAATAGATACAGTGTGATA T <sub>19</sub>
343	385	133	ATGTAC AACATGCATACAGAGACTATAATAGATACAGTGTGATA T <sub>07</sub>
372	417	2	ATAT ATAATGCGTAATGGTATCTTCGGGTAGATATGGTATAAA T <sub>14</sub>
390	414	894	AAATTGTATATAT AATGCGTAATGGTGTGTTG AGCGAGTGTGTGTTAT <sub>15</sub> †
391	431	23	AT ATATATAACAAACATGGGTGTATAATGGTATCTGTAT GTGCAATTAAATTTAAAT <sub>14</sub>
405	434	9	AT AAAACACATAATAATGATGAGTCGTATAGGTATAAT T <sub>13</sub>
411	442	1,950	AAAAAA AAACAACAAGAACACATAGCAGATAGTGGGTG ACGTATAT <sub>12</sub>
411	443	923	A TAAACAACAAGAACACATAGCAGATAGTGGGTG ACGTATATGAT <sub>13</sub>
407	441	215	AAAT AAAACACAAAAACATATAATAATGATGGATGTGA TAAGGTATAAT <sub>15</sub> *
426	466	24,242	ATAAAACCTGGGA CGCTAATAGATATGGTAGATAATGAGAATATAT T <sub>13</sub>
425	466	5,218	ATAAAACCTGGGA CGCTAATAGATATGGTAGATAATGAGAATATATA T <sub>14</sub>
423	466	2,651	ATAAAACCTGGGA CGCTAATAGATATGGTAGATAATGAGAATATATA T <sub>11</sub>
428	466	1,071	ATAAAACCTGGGA CGCTAATAGATATGGTAGATAATGAGAATAT T <sub>10</sub>
423	466	203	ATAAAACCTGGGA CGCTAATAGATATGGTAGATAATGAGAATATGGT T <sub>14</sub>
426	482	13	AT AAAACCTGGGCCTAACATAGATATGGTAGATAATGAGAATATAT T <sub>12</sub>
423	462	29,495	ATACTGGGCACA CAATAGATATGGCTGATAATAGAGATATATAAT T <sub>13</sub>
426	462	11,472	ATACTGGGCACA CAATAGATATGGCTGATAATAGAGATATAT T <sub>18</sub>
425	462	7,969	ATACTGGGCACA CAATAGATATGGCTGATAATAGAGATATATA T <sub>12</sub>
428	462	4,922	ATACTGGGCACA CAATAGATATGGCTGATAATAGAGATAT TCTTT
430	462	582	ATACTGGGCACA CAATAGATATGGCTGATAATAGAGAT T <sub>14</sub>
426	464	526	ATATATAACTGGGCACA TCAATGAGTATATGGTTAAGTGTGATGAAGATATAT T <sub>10</sub>
426	480	127	ATATAT AACTGGGCCTAACATGGTAGATGAGATATAT TTT
427	462	111	ATACTGGGCACA CAATAGATATGGCTGATAATAGAGATATA ATCTTT
465	503	1	ATAT TAAAACAACAATCAAATGATAGAAGCTGGGTG ACTGTGAATATTT
477	512	176	ATATA TAAATAACATAAGACGATAACTGAATAATAGAAATT AAGTGTCAA T <sub>14</sub>
482	510	832	ATATATT AATAACATAGAGCAATGACCGAGTAATGA TAT <sub>12</sub>
477	512	237	ATATA TAAATAACATAAAACGATAACTGAATGATAGAGATT AAGTGTCAA T <sub>09</sub>

5'	3'	Reads	ND8 gRNA Sequences cont.
----	----	-------	--------------------------

489	531	21,576	ATAG ATAAAACACAAATAAAGGTCAAGTGATATAGAGTGATTAA T <sub>14</sub>
489	528	5,835	ATAAATAT AAACACAAATAAAGGTCAAGTAATGTAGAGTGATAATTAA T <sub>12</sub>
500	540	2,861	ATATAT ATAACACGAGACATGAATAGAAATTAGTGATGTAAA T <sub>12</sub>
494	536	1,242	ATATAT ACTACATACACAGATAAGAACAGATAGTGTGAGATAATA T <sub>15</sub>
495	539	884	ATAT ATAACACACAAAATATAGGTAAAAGTTAGATGTGAAATGAT T <sub>11</sub>
495	536	151	ATATAT ACTACATACACAGATAAGAACAGATAGTGTGAGATAAT T <sub>15</sub>
494	539	1,101	ATAT ATAACACAGAACATAGATAAGAGTCAGATAGTAAAGTGATA T <sub>19</sub>
495	539	405	ATAT ATAACACACAGAACATAGATAAGAGTCAGATAGTAAAGTGAT T <sub>15*</sub>
500	541	198	ATA AAATAACTACACAGAACATACGAGTAAAGATTGAATGTGAAA TAAT <sub>18</sub>
495	539	127	ATAT ATAACACACAGAGCACAAATGAAAGTTAAGTAATGTGAAATAGT T <sub>23</sub>
523	567	413	ATAAA TATAAACACAGTAAATCACTCGAGATAGATAGTTATATGAGATAT T <sub>11</sub>
554	598	20	ATATA TAATTCACCGTGAATTCTTAGATTGTAGATATGATAA T <sub>13</sub>
554	598	7	ATATA TAATTCACCGTGGATTCTTAGATTATAGATATGATAA T <sub>06</sub>
554	598	5,311	ATATA TAATTCACCGTGGATTCTTAGATTGTAGATATGATAA T <sub>05</sub>
554	598	386	ATATA TAATTCACCGTGGATTCTTAGATTGTAGATATGGTAA T <sub>15*</sub>

### J) NADH Dehydrogenase 9

5'	3'	Reads	ND9 gRNA Sequences
33	71	18	ATAT AACACATAAACGAAATGAGCATAGAAGTATATGTATGATG ATATTAT <sub>14</sub>
29	72	93	ATAT AAAACATAAACGGAATAAATATAAGAGTGATATGTGATATAAT T <sub>15</sub>
25	72	32	ATAT AAAACATAAACGGAATAAATATAAGAGTGATATGTGATATAATGTT T <sub>11</sub>
60	105	540	A TATAACACAAACAATAGAATAAAGTTAAGTGAGAATATGAGTGA TTTAT <sub>11</sub>
60	101	472	ATAT ATACAAACAATAAAATAGAATTAGATAGAGGTATAAGTGA TGTAAAT <sub>11</sub>
87	124	1	ATATCAAC AAACAAATAGAGCACTGTCTATGATATAAGTGTAA TTCT <sub>09</sub>
87	124	536	AAATATCAAC AAACAGACGAGACATCATCTACAGTATAAGTGTAA TAT <sub>10</sub>
87	124	453	AAATATCAAC AAACAAATAGAGCACTGTCTACGATATAAGTGTAA T <sub>05</sub>
117	160	1,145	ATAT ATAAAACAATAAAGAAATAGAAGGCTACAGTTAATGAGATAAT T <sub>13</sub>
130	167	137	AAAATTAACAAAACAATAAAGAAGCAGAAATTACAGT GAATGAAGTAAAT <sub>07</sub>
117	162	906	ATA TAACAAAACAATAAAGAAATGAGAAACTGTGATTGATAAGATAAT T <sub>12</sub> *
119	162	65	ATATA TAACAAAACAATAAAGAAGCAGAGAGTTACAGTTAATAAGATAA TTAT <sub>19</sub> *
116	162	55	ATA TAACAAAACAATAAAGAAATGAGAAACTGTGATTGATAAGATAATA T <sub>05</sub>
121	162	30	ATA TAACAAAACAATAAAGAAATGAGAAACTGTGATTGATAAGAT T <sub>11</sub>
149	193	1,196	ATAT AATAAAAACATACAATAAAATAGAAGGAAACTAATAAGATGATAG T <sub>15</sub>
147	187	5,984	ATATAT AACATACAATAAGACGAAGAGAAACTAATAGAGTGTAAAGA T <sub>15</sub> *
150	193	463	ATAT AATAAAAACATACAATAAGAATGAAAGGAAACTAATAAGATGATA TAT <sub>11</sub>
176	216	3,869	ATACAATAT AAAACAAAAGTCACAGATTAAGGGATAGAGATGTGATGAA T <sub>14</sub>
177	216	357	ATACAATAT AAAACAAAAGTCACAGATTAAGGGATAGAGATGTGATGAA T <sub>11</sub>
176	216	620	ATACAATAT AAAACAAAATCGTAGATTAAAGATAGAGATATGATGAA T <sub>14</sub>
176	221	63	ATAT ATATAAAAACAAAAGTCACGATTAGAAAGTAAAGATATGATGAA TTTT
201	242	87	ATATAA AATCAATAACAGATCATGATGATAGAGATAGAAGTTATAA T <sub>15</sub>
204	249	652	AAA AAAATCAATCAATGATAGATCACAGTGGTATAGGGATGAGAATT AT <sub>11</sub>
205	249	485	AAAA AAAATCAATCAATGATAGATCACAGTGGTATAGGGATGAGAATT T <sub>11</sub>
202	248	101	ATAT AAAATCAATCAATAGCAAGTTACGATAATATAGAAGTGAGAATTATA T <sub>13</sub>
203	248	94	ATAT AAAATCAATCAATAGCAAGTTACGATAATATAGAAGTGAGAATTAT T <sub>13</sub>
205	248	77	ATAT AAAATCAATCAATAGCAAGTTACGATAATATAGAAGTGAGAATT T <sub>10</sub> *
203	249	58	AAAAA AAAATCAATCAATGATAGATCACAGTGGTATAGGGATGAGAATTAT T <sub>05</sub>
239	279	569	AAATTAT ACAACATAAGCAGCGAAGATAAAGGTATAGAGATTAATT T <sub>12</sub>
239	268	561	AACATAAAT CGGTGAAGATAGAGACTATGAGAATTAAAT T <sub>07</sub>
240	268	346	ATAAAT CGGTGAAGATAGAGACTATGAGAATTAAAT ATGT <sub>15</sub>
238	279	86	AAAATTAT ACAACATAAGCAGCGAAGATAAAGGTATAGAGATTAATT T <sub>11</sub>
239	288	23	AT AAATATAACAACATAATAGAACGATGAAAGTGAAGATTATAAGAGTTAATT T <sub>12</sub>
239	286	26,097	ATATATAACACATAAGCAGTGAAGATAGAGATCATGAAAGTTAATT T <sub>11</sub> *
239	286	7,153	ATAC ATATACAACATAATAGAACGATGAAAGTGAAGATTATAAGAGTTAATT T <sub>11</sub>
238	286	1,828	ATATATAACACATAAGCAGTGAAGATAGAGATCATGAAAGTTAATT T <sub>11</sub>
243	286	1,812	ATATATAACACATAAGCAGTGAAGATAGAGATCATGAAAGTT T <sub>10</sub>
239	279	553	AAATTAT ACAATATAAGCAGCGAAGATAAAGGTATAGAGATTAATT T <sub>10</sub> *
242	286	499	ATAC ATATACAACATAATAGAACGATGAAAGTGAAGATTATAAGAGTTA T <sub>11</sub>
247	289	451	AAAATATACAACATAATGAAGCGACGAAAATAGAGACTATAGA TATTAT <sub>06</sub>
239	268	300	ATAAAT CGATGAAGATAGAGACTATGAGAATTAAAT T <sub>13</sub> *
246	286	262	ATATATAACACATAAGCAGTGAAGATAGAGATCATGAAA T <sub>11</sub>
243	286	236	ATAC ATATACAACATAATAGAACGATGAAAGTGAAGATTATAAGAGTT T <sub>17</sub>
236	282	197	ACAACAAATATAGAACGATGAAAGTGAAGATTATAAGAGTTAATT T <sub>13</sub>
272	312	43	ATAA GAACACACAGAGACAAGCAGAGTAAAGTGATAGTGACATA T <sub>06</sub>
273	314	5,869	ATATAT AGCAACACACAGAAATAAGCAAGGTAGAATATATGATGATAT T <sub>15</sub> *
271	314	250	ATAT ATGAACACACAGAAACAGATGAGATAGAGTATACAGTGATATAA T <sub>17</sub> *
275	314	157	ATATAT AGCAACACACAGAAATAAGCAAGGTAGAATATATGATGAT T <sub>12</sub>
272	314	93	ATAT ATGAACACACAGAAACAGATGAGATAGAGTATACAGTGATATA T <sub>14</sub>
303	339	44,686	ATAT ACAAAACAACACAAGATAGAGCACAAATGAATGCACAG TGATAAT <sub>13</sub>
298	343	2,681	ATAAACAAACAACACAGAGCAGAATACAAGTGAGTATATGAGAATA T <sub>11</sub>
298	343	1,741	GTAAACAAACAACACAGAGCAGAATACAAGTGAGTATATGAGAATA T <sub>13</sub>
297	339	270	ATAT ACAAAACAACACAAGATAGAGCACAAATGAATGCACAGGGATAA TAT <sub>11</sub>
304	339	139	ATAT ACAAAACAACACAAGATAGAGCACAAATGAATGCACAC TGATAAT <sub>13</sub>
305	339	6,792	ATAT ACAAAACAACACAAGATAGAGCACAAATGAATGCACTGAGATAAT T <sub>15</sub>
301	338	260	ATATAT TAAACAAACCGAACGAGCATAGATGGATATATGAGA TTCTTTT

5'	3'	Reads	ND9 gRNA Sequences cont.
319	366	36	AT ATTAAAACACAATCCGAAGAGTACAAATAGACGTATAAGATAAAAT T <sub>11</sub>
319	366	17	AT ATTAAAACACAATCCGAAGAGTATAAATAGACGTATAAGATAAAAT TAAT <sub>19</sub>
343	368	157	AT AAATTAAAACACAACTAAGAAATA GAGATAGACAGTATAAGATAAAAT <sub>08</sub> †
343	369	45	AAAATTAAAACACAACTAAGAAATA GAGATAGACAGTATAAGATAAAAT <sub>08</sub> †
352	392	10,040	ATATT AACGCATAACAGAAATGATTAGAACTGAGATATGATT AAT <sub>14</sub>
349	392	270	ATATAT AGACGCATAATAAAAGCAACTGAGGCTAGAATATGATTAA T <sub>19</sub>
349	392	232	TAT AACCGTATAACAAAAGCAGTTAACGCTAGAATGTGATTAA T <sub>14</sub>
352	392	5,765	ATATT AACGCATAACAGAGACAATTAGAACTGAGATATGATT AAT <sub>10</sub> *
349	392	2,338	ATATAT AACGCATAATAAGGGCAACTGAAACTAGAGTATGATTAA TTTCT <sub>12</sub>
351	392	691	ATATT AACGCATAACAGAGACAATTAGAACTGAGATATGATT T <sub>19</sub>
356	392	575	ATATT AACGCATAACAGAGACAATTAGAACTGAGATAT TTT
351	392	189	ATAT AACGCATAATAAGGGCAACTGAAACTAGAGTATGATT T <sub>16</sub>
358	392	167	ATATT AACGCATAACAGAGACAATTAGAACTGAGAT T <sub>11</sub>
382	421	1	AATCAAAATATTCTGTTCTGATGATAGAGATGTATAAT T <sub>10</sub>
380	418	317	ATATA TAAAACATTCTGTTCTAATGATAGAGATGTGATAA T <sub>12</sub>
409	447	761	ATATAAA ATTACCAACAGAATAGAAGCTAGACAAGTTAAGATATTT T <sub>08</sub>
409	447	667	ATATAAA ATTACCAACAGAATAGAAGCTAGACAAGTTAAGATATTC T <sub>12</sub>
410	447	82	ATATAAA ATTACCAACGAAATAAGACTAAGTGAATTAAGATGTT ATAT <sub>05</sub>
439	484	72	ATAT AACCAATCAACGAGTAAATGATGTAGAGTATTGTCAT T <sub>13</sub>
438	483	185	ATATAT AACCAATCAATAAGTGAGCGATGTAAGATGTTATTGTTAATA T <sub>10</sub>
437	472	130	ATATA TAACAAATAAACCGATGTGAGATATCATTACCACTGA TTTAAT <sub>08</sub>
438	483	46	ATATAT AACCAATCAATAAGTGAGCGGTGTAAGATGTTATTGTTAATA T <sub>14</sub>
439	483	34	ATATAT AACCAATCAATAAGTGAGCGATGTAAGATGTTATTGTTAAT T <sub>13</sub>
447	472	16	ATATA TAACAAATAAACGGTGTAGAGTGTCA GTAT <sub>16</sub>
442	483	25	ATATAT AACCAATCAATAAGTGAGCGATGTAAGATGTTATTGTT T <sub>12</sub>
468	514	234	ATAT ATAACACTTCAACCGGAGAGAGACCAATGAGAGATCAGTTAATA T <sub>13</sub>
469	514	2,106	ATAT ATAACACTTCAACCGAGAAGAACCCAGTGAGAAATCAGTTAAT T <sub>11</sub>
469	514	1,535	ATAT ATAACACTTCAATAGAAAGAAACTGACAGAGAAATTGATTAAT T <sub>13</sub>
472	514	59	ATAT ATAACACTTCAACCGAGAAGAACCCAGTGAGAAATCAGTT TTTAATTGTT <sub>14</sub>
502	549	746	ATATAAAATAACAATACAAGTGAATCAGATAATGGATGATGTTCAAT T <sub>13</sub>
509	542	1,407	ATAT ATAACAATACAATGAACGTAGTAATGGATGTA GTTGACAAT <sub>14</sub> *
497	539	935	ATAT ATAATACAAACAAACTGAGTGTGATGGATAGTGCTTAATGAGAA TAT <sub>14</sub> *
501	547	576	ATAGAATAACAATATAACGAACTAGATGATGGATAGTATTTGATA TAT <sub>16</sub> *
501	547	339	ATAGAATAACAATATAACGAACTAGATGATGGATAGTATTTGATA TATGCCACG
501	547	55	GTAGAATAACAATATAACGAACTAGATGATGGATAGTATTTGATA TATTTT
508	547	36	ATAGAATAACAATATAACGAACTAGATGATGGATAGT CT <sub>17</sub>
534	566	15	AAACGTACATATG ATCTCTCTGCTAGTATATAAGATGATAATA AT <sub>14</sub>
533	566	11	AACGTACATATG ATCTCTCTGCTAGTATATAAGATGATAAT T <sub>11</sub>
535	566	89	ATATG ATCTCTCTGCTAGTATATAAGATGATAAT TTCT <sub>16</sub>
533	566	72	AAATAAACGTACATATG ATCTCTCTGCTAGTATATAAGATGATAAT T <sub>08</sub>
534	566	24	ATG ATCTCTCTGCTAGTATATAAGATGATAATA AT <sub>15</sub>
535	578	13	ATATT AACGTACATACTGTCTCTATTGACATATAAGATGATAAT T <sub>13</sub>
543	580	39	ATAT TAAACGTACATATTATTTCTACTGATATACAAG TGATAAT <sub>09</sub>
568	604	76	ATATA TATGCAACACAGAGATAATTGTAGATGTATATGT GATATCGT <sub>11</sub>
569	600	78	ATATA TAACAACAAAAGTACATTTGTAACGGTATATG ATGAGTTT
582	612	66	TATT AATTGGTATGTGACAATAGAAGTGTATT T <sub>09</sub>
582	611	66	AATGCAGATAATT ATTGGTGTGCAATGATAGAGGTAATATT T <sub>12</sub> *
580	611	17	AATGCAGATAATT ATTGGTGTGCAATGATAGAGGTAATATTGT T <sub>12</sub> *
597	625	7	A AATGCAAATAGAAGTGGTGTGCAAT TATAGAGATGATAT <sub>09</sub>
609	644	274	ATATATAAGAATTACAATGGT GTATTAAGTGAAGTAATGTAATGAGAAATT T <sub>12</sub>
609	640	1	GATA GTTAGATAGAATAATGTAAGTGAAGAAAATT T <sub>12</sub>
615	659	1	ATATATAA AGGATTACAGTGGTGATATTGAATAGAGTGTGAAATG T <sub>12</sub>
618	658	1,673	ATATAT GAATTACAACGGTGATATTGAATAGAATAATGTGA TTGAAAT <sub>14</sub> *

### K) Ribosomal Protein Subunit 12

5'	3'	Reads	RPS12 gRNA Sequences
43	78	12,531	ATAT ACAACAACCATAATGAAGATCATGTACGTAGAAGA TTGATATAT <sub>14</sub>
35	76	2,218	ATA ACAACCGTACAGAACGTTACATATGCAGAGAACGGTGAGAT TTAT <sub>12</sub>
43	78	1,663	ATAT ACAACAACCATAATGAAGATCATGTACGTAGAAGA TTGATATAT <sub>12</sub>
38	78	423	ATAT ATAACAACCATAACAGAAATCGTATATGTGAGAGAAGTGA TTTCT <sub>15</sub>
63	109	5,122	AT ATAATATAAAACAAATAAGACAGAGTGTAGATAGTAATTGTATGA TAT <sub>12</sub>
66	99	32	AT ATAATATAAAACAAACGTAAATGATAACTGTG AGATGAT <sub>07</sub>
74	106	1,212	ATAT ATATAAAACAAATAAGAACGTAGATGAT TACTGTATAAT <sub>12</sub>
73	115	120	A TATATAATAACATAAAACAAATAAGAACGTAGATGATA TCTAT <sub>13</sub>
73	115	1,879	ATA TATATAATAACATAAGAACAGATAGAACGAGATGTAAATGATA T <sub>21</sub>
74	115	1,091	ATA TATATAATAACATAAGAACAGATAGAACGAGATGTAAATGAT T <sub>13</sub>
77	115	466	ATATA TATATAATAACATAAGAACAGATAGAACGAGATGTAAAT T <sub>05*</sub>
79	115	311	ATA TATATAATAACATAAGAACAGATAGAACGAGATGTAA T <sub>14</sub>
96	121	233	AATCA CGGATTATATAGTAACGTAAAATGA TATTAT <sub>12</sub>
92	121	90	AACTGGGC-ATCT CGGATTGTATAGTGTATAAAGTGAATTA TTTT
96	121	1,542	ATATAAACTCTGCAATCGA TGGACTTATATAGTGTAAAGTGA TAAT <sub>12†</sub>
96	121	1,066	ATATAGAACTAGGCAGTC CCGATTGTATAGTGTAAAGTGA TAT <sub>14†</sub>
93	121	793	ATATAGAACTGGCAATT CCGATTATATAGTGACATGAGATAGATA T <sub>15*†</sub>
94	121	219	ATATAGAACTGGCAATT CCGATTATATAGTGACATGAGATAGAT T <sub>10†</sub>
96	131	2	ATATAGAAATTA GGCAATCGGGATTATATAGTAACGTAAAATGA TAT <sub>10</sub>
119	158	3	ACT TACAATACACGTTGGTTATCGGAGTTAGGTGATTGTG ACTTAT <sub>10</sub>
139	170	44	ATATA ACGGCATATAAGTATACGTCGGTTACTAGGATTGTGTAAATTT
132	164	18	CATATAAA GGCATATAGTATACGTCGGTTACTGGATTGTGTAAAT <sub>07</sub>
133	158	56	ATACT TACAATACACGTTGGTTATCGGAGTT AGATGAT <sub>13</sub>
169	208	4,724	ATAAAAT ACAACGCAATATCCGAGTAAGGATTGTATAAGTGTGAGATAT AT <sub>12</sub>
164	195	146	ATAACGCAACA TCAGATGAGATTATATAAGTGTGAGATATG ATATAT <sub>11</sub>
158	207	67	ATAT ATAACGCAACATTGCAATGAGATTATGTAGATGAAATATGGTAT TAT <sub>05</sub>
164	201	896	ATA TAACATCCAACAAAGATTATATAGGTAGAGATG ATGTATAATTTAT <sub>22</sub>
156	207	192	ATAT ACAACGTAACATTGAGATAAGATTGTGTAGATAGAAATATGGTAT T <sub>12</sub>
158	207	104	ATAT ACAACGTAACATTGAGATAAGATTGTGTAGATAGAAATATGGTAT T <sub>06</sub>
198	235	4,950	ATATATAATGAC TAACTAAACTGATAAAGCAGTAGAAGAGATGTAAAT T <sub>11</sub>
194	235	3,025	ATATATAATGAC TAACTAAACTGATAAAGCAGTAGAAGAGATGTAAATTT T <sub>11</sub>
196	235	222	TAATGAC TAACTAAACTGATAAAGCAGTAGAAGAGATGTAAAT AT <sub>14</sub>
198	235	444	ATATATAATGAC TAACTAAACTGATAAAGCAGTAGAAGAGAGCAGTGTAAAT T <sub>12*</sub>
194	235	338	ATATATAATGAC TAACTAAACTGATAAAGCAGTAGAAGAGAGCAGTGTAAATTT T <sub>10</sub>
198	229	36	ATAACTT GACTAATAGAGTAGTGAGAGAGACAGTGTAAAT T <sub>08</sub>
200	235	25	ATATATAATGAC TAACTAAACTGATAAAGCAGTAGAAGAGAGCAGTGTAA T <sub>15</sub>
196	235	20	ATAATGAC TAACTAAACTGATAAAGCAGTAGAAGAGAGCAGTGTAAAT AT <sub>12</sub>

5'	3'	Reads	RPS12 gRNA Sequences cont.
203	246	341,382	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGATAAT AT <sub>14</sub>
203	246	1,324	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGATAAT AT <sub>15</sub>
205	246	1,091	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGATA T <sub>13</sub>
209	246	964	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGATA TAATAT <sub>12</sub>
206	246	922	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGAT T <sub>14</sub>
200	246	609	ATAT ATAATGACATAATTAGACTGATAAGATAACGAGAAAAGTGATGTA T <sub>12</sub>
203	246	544	ATATAT ATAATGACGTAACTGAGCTAATGAAGCAATGAGAGAGATAAT AT <sub>13</sub>
203	246	530	ATATAT ATAATGACGTAACTGAGCTAATGAGACAATGAGAGAGATAAT AT <sub>12</sub>
208	246	484	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAG TTAATAT <sub>11</sub>
204	246	430	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGATAAA AAT <sub>13</sub>
203	246	370	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGGGATAAT AT <sub>12</sub>
207	246	370	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGAGA AATAT <sub>13</sub>
203	246	312	ATATAT ATAATGACGTAACTGAACTAATGAGGCAATGAGAGAGATAAT AT <sub>12</sub>
210	246	282	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAG CGATAATAT <sub>14</sub>
211	246	268	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGA T <sub>12</sub>
203	250	267	AT ATAAATATGACATACTAGGTTAGTAAAGTGACGAAGAGATAAT ATTATTTT
203	246	196	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGATAAT AT <sub>13</sub>
203	246	193	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAACGAGAGAGATAAT AT <sub>21</sub>
205	246	177	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAGAGAGATAAT T <sub>15</sub>
212	246	176	ATATAT ATAATGACGTAACTGAGCTAATGAGGCAATGAG CGAGATAATAT <sub>14</sub>
203	246	173	ATATAT ATAATGACGTAACTGAGCTAATAAGGCAATGAGAGAGATAAT AT <sub>14</sub>
203	246	157	ATATAT ATAATGACGTAACTGAGCTAATGGGCAATGAGAGAGATAAT AT <sub>12</sub>
203	246	141	ATATAT ATAATGACGTAACTGAGCTAATGAGGCATGAGAGAGATAAT AT <sub>19</sub>
203	246	130	ATATAT ATAATGACGTAACTAAGCTAATGAGGCAATGAGAGAGATAAT AT <sub>13</sub>
203	246	119	ATATAT ATAATGACGTAACTGAGCTAATGAGGTAAATGAGAGAGATAAT AT <sub>12</sub>
203	246	117	ATATAT ATAATGACGTAACTGAGCCAATGAGGCAATGAGAGAGATAAT AT <sub>16</sub>
203	246	116	ATATAT ATAATGACGTAACTGGGCTAATGAGGCAATGAGAGAGATAAT AT <sub>13</sub>
203	245	56	ATAAAAT TAATGACATAACTAAATTGATAGGGTAAATGAGAGAGATAAT AT <sub>09</sub>
234	264	35	TAG TGCCTTCTATAGTAGATGATGATATA TGAT <sub>14</sub>
234	280	24	ATATA AGATCAACAAAATGCCATTTCGTAGTAGTGTGATGATATA T <sub>14</sub>
248	281	14	ATA TAAATCAACAGAACTGCCATTTGTAGTA TAGTGTGATATAAT <sub>13</sub>
234	280	434	ATATA AGATCAACAAAATGCCATTTCATAGTAGTGTGATGATATA T <sub>16</sub>
248	282	270	ATAT GTAAATCAACAGAACCGTCATTTGTAGTA TAGTGTGATATAAT <sub>12</sub>
288	322	16,349	ATA TACAATACGTGTATGATATTTTAACT AGGTAGATCAGTCAAATT T <sub>12</sub>
267	322	10	ATA TACAATACGTGTATGATATTTTAACTGGGTAGATCAGTCAAATT T <sub>07</sub>
269	308	3,731	ATA TATAATACCTTACATCGGGTAAATTGACGAGA ACATGAT <sub>11</sub>
288	322	909	ATA TACAATACGTGTAAATTTTAACT AGGTAGATCAATGAAAT <sub>15</sub>
309	349	128	AAATAT AACATATCTTATCTGAATCTAATTGTAATATGTG AAT <sub>16</sub>
309	336	195	AAATAAAACATATCTGAT TCTAAATCTAATTGTAATGTG AAT <sub>25</sub> †

\*Indicates that the tail sequence was shortened where random nucleotides after the poly U tail had been indicated.

†Indicates that the gRNA was identified under conditions of reduced stringency.

## APPENDIX D. Identified CR3 mRNA and gRNA transcripts.

A-C: Major CR3 mRNA and gRNA sequence classes. The CR3 mRNA transcriptome was generated using the TREU667 cell line. Identified sequences were then used to search gRNA transcriptomes from four different cell lines: EATRO 164 Bloodstream (BS), EATRO 164 procyclic (PC), TREU 927 procyclic and TREU 667 procyclic. ORF = previously identified Open Reading Frame (purple protein sequence). ARF = Newly identified Alternative Reading Frame (green protein sequence). Alternatively edited nucleotides are shown in Red. Inserted U-residues are lowercase while deleted U-residues are shown as asterisks. Canonical Watson-Crick base pairs (||); G:U base pairs (:). Previously identified start codons are doubled underlined. Potential upstream AUG start codons are indicated by wave underlines. gRNAs were sorted based on guiding sequence class. Sequence variations observed in the 3'-U-tail were ignored in assigning class. Transcript copy number (Reads), were determined by adding all gRNAs of the same sequence class. Only major sequence classes are shown (defined as containing greater than 100 transcript copies). In the case of rare transcripts, the identified gRNA are shown regardless of copy number. gRNA transcript numbers varied greatly between the different cell lines. Interestingly, the most abundant mRNA (CR3 Form C, 7147 reads) had the fewest identified gRNA reads.

C. CR3 Form C	mRNA Sequence	Reads
AUGUGUAUGAUUAUAAAAACA--A <u>GuGuA*****UGuuGuuGuuuuGuuuuG***AuuuuGGuuGuACAUuuuuuuuG</u>		7147
AUGUGUAUGAUUAUAAAAACA--A-G <u>GuGuA*****UGuuGuuGuuuuGuuuuG***AuuuuGGuuGuACAUuuuuuuuG</u>		889
AUGUGUAUGAUUAUAAAAACA <u>uuA-G<u>GuGuA*****UGuuGuuGuuuuGuuuuG***AuuuuGGuuGuACAUuuuuuuuG</u></u>		565*
AUGUGUAUGAUUAUAAAAACA-u <u>GuGuA*****UGuuGuuGuuuuGuuuuG***AuuuuGGuuGuACAUuuuuuuuG</u>		505
Fully Edited Form		Reading Frame
M C M I Y K N N V Y V V V L F W F W L Y I F F V		
AUGUGUAUGAUUAUAAAAACA <u>GuGuA*****UGuuGuuGuuuuGuuuuG***AuuuuGGuuGuACAUuuuuuuuG</u>  :  :    :  :  :  :  :  :  :  :  :  #      ,UUUAUUAU---AUAGUGAUAAAGAUGGAAU--UGAAGCCGACACGUAAAUAUUAUA		ARF +1
Cell Line	gRNA Sequence	Reads
EATRO 164 PC	ATAATAAAATGCACAACCTAGAATTGAAGTAAAGTGATGATATATAT <sub>N</sub>	122
EATRO 164 PC	AAAATGCACAACCTAGAATTGAAGTAAAGTGATGATATATAT <sub>N</sub>	3
EATRO 164 PC	ATATATAAAATGTACAACCAGAACATTAAGATAAAAGTGATGATGTATAT <sub>N</sub>	1
EATRO 164 BS	ATAATAAAATGCACAACCTAGAATTGAAGTAAATGATGATATATAT <sub>N</sub>	453
EATRO 164 BS	ATATTAAAAATGCACAACCTAGAATTGAAATAAAAGTGATGGTATATAT <sub>N</sub>	224
EATRO 164 BS	ATAATAAAATGCACAACCTAGAATTGAAGTAAAATGATGATATATAT <sub>N</sub>	2
TREU 927 PC	ATAATAAAATGCACAACCTAGAATTGAAATAAAAGTGATGGTATATAT <sub>N</sub>	166
TREU 927 PC	ATATAATTAAATGCACAGCGAACGTTAAGGTAGAATAGTGATATATAT <sub>N</sub>	117
TREU 927 PC	ATATATAAAATGTACAACCAGAACATTAAGATAAAAGTGATGATGTATAT <sub>N</sub>	86
TREU 667 PC	ATAATAAAATGCACAACCTAGAATTGAAATAAAAGTGATGGTATATAT <sub>N</sub>	53
TREU 667 PC	ATATAATTAAATGCACAGCGAACGTTAAGGTAGAATAGTGATATATAT <sub>N</sub>	15
TREU 667 PC	ATATATAAAATGTACAACCAGAACATTAAGATAAAAGTGATGATGTATAT <sub>N</sub>	14
TREU 667 PC	TATAATTAAATGCACAGCGAACGTTAAGGTAGAATAGTGATATATAT <sub>N</sub>	5

\* no gRNAs identified.

## APPENDIX E. ND7 5'-most gRNA populations and the predicted mRNA sequences generated.

A-J: ND7 gRNA major classes and predicted editing patterns. ND7 terminal (5' most) gRNA populations and the predicted mRNA sequence generated. Predicted sequences presented are based on the most abundant gRNAs that generate each reading frame found in the four gRNA transcriptome databases. Initial characterization of the ND7 transcript was done using the EATRO 164 cell line and is unusual in that it is edited in two distinct domains [20]. While the 5' domain was edited in both life cycle stages, complete editing of the 3' domain was only detected in bloodstream stage parasites. Interestingly, the most abundant EATRO 164 PC (procyclic or insect form) gRNA would generate a sequence that brings the 5' most AUG into a +2 frame. The ARF is 65 AA long and involves the entire 5' editing domain. In contrast, the most abundant gRNAs in the EATRO 164 Bloodstream stage library (EATRO 164 BS), would generate sequences that use the originally described ND7 ORF). While gRNA transcript numbers again varied greatly between the different cell lines, all three cell lines had gRNA sequence variants that allowed access to both reading frames.

A. ND7 Form A		Predicted mRNA Sequence	Reading Frame
M T T W ST	M I S I I L C Y F W ST M L F L V V F L H L Y R F T F G P Q AUGACUACAUGAUAA <u>GU</u> AuCAuuuu <u>Gu</u> GuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG :    :   :   :   :   :   :   :   :   :   :   :   :   :   : ,UAUAGUAAGAUGCAA <u>U</u> GAAGAACCGCUCAAGAGAA <u>U</u> GUAAACAUUAAAA		ORF
Cell Line	gRNA Sequence	Reads	
EATRO 164 BS	AAATATACAAATGTAAGAGAACTGCCAAAGTAACGTAGAATGATAT <sub>N</sub>	20487	
TREU 927 PC	AAATATACAAATGTAAGAAA <u>ACT</u> ATCGAGAGTGTAGAATGATAT <sub>N</sub>	10537	
TREU 667 PC	AAATATACAAATGTAAGAAA <u>ACT</u> ATCGAGAGTGTAGAATGATAT <sub>N</sub>	787	

B. ND7 Form B		Predicted mRNA Sequence	Reading Frame
M T T W Y S I I L C Y F W ST	M I ST M L F L V V F L H L Y R F T F G P Q AUGACUACAUGAUAA <u>GU</u> AuCAuuuu <u>Gu</u> GuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG :    :   :   :   :   :   :   :   :   :   :   :   : ,UAUUAUAGUAAGAUGCAA <u>U</u> GAAGAACCGCUCAAGAGAA <u>U</u> GUAAACAUUAAAA		ORF
Cell Line	gRNA Sequence	Reads	
EATRO 164 BS	AAATATACAAATGTAAGAGAACTGCCAAAGTAACGTAGAATGATATTAT <sub>N</sub>	35079	
TREU 927 PC	AAATATACAAATGTAAGAAA <u>ACT</u> ATCGAGAGTGTAGAATGATATTAT <sub>N</sub>	38432	
TREU 667 PC	AAATATACAAATGTAAGAAA <u>ACT</u> ATCGAGAGTGTAGAATGATATTAT <sub>N</sub>	4365	

C. ND7 Form C		Predicted mRNA Sequence	Reading Frame
M T T W ST	M I S T F M L F L V V F L H L Y R F T F G P Q AUGACUACAUGAUAA <u>GU</u> Acuuuu <u>Gu</u> GuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG :    :   :   :   :   :   :   :   :   :   : ,UAUUAUAGUAGAAGAUUCGGAGAA <u>U</u> GUAAACAUAGCAUUAACA		ORF
Cell Line	gRNA Sequence	Reads	
TREU 927 PC	ACATATACGATA <u>CAA</u> ATGTAAGAGGCTGTAGAAGTGTAAAT <sub>N</sub>	75654	

D. ND7 Form D		Predicted mRNA Sequence	Reading Frame
M T T W ST	M I I V S F M L F L V V F L H L Y R F T F G P Q AUGACUACAUGAUAA <u>GU</u> Acuuuu <u>Gu</u> GuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG :    :   :   :   :   :   :   :   :   : ,UAUUAUAGUGAA <u>U</u> AAUUGGAGACUAUCGAAGAA <u>U</u> GUAAACAUUAAAA		ORF
Cell Line	gRNA Sequence	Reads	
EATRO 164 PC	AAATATACAAATGTAAGAGCTATCAGAGGTAATATAAGTGTATAAT <sub>N</sub>	240	
TREU 927 PC	ATATACACAAATGTAAGAGACTATCGAGAGTGTACATAAGTGTATAAT <sub>N</sub>	477	
TREU 667 PC	AAATATACAAATGTAAGAGCTATCAGAGGTAATATAAGTGTATAAT <sub>N</sub>	1152	

E. ND7 Form E	Predicted mRNA Sequence	Reading Frame
M T T W ST M I M T F F M L F L V V F L H L Y R F T F G P Q	AUGACUACAUAGUA <u>AuG*ACAUuuuuuAuGuuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG    :   #  ::   ::   :   :   :   :   :   :   :   :   : _UAU-U<u>AUAGGAAUACGGAGAGUUAU<u>CAGAGAA<u>UGUA<u>AAUAUAUA</u></u></u></u></u>	ORF
Cell Line	gRNA Sequence	Reads
EATRO 164 PC	ATATAATAATGTAAAGAGACTATTGAGAGTGGCATAGGGATATTAT <sub>N</sub>	765
F. ND7 Form F	Predicted mRNA Sequence	Reading Frame
M T T W ST M I S T F M L F L V V F L H L Y R F T F G P Q	AUGACUACAUAGUA <u>ACAUuu<u>AuGuuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG    :   :   :   :   :   :   :   :   :   : _UA<u>AUGUAGUGAA<u>ACGGGAGAGUUACAGAGAA<u>UGUA<u>AAACAUAGCAUA<u>ACA</u></u></u></u></u></u></u>	ORF
Cell Line	gRNA Sequence	Reads
TREU 667 PC	ACATATACGATA <u>CAAATGTAA<u>AGAGGCTGTTAGAGTGTAAAT<sub>N</sub></u></u>	6623
G. ND7 Form G	Predicted mRNA Sequence	Reading Frame
M T T W Y S I I Y V I F G S F F T F V S F Y I W S T A	AUGACUACAUAGUA <u>uAGUA<u>CAuu<u>AuGuuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG    :   :   :   :   :   :   :   : _UA<u>UUAUAGUGAA<u>UACGGGAGAGUUACAGAGAA<u>UGUA<u>AAUAUAUA</u></u></u></u></u></u></u>	ARF +2
Cell Line	gRNA Sequence	Reads
EATRO 164 PC	ATATAATAATGTAAAGAGACTATTGAGAGTGGCATAGTGATATTAT <sub>N</sub>	100761
EATRO 164 BS	ATATAATAATGTAAAGAGACTATTGAGAGTGGCATAGTGATATTAT <sub>N</sub>	354
H. ND7 Form H	Predicted mRNA Sequence	Reading Frame
M T T W ST M I S T F Y V I F G S F F T F V S F Y I W S T A	AUGACUACAUAGUA <u>ACauuu<u>AuGuuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG    :   :   :   :   :   :   : _UA<u>AGAUGCAA<u>AGCCGU<u>CAAGAGAA<u>UGUA<u>AAACAUUA<u>AAA</u></u></u></u></u></u></u></u>	ARF +2
Cell Line	gRNA Sequence	Reads
EATRO 164 BS	AAATATACAA <u>ATGTAA<u>AGAGAA<u>CTGCCAA<u>AGTAACGTA<u>ATN</u></u></u></u></u>	402
I. ND7 Form I	Predicted mRNA Sequence	Reading Frame
M T T W ST M I S T M L F L V V F T F V S F Y I W S T A	AUGACUACAUAGUA <u>ACAU<u>GuuAuuuuuGGuAGuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG    :   :   :   :   :   : _UA<u>AGUA<u>AGAGCU<u>AUUGA<u>AGUGAG<u>GUUA<u>AGUA<u>AAA<u>UGUA<u>AAA<u>UA</u></u></u></u></u></u></u></u></u></u></u></u>	ARF +2
Cell Line	gRNA Sequence	Reads
TREU 667 PC	ATATAAA <u>ATGTAA<u>ATGATATGAGTGTAGAGTTACTGAGA<u>ATGATATN</u></u></u>	12929
J. ND7 Form J	Predicted mRNA Sequence	Reading Frame
M T T W ST M I S T F I V I F G S F F T F V S F Y I W S T A	AUGACUACAUAGUA <u>ACAU<u>uu<u>AuGuuAuuuuuGGuAGuuuuuuuACauuuGuAuCGuuuuACauuuG*GUCCACAG    :   :   :   :   :   : _UA<u>AAUAGUAGUGA<u>AGAU<u>UGUCGG<u>AGAA<u>UGUA<u>AAACAUAGCAUA<u>ACA</u></u></u></u></u></u></u></u></u></u>	ARF +2
Cell Line	gRNA Sequence	Reads
TREU 927 PC	ACATATACGATA <u>CAAATGTAA<u>AGAGGCTGTTAGAGTGTAAAT<sub>N</sub></u></u>	251

## APPENDIX F. RPS12 5'-most gRNA populations and the predicted mRNA sequences generated.

A-E: RPS12 gRNA major classes and predicted editing patterns. RPS12 terminal (5' most) gRNA populations and the predicted mRNA sequence generated. RPS12 differs from both CR3 and ND7 in that the alternative edit that shifts the reading frame occurs just downstream of the previously identified start codon (double-underlined). We do note that the identified alternative gRNAs are rare in all of the gRNA libraries except TREU 667.

A. RPS12 Form A		Predicted mRNA Sequence	Reading Frame
		M W F L Y G C C C L R F V L F V CAAACUAAGUA <u>AAuAu<u>GuGuGA*</u>UUUUUGUAUG*GuuGuuGuuAC*GuuuuGuuu<u>AuuuGu</u> <u>UAUAUAGGUAGAAGAUGCAUGUACU-AGAAGUAUAC-CAACAA<u>CAUAUA</u></u></u>	ORF
Cell Line	gRNA Sequence	Reads	
EATRO 164 PC	ATATACAACAACCATA <u>TGAAGATCATGTACGTAGAAGATTGATATAT</u> N	12531	
EATRO 164 BS	ATATACAACAACCATA <u>TGAAGATCATGTACGTAGAAGATTGATATAT</u> N	1663	
TREU 927 PC	ATATACAACAACCATA <u>TGAAGATCATGTACGTAGAAGATTGATATAT</u> N	505	
TREU 667 PC	ATATACAACAACCATA <u>TGAAGATCGTGTACGTAGAAGATTGATATAT</u> N	936	
B. RPS12 Form B		Predicted mRNA Sequence	Reading Frame
		M W F L Y G C C C L R F V L F V CAAACUAAGUA <u>AA<u>AAAuuuuGuuuuuuuuGCG<u>GuGuGA*</u>UUUUUGUAUG*GuuGuuGuuAC*GuuuuGuuu<u>AuuuGu</u> <u>UAUUJAGAGUGGAAGAGACGUUA<u>CAUJ-GAAGACAU<u>G</u>-CAACAAUA</u></u></u></u>	ORF
Cell Line	gRNA Sequence	Reads	
EATRO 164 PC	ATAAACAA <u>CCGTACAGAAGTTACATATGCAGAGAAGGTGAGATTAT</u> N	2218	
TREU 927 PC	ATAAACAA <u>CCCATACAGAAGTTACATATGCAGAGAAGGTGAGATTAT</u> N	300	
TREU 667 PC	ATAAACAA <u>CCGTACAGAAGTTACATATGCAGAGAAGGTGAGATTAT</u> N	3834	
C. RPS12 Form C		Predicted mRNA Sequence	Reading Frame
		M W F C M V V V Y V L F Y L F CAAACUAAGUA <u>AAAAG<u>uuuuuuuuuuuGCG<u>GuGuGA*</u>UUUUUGUAUG*GuuGuuGuuAC*GuuuuGuuu<u>AuuuGu</u> <u>UUUAGAGAGAA<u>AGUGCAUA<u>ACU-AAGACAUAC-CAAA<u>UAUA</u></u></u></u></u></u>	ARF +1
Cell Line	gRNA Sequence	Reads	
EATRO 164 BS	ATATATAACCATA <u>CAGAACATACGTGAAAGAGAGAGAT</u> N	144	
D. RPS12 Form D		Predicted mRNA Sequence	Reading Frame
		M L F F F R M W F C M V V V Y V L F Y L F CAAACUAAGUA <u>uu<u>AuAu<u>GuGuGuuuuuuGCG<u>GuGuGA*</u>UUUUUGUAUG*GuuGuuGuuAC*GuuuuGuuu<u>AuuuGu</u> <u>UGAUUA<u>UUA<u>AGUA<u>AGAGAGAGAGUA<u>UAUGCU-AAAACAUAC-CAACAGAU<u>AUA</u></u></u></u></u></u></u></u></u>	ARF +1
Cell Line	gRNA Sequence	Reads	
TREU 927 PC	ATATAGAACCA <u>CCATACAGAACATACGTATATGCGAGAGAA<u>ATGATATTATAT</u>N</u>	22	
E. RPS12 Form E		Predicted mRNA Sequence	Reading Frame
		M W F C M V V V Y V L F Y L F CAAACUAAGUA <u>uu<u>AAAuuuuGuuuuuuuuGCG<u>GuGuGA*</u>UUUUUGUAUG*GuuGuuGuuAC*GuuuuGuuu<u>AuuuGu</u> <u>AAA<u>UUUAGAGAGAA<u>AGUGCAUA<u>ACU-AAGACAUAC-CAAA<u>UAUA</u></u></u></u></u></u></u>	ARF +1
Cell Line	gRNA Sequence	Reads	
TREU 667 PC	ATATATAACCATA <u>CAGAACATACGTGAAAGAGATGAGATTAA</u> N	2664	

**APPENDIX G. Alignments of *T. brucei* and *T. vivax* edited mRNAs ATPase 6 (A), COIII (B), CR3 (C), CR4 (D), ND3 (E), ND7 (F), ND8 (G), ND9 (H), and RPS12 (I).**

Uppercase letters indicate nucleotides originally encoded in the DNA, lower case u's indicate uridines inserted during editing and asterisks indicate uridines removed during editing.

**A. ATPase 6 - Pan-edited non-dual coding**

```

AAAAAAUAGUAUUUUGAUAAAAGUAAAaAuGuuuuuAuuuuuuuuuuGuGAuuuA T. brucei
-----AuGuuuuuGuuuuuuuuuuuGuGAuuuG T. vivax

UUUUG-GuuGCGuuuGuuA---uuAuGuAuGuAuuAuuGuGuAuGAuCuAGGuuAuGuu T. brucei
UUUUG*GuuGCGuuuGuuA***UUaUGuGuGuAuuAuuGuGuGuGAuCuAGGuuAuGuu T. vivax
uuAuuGuGuAuuuuAA---uGUuuuAAuGuuGAuuuuuG-AuuuuuuAuuAuuuuGuuuG T. brucei
uuGuuGuGuAuuuuAA***UUGuuuGAuGuuAAuuuuuG*AuuuuuuGuuGuuuGuuuG T. vivax
*UUUGAuuuGuAuuuGuuuGuuGGuuuGuG***UUUGuuuuuAuuGuuGuGGuuuAuGuu T. brucei
-uuuGAuuuGuAuuuGuuuAuuGGuuuAuG---uuuAuuuuuGuuAuuGuGGuuuAuGuu T. vivax
GuuuA---AuuuAuAuAGuuuAAUUUUGuAuuA*UUGuAuuAC---uAUUUG***AA T. brucei
GuuuA****AuuuGuAuAGuuuGAUUUuGuAuuA*UUGuAuuACC****UAuuuG--*AA T. vivax
uuuG*UAuuUGuuGuuuuGuAuuGuuuuuuuAuuGuA----uAuuG-CAuuuuuAuuuu T. brucei
uuuG-uAUUUGuuGUUUuGuAuuGuuuuuuuAuuGuA****UA-uGuCAuuuuuGuuuu T. vivax
uGuuuuGuuuuuuA-uGuGAuuuuuuuuGuuuAAuAAuuuGuUA-GuuGGuGAuA**** T. brucei
uGuuuuGuuuuuuGuuG-GAuuuuuuuuGUUUAAuAGuuGuuGUG-uGGuGAuA--- T. vivax
GuuuuAuGGAuG---uuuuuuuAUUC**GuuuuuuGuuGuGuuuuuuAGAGuGuuuuuCu T. brucei
GuuuuAuGGAuGuuuuuuuuuG-*C--GuuuuuuGuuGuGuuuuuuAGAAuGuuuuuCu T. vivax
uuGuuGuGuC---GuuGuuuGuCGACGuuuuuuGCGuuuG-UUUUGuAAuuuAuuAuCAu T. brucei
uuGuuAuGUC****GuuGuuuGuCACAAuuuuuACGuuuG*UUUUGuAAuuuAuuGuCAu T. vivax
CCCAuUUUUUuGuuGAuGuuuuuuG-A-uuuuuuuUAuuuuA-uuuuuGuuuuuuuuu T. brucei
CCCAUUUUuGuuAAuGuuuuuuGuAuuuuuuuUAuuuuAuuuuuG-uuuuuuuu T. vivax
uuuA-----uG---GuGuuuuuuG-uuA-uuGAuuuAuuuuAuuuAuuuuGuG- T. brucei
uuuGuuuuuuuuGuuuuGuG---uAuuuuAuuuGGuuuAuuuGuuuG--uuuGUGu T. vivax
-uuuuGuuuuuGuuuuAuuAuuuuuAU--G-uGuuuuuAuAuUUGuuuGGGuuAUuGCC-- T. brucei
uuuuGuuuuuGuuuGuuGuuuAuuGuG---uuGuAuuuAuuGGGuuuAuuuGCC** T. vivax
***GC-CA-uAuuAC****AGuuuAuuuAuuuuuuGuAAuAuGAuuuuGCAGuuGAuAAuG T. brucei
***GC**GuuAuuAC---AGuuGuuuAuuuuuuGuAAuAuGAuuuuGCAuuGGuAAuG T. vivax
--G**AuuuuuuGuuGuuuuuGuuG-uuuGuuuAGuuuuGuAuuuGAuuuuuGAuAGuuA T. brucei
**G**AuuuuuuAuuGuuuuuGuuGuuuuG-uuAG----- T. vivax
uuAuAuuGuuGuuGAAuuuG**GuuGUuuA**UUGGAGUUUAAGAAUAAGAUCAAUA T. brucei
----- T. vivax

GUUAUAAUA T. brucei
----- T. vivax

```

## B. COIII - Pan-edited non-dual coding

GGUUAUUGAGGAUUGUUUAAAUAuuAuuAuuuuuuuAuGuuuuuGuuuC\*\*\*\* T. brucei  
 ----- T. vivax  
  
 \*GuuGuAuAuuuGuuGGuGuuA-\*\*\*\*GuGGuGuuuuuGuuuuuuuuAuCuuuACCuGCCA T. brucei  
 ----- T. vivax  
  
 uuGuuAuuGuGuAuuGGuuAuuuuGuuG\*\*\*\*GAuuuAuuuGuuAuuGUUUG\*\* T. brucei  
 ----- T. vivax  
  
 \*\*GuAGuuuuuuAuuuGuuGAuuGuG\*\*\*\*GuuuuAuuuuuuuuuuuGuuGGuuuuGuA T. brucei  
 ----- T. vivax  
  
 uuuGuuuGuuGuuGuuAuuGuuAGAuuuGuuuuGuGAuuuuuuACGuGGuuuAuuGAuu T. brucei  
 ----- T. vivax  
  
 uuuGuGuuuuAuuACGuuGuAuCCAGuAuuGuuuuuuAuGGuuuuuAuGuAG\*UGAGuuu T. brucei  
 ----- T. vivax  
  
 GuuuuAuuuAuGGCGuuuuuuG\*\*UUGuAuuAuuuGGuuuAuGuuuAuuuuuGuGuuGuG T. brucei  
 ----- T. vivax  
  
 AGuuuGCUUUCGuuuuuuGuuuACCuuAuAuGuuuGuuGuuAuuGuGAuuAuGGuu T. brucei  
 ----- T. vivax  
  
 uuGuuuuuuAuuGG\*UAuuuuuuAGAuuuAuuuAAuuuGuuGAuAAuACAuuuuAUUUG T. brucei  
 ----- T. vivax  
  
 uuUGuuAGuGGuuuAuuuGuuAAuuuuuuuGuuuuGuGUUUUJGGuuuAGGuuuuuuGu T. brucei  
 ----- T. vivax  
  
 uG\*\*UUGuuGuuuuGuAuuAuGAuuGAGuuuGuuGuuuG\*\*\*\*G--uuuuuuGuuuuuG- T. brucei  
 -----uuG-uuG--uuGuuuuuuuGuuuuuGu T. vivax  
  
 uGAAACCA--GuuA---UGAGA\*\*GUUGCAuuGuuAuuuAuuACAuuAAGuuGuGG\*\*\* T. brucei  
 uGAA-uCA\*\*GuuG\*\*UGGGA\*\*AuuuACGuuGuuAuuuAuuACGuuGAGuuGuGG--- T. vivax  
  
 \*UG-uuuuuGGuuCuAuuuuAuuuuuAuuG---GAuuuAuUACAuuuuA\*\*UGCAuGuuu T. brucei  
 -uG\*UUUUUGGuuCuAuuuuAuuuuuAuuG\*\*GAuuuGuuGCAuuuuA--uGCAuGuuu T. vivax  
  
 uuuuAGGuGuuuuGuuGuuG-uuuAuuuG-uuuuAuG--CGuuuGuuuAAuuuuuuGuGu T. brucei  
 uuuuAGGuGuuuuAuuGuuGuuuuA--uGuuuUUAuG\*\*CGuuuGuuuAGuuuuuuAuGu T. vivax  
  
 AuGGGuACACGuuuuGuuuuuuuGuAuuGuGuuuGuuuAuAuuGACAuuuuGuuGA-UUU T. brucei  
 AuGGGuACACGuuuuGuuuuuuuGuAuGuuGuuuGuAuuGACAuuuuGuuGAUUU\* T. vivax  
  
 AGuuuGAuuuuuuuuAuuGCGAuuuGuuuAuuuuGAuGuuuuAuG---uGuuAuGuAuu T. brucei  
 GGuuuGGuuuuuuuuGuuGCGAuuuGuuuAuuuuGAuGuuuuAuG\*\*\*UGuuAuGuAuu T. vivax  
  
 uGuGuGuGuAuuuuAuuGGuGuuuuUUUAGUUGuuGAuuA\*GuuAAuuuGuAuuGGUAG T. brucei  
 uGuGuGuGuAG----- T. vivax  
  
 UUUGUAGGAAG-- T. brucei  
 ----- T. vivax

### C. CR3 - Pan-edited dual coding

AGAAAUAUAAAUAUGUGUAUGAUAAUAAAACAAuGuuuGA\*\*\*\*UUGuuuGGuuuuG- T. brucei  
-----AuGuuuGA---\*UUGuuuAGuuuuGU T. vivax  
  
--uuG-uuuuUUUAuuGuuuGuuuGuACauuuuuuuGuuuuuuAuuuGuuuGuG---- T. brucei  
U\*\*\*GuuuuuuuAuuG-uuGuuuGuACauuuuuuuGuUUUUGuuuAuuuGuG\*\*\*\* T. vivax  
  
\*\*\*A\*\*UUUGuuuuuAuGuuuGuuA-\*UUUA----GuuuuuGuuuuuuAuuGGGuuuu T. brucei  
\*\*\*A--uuuGuuuuuAuGuuuGuuG--uuuGuuuuuuG--uuUA---uuuG-uGGGuuuu T. vivax  
  
uGuuuuuuAuuuAA---uA--uGGGuuuA---uUGuuGuGuuuAuuuuuuuuuuuuuAuu T. brucei  
uGuuuuuuGuuuAA\*\*\*\*UA\*\*UGGGuuuG\*\*\*UUGuuGUGuuuAuuuuuuuuuuuGuu T. vivax  
  
uuAuCAuuuGAuAuGuuGuuAuCAuuuUAuuAuuGuAuAuAAGuUUUCGUUAUUAGAUU T. brucei  
uuAuCAuuuGAuAuGuuAuuAuCGuUUUGuuAuuAuAuAAGuUUUCGUUAUUAA--- T. vivax  
  
AAAAAAAGUAUGCAAAUAAAUUUUUGU T. brucei  
----- T. vivax

### D. CR4 - Pan-edited dual coding

UAAUUUAUUGUUUAUCUUUGUGUAUUUAuuAuuuuAuuuAA---uuuuG---GuuGu T. brucei  
-----AUUuuuuuGuuuGuuGu T. vivax  
  
GC---\*\*\*AuuuuuuuuuuuuuuuuAuuuG\*\*\*GuG\*UGuuuGuGuuuuA\*UGuA\*C\*A--- T. brucei  
GC\*\*\*\*AuuuuuuuuuuuuuuAuuuG---GuG-uGuuuGuGuuuuA-uGuA--uA\*\*\*\* T. vivax  
  
\*GuuuAuGGuAuAuuuuAuuGuuGuuuGuuuuuGuuuuuGuuGUUUG-uuUGuG--uG T. brucei  
\*GuuuGuGGuAuAuuuuGUuGuuGuuuuuuuuuGuuuuuuuGuuG\*\*\*UGuGuuuG T. vivax  
  
GGuA--uGuuuuAuuuGuuuGuuAuA---GuuGuuuGuuuuuuuuuGuuGuUUUG\*GG T. brucei  
GGuA\*\*UGuuuuAuuuGuuuGuuAuA\*\*\*GuuGuuuGuuuuuuuuuGuuuuG-GG T. vivax  
  
uuGuG---AuuuuuuAuuG\*\*G--uGuuuuG\*\*\*AuuGuAuA---GuuuAuuuuuuuG T. brucei  
uuGuG---AUUUUuGuuA--GuuuA--UuG---AuuGuAuA\*\*\*GUuuGuuuuuuuuG T. vivax  
  
uGAC-GuuAuAAuuUUGuuuAuuuuuuuuuuuAuuuuGuuuuGuGuuuuuG---uAuu T. brucei  
uGAC\*GuuAuAAUuuuGuuuAUUUUuuuuuuuGUuuuGUuuGuGuuuuuGuuuuGuu T. vivax  
  
G\*UUGuuuuuA-uUUGGuuuGuuuGGuuuuuuuuuG\*\*\*UAuuuuuGUUGuGuuuuGuG T. brucei  
GuuuuuuuuG---GUUU\*GuuuGGuuuuuuuuuG---uAuuuuuuGuuAuGUuuuGuG T. vivax  
  
uuAuuuuuuGAuuuAuuuuuuAuGuUGuuuuuGUuuuG--GG\*\*\*UG\*G-uuuuuuuGu T. brucei  
uuAuuuuuuGAuuuGuuuuuuAuGuuGuuuuuG---uAuuGG---GUGGuuuuuuuGu T. vivax  
  
uuuuGuuuuuuuuuuuuGuuuAuGuuuGuuuuA---uuuGuGGuuGuuG--uuAuuuG T. brucei  
uuuuGuuuuuuuuuuuuGuuuAuGuuuA---uuGuuuuuGuAGuuAUUG\*\*UUGuuuA T. vivax  
  
uuAGuuuGGuuGuuGUUG-uuAuuUGuG---uA--uA\*\*\*\*GGUUUAuuUAuA\*UGCGuu T. brucei  
uuAGuuuGGuuGuuGuuGUU\*AuuuGuGU\*\*\*AUuuA--UUGGuuuAuuuAuA-uGCG-- T. vivax  
  
uuuuAuuuuA-----GAuAAuUAuG\*\*\*\*G\*\*\*\*UA\*\*UUGGUUUUAAAAUGUUUUU T. brucei  
uuuuG---uGUUUU\*\*AA----- T. vivax  
  
UCU T. brucei  
--- T. vivax

#### E. ND3 - Pan-edited dual coding

UCAAAAAAAUCCUCGCCUUUUACUUUA-GUUUGUUAUCAuA--uuuuuAuAuuuGuuuu T. brucei  
-----AUG---UUAuCA--AUuuuuuGuAuuuG---- T. vivax  
  
UG---\*A\*UA-uuGuGGuuuA\*\*UUAuuuuAuuuA-uAGG--uuuuuuuuuAuGuuuuuu T. brucei  
-GUUUuG-uGuuuGuGGuuuA--UuAuuuuA-uuGuuAGGuuuuuuCU\*\*AuGuuuuuu T. vivax  
  
AuGuuuuuuAuuGCAuuuuuuuG---AuuGuuuuCGuuGuuGuuGuGGuuuuuCGuGuG T. brucei  
AuGuuuuuuGuuACAuuuuuuuG\*\*\*AuuGuuuuCGuuGuuGuuAuGAuuuCAuGuG T. vivax  
  
---GuUUGuA---uGAuAuG--A----AuUCACGuuuG\*GUGuuuuuA--uACAuGGAu T. brucei  
\*\*\*GuuuGuA\*\*\*UGAuAuG\*\*A\*\*\*\*\*AuuCACGuuuG-GuGuuuuA\*\*UACAuuAGAu T. vivax  
  
uuAUGuuuuGuuA-----GuUGuUUGuuuuuuGuAuuGuu---AAAuuCCAu T. brucei  
uuAuGuuuuGuuA\*\*\*\*\*GuuGuuuGUuuuuuGuGuuAUU\*\*\*GAAuCuGu T. vivax  
  
UAuuuGu---GuUUUGuuGuuuGuuuuuG--UG-----A-uA\*GuGuuGuuuuA T. brucei  
uAuuuGU\*\*\*GuuuuGuuAUUUG---uAuuuGU\*\*\*\*\*GuuG-GuAuuGuuuuA T. vivax  
  
uuuuuGuuAU---GGuuuuuuuG-uUUUUGuGGuuuuuGuuuuuuGuuG-uA-uGuA--- T. brucei  
uuuuuGUUAU\*\*\*GGuuuuuuuA\*UUUUGuGGuuuuuGuuuuuuG-uGuuGuA--\*\*U T. vivax  
  
uAG\*\*\*\*GAuuUGuGuGGuAuuuuuGGGA-UCAC\*GuAuAuuuGuGuGGuGUAAuuuAu T. brucei  
UAG---GGuuuGuGUGGuAuuuuuGAGA\*UCA-uGuAUUU----- T. vivax  
  
uuuGuuuAuGA\*\*UGuuuUUUGUUGUAUUACAUUAUUAUAAAUAUAAA T. brucei  
----- T. vivax

#### F. ND7 - Pan-edited dual coding

UGAUACAAAAAACAUACAUAGAUAAAGUAuCAuuuuAuG-uuAuuuuuG--GuAGuu T. brucei  
-----AuGuuuAuuuuuGuuGuAGuu T. vivax  
  
uuuuuACAuuuGuAuCGuuuuACAuuuG\*GUCCACAGCAuCCG\*\*\*CAGCACAuG\*\*G- T. brucei  
uuuuuGCAuuuGuAUCGUuuuACAUuuG-GCCCACAGCAuCCG--CAGCACAuG-\*G\* T. vivax  
  
uGuuuuAuGuuGuuuAuuGuAuuuuuGuGGuGA\*AuuuAuuGuuuA\*\*UA---UUGAuUG T. brucei  
UGuuuuAuGUuGuuuAUUGuAuuuuuGUGGuGA-AuuuAuuGuuuA--uA\*\*\*UUGAuG T. vivax  
  
uAuuAuA\*\*\*G\*GuuA--UUUGCAUCGUGGUACAGAAAAGUUAUGUGAAUAAAAGUGU T. brucei  
uAuuAuA---G-GuuA\*\*UUUGCAUCGAGGUACAGAAAAGUUAUGUGAGUUAAGAGCGU T. vivax  
  
AGAACAAUGUCUUCCGuAUUUCGA---CAGGUUAGAuGuuA---\*GuGuuuGuuGuA T. brucei  
AGAGCAGUGUCUUCCGuAUUUuGAU\*\*\*AGAuAGAuA\*\*\*GuGuuuGuuGuA T. vivax  
  
AUGAGCAuuuGuuGuCuuuA\*\*\*UGuuuuGAGuA--uAuGuuGCGAuGuuGuuGuCGuu T. brucei  
AugAACAuuuAuuGuCuuuA--uGuuuuGAGuA\*\*UAuGuuACGGuGuuGuuAuCAuu T. vivax  
  
ACGuuGuGCAuuuAuGCGuuuAuuAAuuGuA\*\*\*\*GAuuuAC\*\*\*CCGuAGuuuAAuG T. brucei  
GCGuGuuGCAuuuAuGCGuuuAuuGAuuGuA---GAGuuuACU\*\*C\*GUAGuuuAAuG T. vivax  
  
GuuuGuuGuGuAuAuCAuGuAuGGuuuuGG\*AuuuAGGuuGuuuGuCUCCGuuG\*UUuGuG T. brucei  
GuuuAuuGuGuGuGuGuAuGGuAuGAuuuuuAG-AuuuAGGuuGuuuAuCCCCGuuA-UuAuG T. vivax  
  
AuCAuuuGAGGAA-\*\*\*CG\*UGA-CAAuuGAuGACAuuuuuuGAuuuAuG\*\*UUGuGGu T. brucei  
GuCAuuuGAGGAG\*\*\*CG-uGAU\*AAGuuAAuGACGuuuuuuGAuuuGuG--uuGuGGu T. vivax  
  
uGuCGuAuGCAuuuGGCUUUCAuGGuuuuAuuA-\*GGuAUUCUUGAUGAuuuuGuuuuuG T. brucei  
uGuCGuAuGCAuuuGGCUUUCAuGGuuuuAuuG\*\*GGuAuuCUUGAuGAuuuuGuuuuuG T. vivax  
  
GuuuGuuGAuuuuuuGuuGuuGuuGA\*\*\*UAuAuCAuGuuGuuGuuAuGGAuGuu T. brucei

GuuuuGuuGAUUUuuuGuuGuuAuuGA---uAAuAuCGuGuuGuuGuuAuGGAuuGuu T. vivax  
 AuGAuuuGuuAuuuG--uGGGuAA---UCGuuuAuuuUAuuuGCGuuuGC\*\*\*GuGGuuu T. brucei  
 AUGAuuuAuuGuuuG\*\*UGGGuAA\*\*\*UCGuuuGuuuuAuuuGCGuuuGC---GuGGuuu T. vivax  
 GuCAuuuuuuuGAuuuA---uAuGAuuuA\*\*GuuuuA\*\*A\*\*UAGuuuAGuGGuGuuuu T. brucei  
 GuCAuuuuuuuGAuuuG\*\*UAuGAuuuG--GuuuuA---A---uAGuuuAGuGGuGuuuu T. vivax  
 GuCuCGuuCGuuAGGuAuGGuGuGAGAuuGUCGuuGuuGuuGuuA\*\*\*\*UGA\*\*\*\*\* T. brucei  
 GuCACGuuCAuuGGGuAuGGuGuGAGAuuGCCGuuGuuGuuGuuA---UGA----- T. vivax  
 GuUG-uAuuuA---uGuuuuGuuAuGAuuAuuGuuuuGuuuuAuA-GGuGAuGCAuuu T. brucei  
 GuuG\*UAUUUU\* \*\*UGuuuGuuAuGAuuAuuGuuuuGuuuuAuA\*GGuGAuGCAuuu T. vivax  
 GA\*UCGuuuAuuuuuACGuuuGuuGAuGCGuAuGAGuuGuuGAuuuGuAAGCAA-u T. brucei  
 GAC\*CGuuuGuuuuGCGuuGuuGAuAugCGuAuGAGuuGuuGAuuuGuAAGCAA\*U T. vivax  
 GuuuuuuGuuGGuuuuuuGuuuuuG\*\*\*\*GuuuuGuuGuuGuuGuuG\*\*AuuAuuA T. brucei  
 GuuuuuuGuuGGuuuuuuGuuuuuG----GAuuuGuuGuuGuuGuuG--AuuAuuG T. vivax  
 uAuuGuGAuAuuACCAuuG\*\*\*\*AGACCAuuAuuAuGuuAuuuAuAGuuuG--uGGuGu T. brucei  
 uAuuGuGAuGuuACCAuuG---AGACuAuuAuuAuGuuGuuuuAuAGuuuA\*\*UGGuGu T. vivax  
 uGuuGuuGCCGGGuAuAU\*----CAuuuGC\*UUGU-GuuAACACCCCAAAG----G T. brucei  
 uGuuGuuuACCAGGuAuAU\*\*\*\*\*CAUUUGC-UUGU\*GuuGAGCAuCCCAAGG\*\*\*\*\*G T. vivax  
 uGA\*\*\*GuAuuGuuGuuAuuAU\*\*\*\*GuuuuGuGuuGGuuuAuGuuCUCGuuACGuu T. brucei  
 uGA---GuAuuGuuGuuAuuAU\*\*\*\*GuuuuGuGuuGGuuuGuGUUCCCGuuGGuu T. vivax  
 uCGuuGuGC GGGuuuuuuGCA---UA--UUUGuuuAuuGGGuGuuGuuGCGuGGuuu T. brucei  
 uCGuuGuGC GGGuuuuuuACA\*\*\*UA\*\*UUUGuuuGuuGGGuGuuGuuACGuGGuuu T. vivax  
 uuuAuuGCAuGAuuuAGuuGC---C\*GuuuuA--GGuAAuAuuGAuGuuGuuuuGGA T. brucei  
 uuuAuuGCAuGAuuuAGuuGC\*\*\*\*C\*G--uuAuuGGuAAuAuuGAuGuuGuuuuGGA T. vivax  
 uCC--GUAGAUCGuuA\*GuuuuAuAuGuG\*\*A\*\*\*\*\*GGUUAUUGuAGGAUUGUUUAAA T. brucei  
 uCU\*\*GuGGAUUCGuuA\*G----- T. vivax  
 AUUGAAUAAAAA T. brucei  
 ----- T. vivax

#### G. ND8 - Pan-edited non-dual coding

-----CAUUUAAUAAUAAAAGUUUUGGUUGAUUAuuAuuuuuuuAuuuuuuuAuuu T. brucei  
 ----- T. vivax  
 uuGuAuGuuuuuuuuGuuuuuuuGuuuuuuuuUUUUUGuuuGuuuuuuAuAuGuGUuuuGu T. brucei  
 ---AuGuuuuuuuuGuuuuuuuGuuuuuuuuuuGuuuGuuuuuuAuAuGuGuuuuGu T. vivax  
 uuGuuGuGuuA\*\*\*\*CuA\*UUUGuuuA-\* \*\*\*CCCAuuGAGGuAACCA--uuGuuAGuuuA T. brucei  
 uuGuuGuGuuA---CC-A-UUuGuuuA\*\*\*\*CCCAuuGAAu AAC-AuuuuG-uAGuuuG T. vivax  
 uuGGuuC--GuGGUAA---C-C---AuuuuuuGCGUUUUuA\*\*\*UUGGuGuGGuuuAGAG T. brucei  
 uuGA-\*CCCGuGGuAA\*\*\*C\*C\*\*\*AuuuuuuGCGuuuuuA--\*UUGGuGuGGuuuAGAA T. vivax  
 CGuuGuAuuGCuuGuCGuuuAuGuGAuuuAAuuuG-C----CCuA\*\*\*\*GuuAGCAuu T. brucei  
 CGUuGuAuuGCuuGuCGuuuAuGuGAuuuGAuuuGuC\*\*\*\*CC-A---GuuAGCAuu T. vivax  
 GGAuG\*\*\*UUCGuGuuGGGuGG---AGuuuuGGuGGuCA\*\*UC\*GuuuuGCG--GAuuG- T. brucei  
 AGAuG---uuCGuGuuGGGuGG\*\*\*AGuuuuGGuGGuCA--uC-GuuuuGCA\*\*GAuuG\* T. vivax

-AuuuACAAuuGAGuuA-\*UC\*\*GU-\*\*C----GuuGuAuuuAuuGuGGuuuuGuAuGCA T. brucei  
 \*AuuuACAAuuGAGuuA\*\*\*C-CG\*\*\*AC\*\*\*\*GuuGuAuuuAuuGuGGuuuuGuAuGCA T. vivax  
 uGuuuGCCGACAGAU\*\*\*\*G---CC---AuuA----CGCA---UUCAuuGuuuGuuA T. brucei  
 uGuuuGuCCAACAGAU---G\*\*\*CC\*\*\*\*AuuA\*\*\*\*\*CACA\*\*\*UUCAuuGuuuGuuA T. vivax  
 uGuGuuuuuGuuGuuuA-----GCC\*\*A\*\*UGuAuuuAuuG\*GCGC\*\*\*C\*\*\*C---- T. brucei  
 uGuGuuuuuGuuGuuuA\*\*\*\*\*GCC\*\*A--UGuAuuuAuuG-GCGC---C---C\*\*\*\* T. vivax  
 AAGuuuuuAuuGuuuGG---uuGuuGuuuuAuGuuAuuuGAuuuuuAuuuGuGuuuuGuG T. brucei  
 AAGuuuuuGuuAuuuGG\*\*\*UUGuuGuuuuAuGuuGuuuGAuuuuuAuuuGuGuuuuGuG T. vivax  
 uAGuuAuuuAuuuGGGuGAuuuAuuGUGuuuAuGAuuuAA\*\*\*AGAA\*\*AuuCACGGUG T. brucei  
 uAG----- T. vivax  
 AAAUUAAAUUUUGACUAAAU T. brucei  
 ----- T. vivax

#### H. ND9 - Pan-edited dual coding

UUAAUAUCAACUUAAUCCCCCCCCAUAAACAAuuAuAuuAUGuGuA--uAuUUUUuGuuuA T. brucei  
 -----AuG-G-GuuuGuuGuGuuuA T. vivax  
 uuuCGuuuAuGuuuuuGuuuAAuuUUuuuA\*\*UUGuuuGuGuGuA-----GA T. brucei  
 uuuCGUUuGuGUUUUUGuuuGAuuuuGuuuuA--UuGuuuGUGuuGuA\*\*\*\*\*GG T. vivax  
 uGGuGuuUUGuuuGuuuuGuuGA-uuGuAGuuuuuGuuuuuuAuuGuuuuGuuA-Guu T. brucei  
 UGGuGuuuuGuuuGuuuuGuuGA\*UUGuAGuuuuuGuuuuuuAUUGUUUuGuuA\*Guu T. vivax  
 uuuuuuuGuuuuAUUGuAuGuuuuuAuuuuuAAuuuGuG-AuuuuuGuuuuuAuAuuGU T. brucei  
 uuuuuuuGuuuuAuuGuAuGuuuuuAUUuuuuAAuuuAuG\*GuuuuuGuuuuuGuAuuGu T. vivax  
 UGuG-AuUUGuuAuuGAuuGAuuuuuGuGGuuuuuGuuuuuGuCGuuuuAuGuuGuuGUA T. brucei  
 uGuG\*AuuuAuuGuuGAuuGAuuuuuGuGGUuuuuGuuuuuGuCGuuuuAuGuuAuuAuA T. vivax  
 uAuuuuAuuuuGuuuGuuuuGuGuGuuCGuuuGuGuuuuGuuuuuGuGuuGUUGuuuGU T. brucei  
 uAuuuuGuuuuGuuuGuuuuGuGuuuuGuGuuuuCGuuuAuGuuuuGuuuuGuGUUGUUuGuuuuu T. vivax  
 AuuuuuGGGuuG---uGuuuuA-\*GuuuuA\*\*GuuGuuuuuGuUAuGC---GuuuuuG T. brucei  
 GUUUUuuGGGuuG\*\*\*UGuuuuA\*\*GuuuuA--GuuGuuuuGuuAuGC\*\*\*GuuuuuA T. vivax  
 uuGuuGGA---ACGC\*GAuGuuUGAUUUGuuuGGuuuuUAuuuuG--uuGGuAAuGA T. brucei  
 UUGuuAGA\*\*\*ACG-uGAGuGuuuuGAuuuGuuuGGuuuuuAuuuuG\*\*UUGGuAAuGA T. vivax  
 uAuuuuACAUCCGuuuAuuuGuuG--AuuG\*\*\*\*GuuuuuGuuG-GuuuuuuuuGuuGA T. brucei  
 uGuuuuACACCGuuuAuuuGuuG\*\*AuuG---AUUuuuuGuuG\*GuuuuuuuuGuuGA T. vivax  
 -AGuGuuAUCCA--uuAuuuGGuuuGuuuGuAuuGuuAuuuuGuG---uGuuG\*\*GuG-G T. brucei  
 \*AGuGuuAuCCA\*\*UAuuuGGuuuGuuuGuAuuAuuGuuuuGuGuuuuA--GuuA-AuG T. vivax  
 A--GGA-GAUA-GuAuGuACGuuuACAAuGuuA--uuuuuGuuGuuGCAuACC\*\*AAuuU T. brucei  
 A\*\*GAUGauAuGuA---CGuuuACAAuG--GuuuuuuGuuGuuGCAuACC\*\*AAuuu T. vivax  
 UUAuuuG\*CA----uuAuuuuAuuuA\*\*\*AuA\*\*UCACCGuUGUAAUUCUAAUUCUC T. brucei  
 uuAuuuG-CAUUuuuuuG-uuuA--\*G----- T. vivax  
 ACUUCC T. brucei  
 ----- T. vivax

### I. RPS12 - Pan-edited dual coding

CUAAUACACUUUUGAUAAACAAACUAAGUAAAuAuAuuuuGuuuuuuuuGCGuAuGuGA\* T. brucei  
-----AuGuGA\* T. vivax

UUUUUGUAUG\*GuuGuuGuuuAC----\*GuuuuGuuuuAuuuGuuuuAuGuuAuuAuAuG T. brucei  
UUUUUGuAuG-GuuGuuGUUUGC\*\*\*\*GuuuuGuuuGuuuuAuGuuAuuAuAuG T. vivax

AGuCC---G\*\*CGAuuGCCAGuuCCGGuAACCGACGuGuAuuGuAuGC\*\*C\*\*\*\*GuA T. brucei  
AGUCC\*\*\*\*C--CGAuuGCCAGuuCCGGuAAuCGACGuGuGuuGuAuGC--C--\*GuG T. vivax

uuuuAuuUAuAuAAuuuuGuuuG-GA-uGuuGGuuGuuuuuuuGuuGuuuuAuuG--- T. brucei  
uuuuAuuuGuAuAAuuuuG--uGuGGuuGuuGGuuGuuuuuuuGuuG---uG-uGUUU T. vivax

---GuuuA---GuuA--uG\*\*UCAuuAuuuAuuAuAGA--\*\*\*G----GGUGGGuGGuuuu T. brucei  
UuuG---GuuuG-CAUUUG--UCGuuAuuuAuuAuAGA\*\*\*\*\*GGuGGuGGuuuu T. vivax

GuuGAuuuACCC--\*\*\*G\*\*\*\*GuG\*UAAAAGuAuuAuACA\*CG\*\*UAuuG--uA--AGuu T. brucei  
GuuGAuuuACCC\*\*\*\*G--\*\*GuA-UAAAAGuAuuAuACA-CG--uA-uGuuAuuAAuu T. vivax

AGA\*UUUAGAuAUAAGAUUAUGUUUUU T. brucei  
AA----- T. vivax

**APPENDIX H. Alignments of protein sequences of pan-edited dual-coding genes in *L. tarentolae*, *L. amazonensis*, *P. serpens*, and *Perkinsela CCAP1560/4* with *T. brucei* and *T. vivax* sequences.**

A: CR3, B: CR4, C: ND3, D: ND7 5' Editing Domain E: ND9, F: RPS12, G: ND8 (Nondual-coding). Absent sequences were unavailable. All ORF alignments show published protein sequences. All ARF alignments show +1 or +2 (ND7 only) reading frame translations of the full length mRNA sequences. In the ARF alignments of CR3, ND7 and RPS12, translations were made using the alternative *T. brucei* mRNA sequences shown in Figure 8. Two alignments of CR3 ARFs are presented to display the *P. serpens* +2 reading frame which has no stop codons. *L. tarentolae* CR4 published protein sequence shows limited homology in the C terminus to all other CR4 protein sequences. The edited mRNA has two editing sites where 13 U residues are inserted. If the second of these insertion sites is shortened to 12 U residues, the translation of this mRNA has much better homology to other CR4 proteins. Alignments with translations of the two different sequences (13U and 12U) are both shown, with the location of the altered site highlighted in red. While ND8 does not appear to be dual-coding, this alignment was included as well, for comparison of the conservation of a nondual-coding gene with that of the dual-coding genes. It should be noted that ND8 is the only nondual-coding gene that is pan-edited in *L. tarentolae*, *L. amazonensis*, and *P. serpens*, and ND7, A6 and COIII are only partially edited in these species. [Termination codon]

A. CR3

## CR3 ORF Alignment

CR30RF *T. brucei* --MFDCLVLL-FFYCLVFHFFCFLVCDLFLCLLFSFCFLDFCFLFNMGLLLCLFFFFFFI  
CR30RF *T. vivax* --MFDCLVLL-FFLFFFVHFFCFLFICDLFLCLLFFVFCFLFDVFCFLFNMGLLLCLFFFFFFV  
CR30RF *L. amazonensis* --MFDFVIIMFL-FMSFVFHFFCFLFIVDLLFCLMFFFVFLYDFCFVCNLGFCCCLFFFFFFL

CR30RF P. serpens IFLFDFVLFVLFLFLLFFVH

: \*\* : : : : : \* \* \*

CD30PNT-huvec1 LSCDM1-LSCM1X1CSDY1

CR30RF T. bluetel LSFDMLLSFLLYISFRI!  
CR30RE T. vivax I SEDMILISELIYISERYI

CB3 ARE Alignment with *R. serpens* +1 RE

CRS ARE Alignment with *P. serpens* +1 RF  
CB3APeT\_brucei - - - - - BNTNMCMTYKNVIVVVVLFEWELXIEEVEVYI EVICEXVCVILNEVEVWIEVE

-----CL-----V~~EC~~<sup>CC</sup>FFY~~C~~<sup>C</sup>CL~~I~~<sup>I</sup>FF~~V~~<sup>V</sup>CL~~E~~<sup>E</sup>FC~~T~~<sup>T</sup>IC~~F~~<sup>F</sup>Y~~V~~<sup>V</sup>CC~~I~~<sup>I</sup>FF~~V~~<sup>V</sup>LY~~W~~<sup>W</sup>EV~~F~~<sup>F</sup>V~~F~~<sup>F</sup>

CR3ARF P. serpens +1 IIVYNIKHNILYFCCLI---LFCECYFYCFLCIFVFVYLLLICFVVVFYYCFFCLIFVF

CR3ARF L. amazonensis --K|NNMY|V|IYICLI---SLL|CFCLWVLYIFFFVYLLLICYFVWCFLFFFYMIFL

\* : \* . \* \* \* \* \* : \* : \* : \* : \* . \* \* :

CR3ARE T. brucei YIIWVYCCVFFFLFYHLICCYHFYYCIVFVIRLKKYANNFC-  
CR3ARE T. brucei GIIWVYCCVFFFLFYHLICCYHFYYCIVFVIRLKKYANNFC-

CR3ARE T. vivax  
CB3APE P. serpens +1  
CLIVWCCSYVVFLLFYHLYCYYRGCYYI  
VAVVLTWVYEVETFETTITWVYI SVVVI  
VSVVTKST~~E~~VKHTTS-

CRS4F1 P. serpens + I  
CR3ABF I. amazonensis

ORIGIN: 24. *Amazoneurus*

### CR3 ARF Alignment with *P. serpens* +2 RF

CR3ARF T. brucei -----RNINMCMIYKNNVYVVVLFWFWLWYIFFVFYLFVICFYCYLVFVFYWFVFYL

CR3ARF T. vivax -----CL-----VI~~F~~CCFFYCCLYI~~F~~FFVFCFLFVICFYVCCLFFFVYLWIFVFCL

CR3ARF P. serpens +2 -----SKCIYYKNIIFYI-FVWFC----FVFSVIFIVFCAFFLFFIYYWFVLLFF

CR3ARF L. *amazonensis* KUNNMYIVIYICLID--SLLCFCFLWVLYIFFVFYLLLICYFVWCFLFFYMFIVLCV

$\times \quad \bullet \quad \times$        $\bullet \quad \bullet \quad \bullet$        $\times \quad \times \quad \times$        $\bullet \quad \bullet \quad \bullet$        $\times \quad \times \quad \times$        $\bullet \quad \bullet \quad \bullet$

CB3ABE\_T\_brucei TWVYCCVYEEFLFYLHICCYHEYYCIL---VEVIRLKKYA--NNFC---

CR3ARF T. vivax IWWVCCCVYFFFLFYHLICYYRFCYYI VI VIREMIA NTC

CR3ARF P. serpens +2 IIVFFSVWFL--FLLLFWFCLELFIFIFCFCSFWYGFIFHIIICKFPPLLKAFKNISLIV

CR3ARF L. amazonensis IIVFVVVCF~~FFF~~CYPLIWF~~C~~R~~L~~FYYML-----VSDIKIILLL--FL~~I~~K---

B. CR4

CR4 ORF Alignment with *L. tarentolae* 13U translation (Published Sequence)

CR4ORF	T. brucei	LGCDFLLVFWLWLSLFFLWRYNFVYFFFLLFCVFVFFVLLF-LFGLFGFFFYFLLCFVLFFDL
CR4ORF	T. vivax	LGCDFLLVWLYSLFFLWRYNFVYFFFLLFCVFVFFVLLF-FFGLFGFFFYFLLCFVLFFDL
CR4ORF	L. tarentolae 13U	LCCDFVVVVFWLVSVFVYRYNFFFVYFLGVYFFVIIILICIWWFI <del>FF</del> LCLCFDL--F
CR4ORF	L. amazonensis	LCCDFVVVVFWLVSVFVYRYNFFFVFFFLWFVFIFLIIIFI <del>G</del> FGFLFFFVLLVCLVYFEF

CR4ORF	T. brucei	FFMLFFVLLGGFFVFVFFF-----FCLCLFLFVVVVVILLWLL----LIFVYRFIYM
CR4ORF	T. vivax	FFMLFFVLLGGFFVFVFFF-----FCLCLLFVVFIVVLLWLL----LIFVYRFIYM
CR4ORF	L. tarentolae 13U	WIFVVVVFCFLWIFVVCDCVYFIFYIIFCFNCVGVLLVVVYICVSIFLYDVLYFNFNWIIL
CR4ORF	L. amazonensis	LFMLFFVFCFGFLLFVMFILFFVSFF-----VLIVLILFCWCMLF----IFVFRFCIM ..... * . * * . . * . * . * . * .

CR4ORF	<i>T. brucei</i>	RFLF	
CR4ORF	<i>T. vivax</i>	RFVF	
CR4ORF	<i>L. tarentolae</i>	13U	KF---
CR4ORF	<i>L. amazonensis</i>	RFVF	

CR4 ORF Alignment with *L. tarentolae* 12U translation (Hypothetical Sequence)

CR4 ORF	CR4 ORF Alignment with L.	tarentolae 12U translation (hypothetical sequence)
CR4ORF T.	brucei	IILIVVHFFFFYLVCLC----
CR4ORF T.	vivax	IFLFVVHFFFFYLVCLC----
CR4ORF L.	tarentolae 12U	-----KCCCFWFFYVLFCVLYIILFLFFFLFVLCGMFYLCFCYSCLFFFFFFV
CR4ORF L.	amazonensis	---ISNILLFLYIFIYICWLIF-MYSWCYIILFLFFFLFVYYGLFYLYCICLFLICFSL

CR4ORF	T. brucei	LGCDFLLVFWLWLSLFFLWRYNFVYFFFLLFCFVFFFVLLFL-FGLFGFFLYFLLCFVLFFDL
CR4ORF	T. vivax	LGCDFLLVWLYSLFFLWRYNFVYFFFLLFCFVFFFVLLFF-FGLFGFFLYFLLCFVLFFDL
CR4ORF	L. tarentolae 12U	LCCDFVVVFWLWLSVFFFVYRNYFFFFVYFLGVYFVIIIC1IWIFFI <b>EF</b> YVCVLIFYFEF
CR4ORF	L. amazonensis	LCCDFVVVFWLWLSVFFFVYRNYFFFFVYFLWFVIFLIIIF1FGFGLFFFLVLCVLFYFEF

CR4ORF	T. brucei	FFMLFFFVLLGGFFVFVFFFCCLCLFLFVVVVVILLVWLLLLFLFVYRFIYMRFLF <span style="background-color: green;">█</span> -----
CR4ORF	T. vivax	FFMLFFFVLLGGFFVFVFFFCCLCLLFVVIVVLLVWLLLLFLFVYRFIYMRFV <span style="background-color: green;">█</span> -----
CR4ORF	L. tarentolae 12U	LFMLFFFVFCGFLLFVMFILFFISFFVLIVLVCWLFLFIFVFRFFCMTFCILILIGLF <span style="background-color: green;">█</span> NL
CR4ORF	L. amazonensis	LFMLFFFVFCGFLLFVMFILFFVSFFVLIVLVCWLFLFIFVFRFFCMTFCILILIGLF <span style="background-color: green;">█</span> -----

CR4 ARF Alignment with *L. tarentolae* 13U translation (Published Sequence)

<b>CR4 ARF Alignment with L.</b>	<b>tarentolae 13U translation (Published sequence)</b>
CR4ARF T. brucei	-----IYCYLCVFIIILF <span style="background-color: green;">I</span> FWLCIFFFFF <span style="background-color: green;">I</span> FWCVCVCLCTVY <span style="background-color: green;">G</span> IIFYCCFVFCFCCLFWVVCFI
CR4ARF T. vivax	----- <span style="background-color: green;">I</span> FCLLCIFFFFF <span style="background-color: green;">I</span> WCVCVCLCIVCGIFCCCCFFCFCLCVWVCFI
CR4ARF L. tarentolae 13U	QIH <span style="background-color: green;">I</span> NTYMYNCKSVV-VFGFFMY-----YFVC---CIFYCFFFCLFDLVCVMVYFI
CR4ARF L. amazonensis	-----DIKNIK <span style="background-color: green;">I</span> VI-FYYFYIFL---FTFVGWFCLIVVVGIFWFYFYFCYL <span style="background-color: green;">I</span> FTVYFI

CR4ARF	<i>T. brucei</i>	CFVIVVCFFLFWVVIFYWCFCDCIVYFFCDVIILFIFFFFYFLVLCFLYCCFYLVLVFFC-
CR4ARF	<i>T. vivax</i>	CFVIVVCFFLFWVVIFCFI <b>I</b> FIDCIVCFCDVIILFIFFFFCFVLCLFLCCFFLVLVFFC-
CR4ARF	<i>L. tarentolae</i> 13U	YIAFVCLFLVLCYVVILLLCFDCIVFFLFTVIIIFFFLFIWVFIFLFLFWFVGFL <b>FFF</b>
CR4ARF	<i>L. amazonensis</i>	CIALFVCLFYCYVVILLLCFDCIVFFLFI <b>D</b> IILFFFFFFYGLCLFF <b>I</b> LFLYLDLVFYFF

CR4ARF	T. brucei	IFCCVLCYFLIYFLCCFLFWVVFLFLFFFFVYVCFLWLFLFC[FGCCCCYL
CR4ARF	T. vivax	IFCYVLCLYFLICFLCCFLWVVFLFLFFFFVYVCFL[LLLFLY]FGCCCCYL
CR4ARF	L. tarentolae	13U [YVCVLIFYF-----EFLFMFLFFVFCGFLLFVMFILFFFISF
CR4ARF	L. amazonensis	F[FCVWFFFLNFCLCYFLFFADFCCLWLCLFYFLCHFLF[LFYFFFVG

```

CR4ARF T. brucei          --LFICVFYFR|LWYWFYKMFF-----
CR4ARF T. vivax           --LFICVLCF-----
CR4ARF L. tarentolae 13U  FVLIVLVFCWLIFIV-FRFFCMTFCILILIGFL|NLD
CR4ARF L. amazonensis     --LFVCVLCKF--VG-YEFIFI-----

```

CR4 ARF Alignment with *L. tarentolae* 12U translation (Hypothetical Sequence)

CR4ARF T. brucei	-----IYCYLCVFIILF <span style="background-color: green;">I</span> FWLCIFFFIWCVCVLCTVYGI <del>FYCCFVFC</del> <del>CCFLFWVCFI</del>
CR4ARF T. vivax	-----FFCLLCIFFFIWCVCVLICVGIF <del>CCCFFF</del> <del>CCFLCWVCFI</del>
CR4ARF L. tarentolae 12U	QIH <span style="background-color: green;">I</span> NTYMYNCKSVV-VFGFFMY-----YFVC---CIFYFCFFFCLFLDCVMVYFI
CR4ARF L. amazonensis	-----DIKNI <span style="background-color: green;">K</span> VI-FYYFYI <del>FL</del> ---FTFVGWFLCIVVGIFWFY <del>FFFFYFCYL</del> <span style="background-color: green;">I</span> FTVYFI
	: : .: * ** *.* : * . * ***
CR4ARF T. brucei	CFVIVVCF <del>FLL</del> FWVV <del>I</del> YWC <del>DC</del> CIVYFFCDVI <del>I</del> LFI <del>FFF</del> YFVL <del>CF</del> LYCC <del>F</del> LVCLVFFCI
CR4ARF T. vivax	CFVIVVCF <del>FFF</del> FWVV <del>I</del> FC <span style="background-color: green;">I</span> FD <del>I</del> CI <del>V</del> CC <del>F</del> CDVI <del>I</del> LFI <del>FFF</del> CF <del>V</del> LC <del>L</del> FC <del>CC</del> FLVCLVFFCI
CR4ARF L. tarentolae 12U	YIAFVCLFVL <del>CYVV</del> VILL <del>LC</del> DCIV <del>FFL</del> FTVII <del>FFF</del> LF <del>I</del> WVF <del>V</del> F <del>I</del> L <del>LL</del> FWFVG <del>G</del> FL <span style="background-color: red;">FFF</span>
CR4ARF L. amazonensis	CIALFVCLFYCYVV <del>V</del> ILL <del>LC</del> DCIV <del>FFL</del> FI <del>D</del> I <del>I</del> LF <del>FFFFF</del> YGLCL <del>FF</del> <span style="background-color: green;">I</span> LFLYLDLV <del>Y</del> FF
	: . . : :***: :**** *: **: *.*: : : : * : . . : :
CR4ARF T. brucei	FCCVL-CYFLIYFLCCFLFWVVFLFLFFFVYVCY <del>L</del> WLLLFC <span style="background-color: green;">I</span> FGCCCYLCIGL <del>F</del> IC <del>V</del>
CR4ARF T. vivax	FCYVL-CYFLICFLCCFLWVVFLFLFFFVYVCY <del>L</del> FLFY <span style="background-color: green;">I</span> FGCCCYLCIGL <del>F</del> IC <del>V</del>
CR4ARF L. tarentolae 12U	<del>FM</del> FVFWFFILNFCLCCFLFVD <del>FC</del> CLWCLFYFLHFLP <del>L</del> WC <del>C</del> FG <del>CY</del> LYLC <del>F</del> DF <del>V</del> WRF
CR4ARF L. amazonensis	F <span style="background-color: green;">I</span> FCVWFFILNFCLCYFLFFAD <del>FC</del> CLWCLFYFLCHFLP <del>L</del> FYYFVG <del>CY</del> LYLYFVL <del>V</del> FCV <del>C</del> VL
	* . :* ** *.*: * *: * : * : .** * * : *:
CR4ARF T. brucei	YFR <span style="background-color: green;">I</span> LWYWFYKM <del>FF</del>
CR4ARF T. vivax	CF-----
CR4ARF L. tarentolae 12U	VF <span style="background-color: green;">I</span> FLDYFKI---
CR4ARF L. amazonensis	CFKVGYEFIFI---
	*

C. ND3

## ND3 ORF Alignment

## D. ND7 5' Editing Domain

### ND7 5' Editing Domain ORF Alignment

ND7ORF T. brucei	-----MLFLVVFLHLYRFTFGPQHPAAHGVLCCLLYFCGEFIVYIDCIIGYLHR
ND7ORF T. vivax	-----MFIFVVVFLHLYRFTFGPQHPAAHGVLCCLLYFCGEFIVYIDCIIGYLHR
ND7ORF L. tarentolae	ILFSRLHDNYIILYLLIVFLHLYRFTFGPQHPAAHGVLCCLLYISGEFITYIDVIIGYLHR
ND7ORF P. serpens	-----IIFIFFIFVVVFLHLYRFTFGPQHPAAHGVLCCLLYFSGEYITYIDVIIGYLHR : .:*****:*****:.*:*.*** *****

ND7ORF T. brucei	GTEKLCE
ND7ORF T. vivax	GTEKLCE
ND7ORF L. tarentolae	GTEKLCE
ND7ORF P. serpens	GTEKLCE *****

### ND7 5' Editing Domain ARF Alignment

ND7ARF T. brucei	-IQKNMTTWYSIIIVFGSFFTFSFYIWSTASRSTWCFMLFIVFLWWIYCLYWLYRLFA
ND7ARF T. vivax	-----VYFCCSFFAFVSFYIWPTASRSTWCFMLFIVFLWWIYCLYWLYRLFA
ND7ARF L. tarentolae	LNFILPTTR[LYFIFINCFFTLV]IYFRTPASSSPWRIMLFIIISFWRIYNVYRCNYWVFT
ND7ARF P. serpens	-----SNNFYFFYFCCFFAFVSFYVWSTASSRTWCVMLFVIFFRWVYNIYWRNRYRLLT . .*: :* :*. ** * .*: : : * :* :* * :::

ND7ARF T. brucei	SWYRKVMWI
ND7ARF T. vivax	SRYRKVMWV
ND7ARF L. tarentolae	SRYRKVMWI
ND7ARF P. serpens	PWNGKIVWI *: :* :*

E. ND9

## ND9 ORF Alignment

## F. RPS12

### RPS12 ORF Alignment

RPS12ORF T. brucei	-----MWFLYGCCLRFVLFVLCYYMSPRLPSSGNRRVLYAVFYLYNFVWMLRCFF
RPS12ORF T. vivax	-----MWFLYGCCLRFVLFVLCYYMSPRLPSSGNRRVLYAVFYLYNFVWLLRCFF
RPS12ORF L. tarentolae	-----MRVLFLYGLCVRFLFLCCLVLYLSPRLPSSGNRRCLYAICYMFNILWEFCVF-
RPS12ORF L. amazonensis	TFLNLIYFVRVLYLYGLCVRFLFLCCLVLYLSPRLPSSGNRRCLYAISIMFNILWYFLV-
RPS12ORF P. serpens	-----MFFVRSYCLYGFcvrFCvFLCIYVSPRLPSSGNRRVYVVCFNLYSFVIYCFLFG
RPS12ORF Perkinsela	-----MLFGFLVRYGFIEFFFFVSRLPSSGNRFCYELDMRFFFVCYDFVLLG

\*: \* : \* . : :\*\*\*\*\* : . : :

RPS12ORF T. brucei	CC-FIGLVMSLFIIEGGGF---VDLP-GVKYYTRIVS[RE]
RPS12ORF T. vivax	CCVFFGLHLSLFIIEGGGF---VDLP-GIKYYTRMFIN[RE]
RPS12ORF L. tarentolae	CCVCF-LNHLLFIVEGGGF---IDLP-GVKYFSRFFLN[RE]
RPS12ORF L. amazonensis	CCFVF-VIFQLFIVEGGGF---IDLP-GVKYFSRFCNV[RE]
RPS12ORF P. serpens	CCVICYSQSFYFLCEGGGF---VDLP-CIKLVVRVPIA[RE]
RPS12ORF Perkinsela	FSV---LLSSLVVFYEGFGFWLFMDVPFGLYYFSRG[RE]---

. : \*\* \*\* : \* : \* : : \*

### RPS12 ARF Alignment

RPS12ARF T. brucei	NTLLITN[RE]SK-----YILFFLRMWFCMVVYVLFYVIIWVRDCPVPVTDVYCMP
RPS12ARF T. vivax	-----CDFCMVVVVCVLFCFLFYVIIWVPDCPVPVIDVCCMP
RPS12ARF L. tarentolae	---NTYRPI[RE]-----IIFILCVYYFCMVVVFVYIFVWFYI[RE]VHDYLVPVIDVVYMQ
RPS12ARF L. amazonensis	HK-YLFRPF[RE]-----I[RE]FILFVFYICMVYVVFVYFYVWFYI[RE]VHDYQAPVIDVVYMQ
RPS12ARF P. serpens	----LKPIFI[RE]LS[RE]YFYLYLCFLFVVIVYMFVYVFVLYFYVYMLVPVYPVQVIVVFLF
RPS12ARF Perkinsela	-----CCLVFWFVMV[RE]LSFFFLLALVCPVLVIGFVMSW

\*: \* : . : . . \* .

RPS12ARF T. brucei	YF--IYIILFGCCVVFVVL-LV[RE]LCHYLL[RE]RVVVLILIYPV[RE]SIIHVL[RE]VRFYKICF--
RPS12ARF T. vivax	CF--ICIILGCCVVFVVCFLVCICRYLL[RE]RV-VLLIYPV[RE]SIIHVCLLI-----
RPS12ARF L. tarentolae	YV--ICLIFYDFVFFVV-FVFWIIC-CL[RE]LKVVVLLICQE[RE]SIFHVFVFWMRQ[RE]VIKI
RPS12ARF L. amazonensis	LV--LCLIFYDIFWFFAV-LFLWFFS-CL[RE]LKVVVLLICQE[RE]SIFRVFVMCRKFNYLYFY
RPS12ARF P. serpens	VL--ICIVLLFIVFYLVVVLFVILRVFIFYVRRVLLIYHV[RE]SYMSVCQ[RE]PK[RE]IIAS--
RPS12ARF Perkinsela	IWGFFLFVMI[RE]LCCWVFPFCCQVWFFMKVLV-----FGCLWMYRLDCIIFP

: : : . . . : :

RPS12ARF T. brucei	----
RPS12ARF T. vivax	----
RPS12ARF L. tarentolae	ILFR
RPS12ARF L. amazonensis	KN--
RPS12ARF P. serpens	----
RPS12ARF Perkinsela	VV--

## G. ND8 (Nondual-coding)

ND8 ORF Alignment

ND8ORF T. brucei	MFDFDFLFFFVFCYMCFCVCTICLPIELTIVSLLVRGNHFLRFYWCGLERCIACRLCD
ND8ORF T. vivax	MFFFDLFFFVFCYMCFCVCTICLPIELTFCSSLTRGNHFLRFYWCGLERCIACRLCD
ND8ORF L. tarentolae	MFVYDFCFSFVFVFCYMCFCVCTILVLPLEITIVSICVRGNHFLRFYWCGLERCIACRLCD
ND8ORF L. amazonensis	MFCYDFVSFVFVFCYMCFCVCTILPCEITIVSICARGHHFLRFYWCGLERCIACRLCD
ND8ORF P. serpens	MFFFDFFFCFCFVYMCFCCTIVVPCEVSLCSFLVRGTHFLRFYWCGLERCIACRMCD

\*\*\* : \* \* \* . \* \* \* \* \* : \* \* :: \* . \* . \* \* \* \* \* \* \* \* \* \* \* \* \* \* :

LICPSLALDV RVGWSFGGHRFADWFTLSYRRCIYCGFCMHVCPTDAITHSLFVMCFCCLA

LICPSLALDV RVGWSFGGHRFADWFTLSYRRCIYCGFCMHVCPTDAITHSLFVMCFCCLA

FICPSLALDV RCVRSLCGYRFSDVNISYRRCIYCGFCMHVCPTDAITHSCFLFCCIA

FICPSLAI DVRCIRSLCGYRSDFLYRRCIYCGFCMHVCPTDAITHSCFLFCCIA

YICPSVAIDVRCGVSLIGHRAFLHFFISYRRCIYCGFCMHVCPTDAITHSFVVLFSVLLS

\*\*\*\*\*: \* \* \* \* \* : \* \* \* :: \* : \* :

ND8ORF T. brucei

ND8ORF T. vivax

ND8ORF L. tarentolae

ND8ORF L. amazonensis

ND8ORF P. serpens

MYLLAPKFLLFGCCFMFLDFYLCFV

MYLLAPKFLLFGCCFMFLDFYLCFV

MYLCAPKFVLFGCCFMFLDFYLCFV

MYLCAPKFVLFGCCFMFLDFYLCFV

SYLVAPKFILFGCCFMVFDLFCF

\*\*\* :

ND8 ARF Alignment

ND8RF2 T. brucei	NLIILSGWLLFFYFFFIFVCFFLIFCFFFLFLVFLCIVLFVVLLFVYPLS
ND8RF2 T. vivax	-----CFLIFCFFFLFLVFLCIVLFVVLLFVYPLN
ND8RF2 L. tarentolae	-----KHI CIRLKECLFMIFVFLFLVFLCIVFVYVLLLWFYHWSWPLLFLVFLVVTI
ND8RF2 L. amazonensis	-----NIIRSILIKICFVMLFLFLVFLCIVFVYVLLLWFYHVRFLPLLFLVFLVVI
ND8RF2 P. serpens	-----II CVNVVIKCFFLIFFFVFFVLFLCIVFVVVLPLLFLHVKYHCVVFWFVVL

\* :

ND8RF2 T. brucei

ND8RF2 T. vivax

ND8RF2 L. tarentolae

ND8RF2 L. amazonensis

ND8RF2 P. serpens

FCVFIGVVIVNVVLLVVYVI FALV HWMFV LGGV LVV VLRIDLHW VVV FIVV FVCMF

FCVFIGVVIVNVVLLVVYVI WFVPV HMFV LGGV LVV VLRIDLHW VTDVV FIVV FVCMF

FCVFIGVVIVNVVLP AVV YI LYAQV L MFV LEV YV VIG FPMCLILV I VVV FIVV FVCMF

FCVFIGVVIVNVVLP AVV YI LYALVWP FMV LEV YV VIV I PIY FILV I VVV FIVV FVCMF

FCVFIGVVIVNVVLLVVCVII FV LFLVLF MFV VV L V L V L HICFLV IDV VFIVV FVCMF

\* :

ND8RF2 T. brucei

ND8RF2 T. vivax

ND8RF2 L. tarentolae

ND8RF2 L. amazonensis

ND8RF2 P. serpens

ARQMPPLRIHCLLCV FV PCIYWRPSFYCLVVVLCYLIFICVLC SYLFWI ----- YCV

VQQMPLRIHCLLCV FV PCIYWRPSFYCLVVVLCYLIFICVLC -----

VQQTPLRIHVFCYFV VV LP CIYAHNLFLV VV L CYLIFICVLC FLSVYLEK Y - IWLII !!

VPPMLLRIHVFCYFV VV LP CIYAHNLFLV VV L CYLIFICVLC FNLF EYFFY IFYVVC S

VQPMQLPIHLLFYL VCCY PVI WLHPNLF CLVVV L WFLIYFCV FVSCIV FI ----- LFV

. \* \* \* : .. \* \* : : .. \* \* \* \* \* \* \* \* \* \* \* \* \* \* :

ND8RF2 T. brucei

ND8RF2 T. vivax

ND8RF2 L. tarentolae

ND8RF2 L. amazonensis

ND8RF2 P. serpens

YDLKKFTVKLNFD

-----

INF-----

KNLVNV-----

!!!!K RTK! Y --

## APPENDIX I. RPS12 gRNA Alignments for TREU 667 SDM79 (A) and EATRO 164 SDM79 cells (B), and all editing variants (C and D).

Amino acid translations are shown above the mRNA sequences. The cDNA sequence of the most abundant gRNA in its sequence class is shown aligned beneath the fully edited mRNA. Lowercase u's indicate uridines added by editing, asterisks indicate encoded uridines deleted during editing. Nucleotides and deletion sites in the fully edited mRNA were numbered starting from the 5' end (+1=0). gRNAs are colored based on transcript abundance as follows: Blue<100; Green<1,000; Purple<10,000; Orange<100,000; Red>100,000; Black=not quantified. Watson-Crick (|) and G:U base pairs (:) are indicated. Mismatches are indicated by an octothorpe (#). Highlighted sequence represents sequences where multiple CU configurations are possible.

#### A. TREU SDM79 gRNAs - Jv2, I, H, G, F, E, D, C, B1, A

M W F L Y G C C L R F V L F V  
 M V V V Y V L F Y L F  
 CUAAUACACUUUUGAUAAACAAACUAAG\*AAuAAAuuuuGuuuuuuuuuGCGuAuGuGuGAUUUUU\*GuAuG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGuu  
 : ||| ||| : | :: | : ||| : ||| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 11T-TTATTAGAGTGGAGAGACGTATACATTGAAGA-CATGC-CAACAAATA gJv2 gH 13TATAGTGA  
 || : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 gI 14TAA-TATGT-CAGTGATAGATG-CAAAGTAAGATGAACAA  
  
 L C Y Y M S P R L P S S G N R R V L Y A V F Y L Y N F V W M L  
 Y V I I W V R D C P V P V T D V Y C M P Y F I Y I I L F G C C  
 uuAuGuuAuuAuAuGAGuCCG\*\*CGAuuGCCAGuuCCGGuAACCGACGuGuAuGuAuGC\*\*C\*\*\*GuAuuuuAuuUAuAuAAuuuGuuuGGAuGuu  
 ||| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 AATATAATATA gI gF 18TCAATTAAATATATG--G---TATAAGATAAGTGTATTAAGACAAACCTACAAATATA  
 ||| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 AATGTAGTAGATATACTTAGAT--GTTAACGTGTCAGATATA gH gE 12TATATAGAGTGAATATGTTAGAATGAGCCTATAA  
 : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 gG 13TATTCAGT--GTTAGTAGATTGAGGCTATTGGTTGCACATAACATTATA gG  
 ||| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 12TAAATTAGTGACCGAAGGCTAGTGGTT-CATATAACATACG--G---TAATATA gF<sup>P</sup>  
  
 R C F F C C F I G L V M S L F I I E G G G F V D L P G V K  
 V V F F V V L L V STOP  
 CGGuuGuuuuuuuuuGuuGuuuuAuuGGuuuAGuuAuG\*\*UCAuuAuuuAuuAuAGA\*\*\*GGGuGGuGGuuuuGuuGAuuuACCC\*\*\*G\*\*\*\*GuG\*UAAA  
 ||| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 CGCAACAAATAATA gE 14TATTAT--AGTGTAGATGATGTCT--CCTGTAGTCAAAACAACATAAATATA gC  
 : # : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 11TATAATAAGAGAGTAGCAAGATAGTTAAGTCAATAC--AGTAATAAATA gD gB1 11TTAAAGTGAATGGG---T---CAT-ATTT  
  
 Y Y T R I V S STOP  
 GuAuuAuACA\*CG\*UuuGuAAGuuAGA\*UUUAGAuUAAGAUAGUUUUU  
 : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 TATAGTATGT-GC--ATAACATATA gB1  
 ||| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
 16TAAGT-GT--ATAATGTTCAATCT-AAGTCTATATCTATACAATATAAA gA

B. EATRO SDM79 gRNAs - Jv2, I, H, G, F, E, D, C, B1, A

M W F L Y G C C L R F V L F V  
 M V V V Y V L F Y L F

CUAAUACACUUUUGAUAAACAAACUAAG\*AAuAAAuuuuGuuuuuuuuGCGuAuGuGAUUUUU\*GuAuG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGuu  
 : |||:|||:||:||:||:|||:|||:|||:|||:|||:|||:  
<sub>10</sub>TT-TTATTAGTAGTGGAGAGACGTATACTGAAGA-CATGC-CAACAAATA gJv2 gH <sub>13</sub>TATAGTGA  
 : :||: :||: :||: :||: :||: :||: :||: :||:  
 gI <sub>12</sub>TATAG-TATGT-TAATGATAGATG-TGAGACAGAATAAACAA

L C Y Y M S P R L P S S G N R R V L Y A V F Y L Y N F V W M L  
 Y V I I W V R D C P V P V T D V Y C M P Y F I Y I I L F G C C  
 uuAuGuuAuuAuAuGuAGuCCG\*\*CGAuuGCCAGCCGGuACCGACGuGuAuuGuAuGC\*\*\*C\*\*\*\*GuAuuuuAuuUAuAuAuuuuGuuuGGGuu  
 :||:||:  
<sub>12</sub>TAATGTTAGTGAGTTAAATGGGATTAGATTGAT-----ACCTACAAGCATATA  
 gF  
 :||:||:#|||:||#|||:|||:|||:|||:  
<sub>14</sub>TAAATTTAGTGACCGAAGGCTAGTGGCT-CATATAACATACG--G---TAATATA gF<sup>P</sup>  
 :||:||:  
 AATATAATATA gI  
 :||:||:||:||:#:||:|||:#|||:  
 AATGTTAGTGATATACTTAGAT--GTTAACGTGTCAAGATATA gH gE <sub>05</sub>TATTATG--G---TATAAAGTAGATGTATTAGAGTAAGCTTACAA  
 :||:||:#||#|||:  
<sub>13</sub>TATTCAGT--GTTAGTAGATTGAGGCTATTGGTTGCACATAACATTATA gG  
 :||:||:#|||:  
<sub>14</sub>TAATGTTAGGATCATTGGCTGCATATGATATACG--GAAATATA gG<sup>E</sup>

R C F F C C F I G L V M S L F I I E G G G F V D L P G V K  
 V V F F V V L L V STOP  
 GCGuuGuuuuuuuGuuGuuuuAuuGGuuuAGuuAuG\*\*UCauuAuuuAuuAuAGA\*\*\*GGGuGGuuuGuuGAuuuACCC\*\*\*G\*\*\*\*GuG\*UAAA  
 :||:||:  
 CGCAATATATA gE <sub>14</sub>TATAT--AGTAGTGAATGATGTCT---TTTACCGTCAAACAACTAGAATATA gC  
 CGCAACAATAATA gE  
 :#:||:||:||:||:||:||:||:||:||:||:  
<sub>14</sub>TATAATAGAGAGAGTAACGGAGTAATCGAGTCATGC--AGTAATATATATA gD  
 gB1 <sub>11</sub>TTAAAGTAGACTAGATGGA---T---CAT-ATTT

Y Y T R I V S STOP  
 GuAuuAuACA\*CG\*\*UuuGuAGuAGA\*UUUAGuAUAGAUAUGUUUU  
 :||:||:  
 TATAGTATGT-GC--ATAACATATA gB1  
 :||:||:||:||:||:||:  
<sub>16</sub>TAAGT-GT--ATAATGTTCAATCT-AAGTCTATATTCTATACAATATAAA gA

### C. TREU SDM79 Variants

#### J Variants

M W F L Y G C C L R F V L F V  
 M V V V Y V L F Y L F

CUAAUACACUUUUGAUAAACAAACUAAAG\*AAuAAAuuuuGuuuuuuuuGCGuAuGuGAUUUUU\*GuAuG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGuu  
 : |||:|||:|||:|||:|||:|||:|||:|||:|||:|||:  
<sub>11</sub>T-TTATTTAGAGTGGAAAGAGACGTATACTGAAGA-CATGC-CAACAAATA gJv2

M W F L Y G C C L R F V L F V  
 M V V V Y V L F Y L F

CUAAUACACUUUUGAUAAACAAACUAAAG\*AAuAuAuuAGuuuuuuGCGuAuGuGAUUUUU\*GuAuG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGuu  
 : |||:|||:|||:|||:|||:|||:|||:  
<sub>11</sub>T-TTATATAGTTAGAAGAGATGCATGTGCTAGAAG-TATAC-CAACAAACATATA gJv3

M W F L Y G C C L R F V L F V  
 M V V V Y V L F Y L F

CUAAUACACUUUUGAUAAACAAACUAAAG\*AAuAuAuAuuuuGuuuuuuuuGCGuAuGuGAUUUUU\*GuAuG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGuu  
 : |||:|||:|||:|||:|||:  
<sub>12</sub>TATATAGAGTGGAAAGAGACGTATACTGAAGA-CATGC-CAACAAATA gJv1

M W F L Y G C C L R F V L F V  
 M V V V Y V L F Y L F

CUAAUACACUUUUGAUAAACAAACUAAAG\*AAuAuAuuGuuuuuuuuGCGuAuGuGAUUUUU\*GuAuG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGuu  
 : |||:|||:|||:|||:  
<sub>12</sub>TT-TTATGTGATAGTGAAGAGAGTGTATACTAAAGA-CATAC-CAAATATATA gJv4

#### D Variants

##### D

GC\*\*C\*\*\*\*GuAuuuuAuuUAuAuAAuuuuGuuuGGGuuGuuGCGuuGuuuuuuuGuuGuuuuAuuGGuuuAGuuAuG\*\*UCGuuAuuAuuAuu  
 : #:|||:|||:|||:|||:  
<sub>11</sub>TATAATAAACAGAGTAGCAAGATAGTTAAGTCAATAC--AGTAATAAATA gD  
<sub>12</sub>TATATAGAGTGAATATGTTAGAATGAGCCTATAACGCAACAATAATAAATA gE

##### Dx

GCUUCUUUUGAAUAAAuuuGGGuuAuuGGuuuuCGGuuGuuGAuGuAuGuAuG\*\*UCGuuAuuAuuAuu  
 : |||:|||:|||:|||:  
<sub>09</sub>TTTTAAATTAGTGACCGAAGGCTAGTGGTTCATATAACATAC--GGTAATATAAATA gF<sup>p</sup>

## BC Variants

B1, C1

B1\*, C1

B3<sup>t</sup>, C<sup>t</sup>

B4<sup>t</sup>, C<sup>t</sup>

AAGA**UUU**GuGuGuGuuuGuuuuGuGuuuGGuuuuAuACCC\*\*\*G\*\*UUGuuuuG\*UAuuuuuAuAuuuGuAuAuAuuuCA\*CG\*\*UAuuGuAAGuuAGA\*\*\*UAGAu  
 :|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
<sup>18</sup>TATATATGGTAAAGTGTAAATTAGAATATGGG--C--**AACAAAAC-ATATATATA** gC<sup>t</sup>  
 :|:|:|:|:|:|:|:|:|:|:  
<sup>14</sup>T--GACAGAGT-ATAGAGATGTAGATATATAAGAGT-GC--**ATAACATATA** gB4<sup>t</sup>  
 |||:|:|:|:|:|:  
 gA <sup>16</sup>TAAGT-GT--**ATAATGTTCAATCT-AAGTCTA**

## D. EATRO SDM79 Variants

### J Variants

M W F L Y G C C L R F V L F V  
M V V V Y V L F Y L F

CUAAUACACUUUUGAUAAACAAACUAAG\*AAuAAAuuuuGuuuuuuuuGCGuAuGuGAUUUUU\*GuAuG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGuu  
| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>10</sub>TT-TTATTAGTAGTGGAGAGACGTATACATTGAAGA-CATGC-CAACAAATA gJv2

M W F L Y G C C L R F V L F V  
M V V V Y V L F Y L F

CUAAUACACUUUUGAUAAACAAACUAAG\*AAuAuAuuAGuuuuuuuGCGuAuGuGAUUUUU\*GuAuG\*GuuGuuGuuuAC\*GuuuuGuuuuAuuuGuu  
| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>11</sub>TT-TTATATAGTTAGAAGATGCATGTACTAGAAG-TATAC-CAACACATATA gJv3

### G Variants

#### Multiple G sequences from gG

L C Y Y M S P R L P S S G N R R V L Y A	V F Y L Y N F V W M L
Y V I I W V R D C P V P V T D V Y C M P	Y F I Y I I L F G C C
uuAuGuuAuuAuAuGAGuCCG**CGAuuGCCAGuuCCGGuAACCAGuGuAuGC***C****GuAuuuuAuuUAuAuAAuuuuGuuuGGAuGuu	
L C Y Y M S P R L P S F G N R R V L Y A	V F Y L Y N F V W M L
Y V I I W V R D C P A S V T D V Y C M P	Y F I Y I I L F G C C
uuAuGuuAuuAuAuGAGuCCG**CGAuuGCCAGCuuCCGGuAACCAGuGuAuGC***C****GuAuuuuAuuUAuAuAAuuuuGuuuGGAuGuu	
L C Y Y M S P R L P S S G N R R V L Y A	V F Y L Y N F V W M L
Y V I I W V R D C P A L V T D V Y C M P	Y F I Y I I L F G C C
uuAuGuuAuuAuAuGAGuCCG**CGAuuGCCAGCuGGuAACCAGuGuAuGC***C****GuAuuuuAuuUAuAuAAuuuuGuuuGGAuGuu	
:   :   : #   :   :   :   :   :   :   :   :   :   :   :   :   :	
<b>gG1</b> <sub>13</sub> TATTCAGT--GTTAGTAGATTGAGGCTATTGGTTGCACATAACATTATA gG	
: : : : : : # : : : : # : : : : : # :	
AATGTAGTGATATACTTAGAT--GTTAACGTGTCAAGATATA gH	

### G<sup>e</sup>

L C Y Y M S P R L P S P G N R R V L Y A	V F Y L Y N F V W M L
Y V I I W V R D C P V L V T D V Y C M P	Y F I Y I I L F G C C
uuAuGuuAuuAuAuGAGuCCG**CGAuuGCCAGuCCuGGuAACCAGuGuAuGC***C****GuAuuuuAuuUAuAuAAuuuuGuuuGGAuGuu	
:   :   #   :   :   :   :   :   :   :   :   :   :   :   :   :	
<b>gG<sup>e</sup></b> <sub>14</sub> TAATGTGTTAGGATCATTGGCTGCATATGATATACG--GAAATATA gG <sup>e</sup>	
: : : : : : # : : : : # : : : : : : : : : : : : : : :	
AATGTAGTGATATACTTAGAT--GTTAACGTGTCAAGATATA gH	
gE <sub>15</sub> TATTATG--G---TATAAAGTAGATGTATTAGAGTAAGCTTACAA	

## DEF Variants

D, E, F

GACGuGuAuuGuAuGC\*\*\*C\*\*\*\*GuAuuuuAuuUAuAuAAuuuuGuuuGGAuGuuGCGGuuGuuuuuuuGuuGuuuuAuuGGuuuAGuAuG\*\*UCAuuAuuuAuu  
| # | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
C-TCATATAAACATACG--G---TAATA gF<sup>p</sup> gD <sub>14</sub>TATAATAGAGAGACTAACGGAGTAATCGAGTCATGC--AGTAATATATAT  
| : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : | : |  
gE <sub>65</sub>TATTATG--G---TATAAGTAGATGTATTAGAGTAAGCTTACAACGCAATATATA gE  
<sub>12</sub>TATATAGAGTGAATATGTTAGAATGAGCCTATAACGCAACAAATAAATA gE

Dx

**AGCCGGAACCGACGGAGAGCUUCUUUUGAAUAAA**AuuuGGGuuAuuGGuuuuCGGuuGuuGAGuGuAuGuAuG\*\*UCAuuAuuuAuu  
 |||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
<sup>14</sup>**TAAATTTAGTGACCGAAGGCTAGTGGCTATATAACATAC--GGTAATATA** qF<sup>P</sup>

E<sub>e</sub>

**CCGGAAACCACGGAGAuGuC\*\*C\*\*\*\*GuAuA\*AuuuuAAAuuuGGGuuAuuGGGuuCGGuuGuuuuuuuGuuGuuuuAuuGGGuuAGuuAuG\*\*UCAuuAuuuAuu  
#::|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
gD<sub>14</sub>TATAATAGAGAGACTAACGGAGTAATCGAGTCAATGC--AGTAATATATAT  
| :|:|:|:|:|:|:|:|:|:|:|:|:|:  
<sub>12</sub>TATAGTG---TATAT-TGAAGTTAGATCTAATAGACAGAGCAACAATATATA gE<sup>ep</sup>**

### D, E Misanchored, F<sup>e</sup>

## B Variants

B1

B1'

B2FSe

**APPENDIX J. gRNAs identified to edit the RPS12 mRNAs of found in both TREU 667 and EATRO 164 gRNA transcriptomes.**

Editing Region	Population	Sequence	Reads TREU 667	Reads EATRO 164
J	J variant 1	ATA ACAACCGTACAGAAGTTACATATGCAGAGAACGGTGAGATATATT TTTTTTATT	1	
	J variant 2	ATA ACAACCGTACAGAAGTTACATATGCAGAGAACGGTGAGATTTTTTT TTTTTT*	3,190	1,731
	J variant 2	ATA ACAACCGTACAGAAGTTACATATGCAGAGAACGGTGAGATTTAT ATTTTTTTTTT	848	401
	J variant 3	ATAT ACAACAACCATAATGAAGATCATGTACGTAGAACGATTGATATAT TTTTTTTTTTT	751	8,006
	J variant 3	ATAT ACAACAACCATAATGAAGATCATGTACGTAGAACGATTGATATATA TTTTTTTTTTT	628	3,974
	J variant 3	ATAT ACAACAACCATAATGAAGATCATGTACGTAGAACGATTGATATAATT TTTTTTTTT		2,969
	J variant 3	ATAT ACAACAACCATAATGAAGATCATGTACGTAGAACGATTGATAT TTTTTTTTT	72	282
	J variant 3	TAT ACAACAACCATAATGAAGATCATGTACGTAGAACGATTGATATATAAT ATT		184
	J variant 3	AT ACAACAACCATAATGAAGATCATGTACGTAGAACGATTGATATATAATT ATT	28	144
	J variant 4	ATATATA AACCATAACAGAAATTACATATGTGAGAGAACGTTGATAGTGTATT TTTTTAATT	8	
	J variant 4	ATATATA AACCATAACAGAAATTACATATGTGAGAGAACGTTGATAGTGT TTTT	6	
	J variant 4	ATATATA AACCATAACAGAAATTACATATGTGAGAGAACGTTGATAGTGTAT ATT	2	
	I	AT ATAATATAAAAACAATAAGACAGACTGTAGATAGTAATTGTATGA TTTTTTTTTTT	5,762	6,201
	I	ATAT ATAATAAAACAATAAGAACGTTGATAGTGTACTGTATAATT		2,718
I	I	A TATATAATAACATAAAACAATAAGAACGAGATGTAAATGATA TCTATT	503	133
	I	AT ATAATATAAAAACAAGTAGAACGTTGACTGTATAAA TTTTTTTTTCT	6,937	
	I	ATAATATAAAAACAATGAAACGAGACAGTAAATTATATGA TTTTTTTCTTT	687	
	I	AT ATAATATAAAAACAAGTAAACGTTGACTGTATAATT	72	
	H	ATATAGAACTGTCAATTGTA GATTCAATAGTGTGATGAAAGTGA TTTTTTTTTTT*	12,736	4,786
	H	ATATAGAATTAGGCAATCA CGGATTTATATAGTAACTGAAATGA TTTTTTTTTTTG		2,326
H	H	ATATATAACTGGCATTCT CGGATTTGTATAGTGTATGATAAAAGTGAATAA TTTTTTTTT		979
	H	ATATATAACTGGACAATCGTA GGCTTGATGATGAGATGAGTAA TTTTTTTTT		162
	H	ATATATAACTGGACAATCGTA GGCTTGATGATGAGATGAGTAA TTTTTTTTTCTT		90
	H	ATATATAACTGGACAATCGA TGGGCTTGATGATGAGATGAGTAA TTTTTTTTTTG	17,207	
	H	ATATATAACTGGACAATCGA TGGGCTTGATGATGAGATGAGTAA TTTTTTTTT	1,140	
	H	ATATAAACTGTCAATCGA TGGACTTATGTAGTGTATAAGATGA TAATT	828	
	H	TATATAACTGGACAATCGA TGGGCTTGATGATGAGATGAG GAAATT	505	
	H	ATATATAACTGGCAATAT CGGACTCATATAGTGTGAAAGTAAATA TTTTTTTT	169	
G	G	ATACT TACAATACACGTTGGTATCGGAGTT AGATGTTGACTTATT	791	15,732
	G <sup>e</sup>	ATATAAA GGCATATAGTATACTCGCGGTACT AGGATTGTGTAAATT	94,870	90,998
	G <sup>e</sup>	ATATAAAAGC CATATAGTATACTCGCGGTACT AGGATTGTGTAAATT	196	142
F	F <sup>p</sup>	AT ATAATGGCATAACAATATACTCGGTGATCGGAAGCCAGTGATTAAATT TTTTTTTT	2,288	2,552
	F <sup>p</sup>	AT ATAATGGCATAACAATATACTCGGTGATCGGAAGCCAGTGATT TTTTTTTT	599	346
	F <sup>p</sup>	AT ATAATGGCATAACAATATACTCGGTGATCGGAAGCCAGTGATTAA ATT TTTTTTTT	181	256
	F <sup>p</sup>	T ATAATGGCATAACAATATACTCGGTGATCGGAAGCCAGTGATTAAATT AAATT	224	172
	F <sup>p</sup>	AT ATAATGGCATAACAATATACTCGGTGATCGGAAGCCAGTGATTAAATT TTTTTTTGTT	87	120
	F <sup>p</sup>	AT ATAATGGCATAACAATATACTCGGTGATCGGAAGCCAGTGATTAAATT TTTTTTTT	156	
	F <sup>p</sup>	AT ATAATGGCATAACAATATACTCGGTGATCGGAAGCCAGTGATTAAATT T	60	
	F	ATATAAA AACATCCAAACAGAATTATGTGAATAGAATATGGTATATAAT TAACT	63	
	F <sup>e,p</sup>	ATATATAAAATAC AACGGTGATAATTGATGTGATG ATATCTGTAAATT		1,823

E	E	ATAAAT AACACGCAATATCCGAGTAAGATTGTATAAGTGAGATAT ATTTTTTTTTT	10,751	4,860
	E	ATAT ATAACGCAACATTGAATGAGATTATGTAGATGAAATATGGTAT TTTTTT	125	64
	E	ACAAATAACGCAACA TCAGATGAGATTATATAAGTGAGATATG ATATATTTTTTTT		706
	E <sup>ep</sup>	ATATATAACAAACGAGACA GATAATCTAGATTTGAAGTTATATG TGATATTTTTTTT	634	234
	E <sup>ep</sup>	ATACATAACAAACGAGACA GATAATCTAGATTTAGAGTTATATG TGATATTTTTT	113	
D	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAGATAAT ATTTTTTTTTT		344,244
	D	ATATATAATGAC TAACTAAACTGATAAAAGCAGTAGAAGAGATGATGTAAT TTTTTTTTTT	3,627	5,997
	D	ATATATAATGAC TAACTAAACTGATAAAAGCAGTAGAAGAGATGATGTAATATTT TTTTTTTTTT	1,168	3,292
	D	ATATATATGACATAACTT GGCAATGAGATAATGAAAGAGATGGTGTAAAT TTTTTTTTTTCT	10,277	2,172
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAGATA TTTTTTTTTT		1,088
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAGAT TTTTTTTTTT		924
	D	ATATATATAATGACT TAACTGAGCTAATGAGGAATGAGAGAGATAAT ATTTTTTTTTT		809
	D	ATATATATAATGACGTAACTGAT CTAATGAGGAATGAGAGAGATAAT ATTTTTTTTTT		747
	D	ATAAATTAAATGACATAACTT GACTAATAGGACAGTGAAGAGGCCAGTGTAAAT TCTTTTTTTT		558
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGCAATGAGAGAGATAAT ATTTTTTTTTT		534
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGACAATGAGAGAGATAAT ATTTTTTTTTG		521
	D	ATATATAT ATGACGTAACTGAGCTAATGAGGAATGAGAGAGATAAT ATTTTTTTT		512
	D	ATAT ATAATGACATAATTAGACTGATAAGATAACGGAAAAAGTGTATT TTTTTTTT	2,340	382
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAGATAAT ATTTTTTTT		353
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAGATAAT ATTTTTTTT		308
	D	AT ATAAATAATGACATAACTAGTTAGTAAAGTGACGAAGAAGATAAT ATTATTT		261
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAAATAAT ATTTTTTTT		258
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGGGAGATAAT ATTTTTTTT		250
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAGATAAT ATTTTTTTT		193
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAGGTAAAT ATTTTTTTT		178
	D	ATATAT ATAATGACGTAACTGAGCTAATGAGGAATGAGAGAGATAAT TTTTTTTT		160
	D	ATAAAATAATGACATAACTGAATTGATAGAACGATGAGAGAAATAAT ATTTTTTTT	39,218	
	D	ATAAAAT TAATGACATAACTGAATCGATAGAATAATGAGAGAGATAAT ATTTTTTTT	23,327	
	D	ATAAAATA AATGACATAACTGGATTAGTAAAGTGGTGAAGAAAGATAAT ATTTTTTTT	6,953	
	D	A AAAAATGACATAACTGAATTGATAGAACGATGAGAGAAATAAT TTTTTTTTCC	5,004	
	D	ATATATAATA ACATAATTGGATCAGTGAATGAGATAACGA TAGAAATGATATATTTTTTTTTGTT	3,253	
	D	ATAT ATAATGACATAATTAGACTGATAAGATGACGAGAAAAGTGTGTAA TTTTTTTT	2,627	
	D	ATATATA AATGACATGACTAAACTAATAGGGCAGTGAAGAAGACAATG ATTTTTTTT	1,811	
	D	ATATATATC AATGACATAACTGAACACTAGTAAGATAATGAGAGAAAGT TTTTTTTTT	1,594	
	D	ATAT ATAATAATGACGTAATTAAACTGGTAAGATGATAGAAAAAGT TTTTTTTTT	832	
	D	ATATATA AATGACATGACTAAACTGATAGGACAGTAAAGAGGACAATG ATTTTTTTTTGTT	815	
	D	ATATATATTC AATGACATAACTGAACACTAGTAAGATAATGAGAGAA TTTTTTTTCTTT	756	
	D	ATAT ATAATGACATAACTGAACCT GTAAAGATAGCGAGATTTTTTTT	636	
	D	ATATATATC AATGACATAACTGAACACTAGTAAGATAATGAGAGAAAGT TTTTTTTT	220	
	D	ATAAAATAATGACATAACTGAATTGATAGAACGATGAGAGAAAT TTTTTTTT	183	
	D	T TAAATAATGACATAACTGAATTGATAGAACGATGAGAGAAATAAT ATTAATTTTTTTT	159	
	D	ATAAAATAATGACATAACTAAATTGATAGAACGATGAGAGAAATAAT ATTTT	148	
	D	ATAT ATAATAATGACGTAATTAAACTGGTAAGATGATAGAAAAAGTGTAT TTTAAATTTTTT	129	

C	C	ATAC AAATCAACAAAACTACTACTCTCTATAGTGA TGATGATATGTATTTTTGTTTT	80	41
	C	ATATA AGATCAACAAAACTGCCATTCTCTGTAGTAAGTGATGATATAT TTTTTTTTATTTT		28
	C	ATATAAATCAACAAAACTGA TGCCTCTGTAGTAGATAGTGATAT TATTTTTTTCTTT	1,587	
	C	ATATAAATCAACAAAACTAG TGCCTCTATAGTAGATGATGATATATGAT TTTTTTTTTTT	338	333
	C <sup>t</sup>	ATAT ATATACAAAACAACGGGTATAAGATTAAATGTGAAATGGTATATAT TTTTTTTTTTTTTT	8	
	C <sup>t</sup>	ATAT ATATACAAAACAACGGGTATAAGATTAAATGTG TTTTTTTT	1	
B	B1	ATA TACAATACGTGTATGATTTTACTAGGTAGATCAGTGAAATTTTTTTTTT	1,699	17,696
	B1	ATA TACAATACGTGTATGATTTTACTGGTAGATCAGTGAAATT TTTTTT		10
	B1*	ATA TACAATACGTGTATGATTTTACT AGGTAGATCAGTGAAATTTTTTTTTT	69,971	15,488
	B2 <sup>fSe</sup>	ATA TACAATACGTGAATGATTTTACT AGGTAGATCAGTGAAATTTTTTTT		4
	B2 <sup>fSe</sup>	ATA TACAATACGTGGATGATTTTACT AGGTAGATCAGTGAAATTTTTTTTTT	1	2
	B3 <sup>t</sup>	ATATA ACAATACGTGAATGAAATAATTATAGGGATATATGAGATAAT TTTTTTTTTT	1,941	
	B4 <sup>t</sup>	ATA TACAATACGTGAGAATATATAGTAGAGATATGAGACAGT TTTTTTTTTT	3,453	
A	A	AAATAT AACATATCTTATATCTGAATCTAACTGTAAATATGTG AATTTTTTTTTTTT	758	158

## APPENDIX K. ND7 gRNA Alignments for TREU 667 SDM79 (A) and EATRO 164 SDM79 cells (B), and all editing variants (C and D).

Amino acid translations are shown above the mRNA sequences. The cDNA sequence of the most abundant gRNA in its sequence class is shown aligned beneath the fully edited mRNA. Lowercase u's indicate uridines added by editing, asterisks indicate encoded uridines deleted during editing. Nucleotides and deletion sites in the fully edited mRNA were numbered starting from the 5' end (+1=0). gRNAs are colored based on transcript abundance as follows: Blue<100; Green<1,000; Purple<10,000; Orange<100,000; Red>100,000; Black=not quantified. Watson-Crick (|) and G:U base pairs (:) are indicated. Mismatches are indicated by an octothorpe (#). Highlighted sequence represents sequences where multiple CU configurations are possible.

A. TREU 667 ND7

E1v1, D, C, B, A

Y K K T W L H D K Y H F M L F L V V F L H L Y R F T F G P Q H P A  
 I Q K N M T T W S T V S F Y V I F G S F F T F V S F Y I W S T A S R  
 GAUACAAAAAAACAUGACUACAUGUAAGUAuCAuuuuAuGuuAuuuuuGGuAGuuuuuuuACGuuuGuAuCGuuuuACGuuuG\*GUCCACAGCAuCCCG  
 : |#|||||##||| #|:  
 gC<sub>17</sub> TAAGTGTATTAGAGT  
 ||:||||:|:|||:|:|||:||:|:|||:||:|||:||:|||:||:|||:  
 gD<sub>13</sub> TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA  
 ||:||||:|:||:|:||:|:||:||:|||:||:|||:||:|||:  
<sub>09</sub>TATTATAGTAAGATGTTAGTGAGAGCTATCAAAGAACATATAAA gE1v1  
  
 A H G V L C C L L Y F C G E F I V Y I D C I I G Y L H R G  
 S T W C F M L F I V F L W W I Y C L Y W L Y Y R L F A S W  
\*\*\*CAGCACAuG\*\*GuGuuuuAuGuuGuuuAuuGuAuuuuuGuGGuGA\*AuuuAuuGuuuA\*\*UAUUGAuUGuAuuAuA\*\*\*G\*GuuAUUUGCAUCGUGG  
 |||:|||:||:|| :||:||:||:||:|||:|||:|||:  
 gA<sub>13</sub> TAAATAGTAGAT--GTAGTTAGCATAGTAT---T-CAATAAACGTAGTATATA  
 |||:|||:||:||:|||:|||:|||:|||:|||:  
<sub>16</sub>TAATTAAATAGTATGAGAGTGTCACT-TAAGTAGTAAAT--ATAACTAACATA gB  
 |##|||||:||:|||:||:|||:|||:  
---GAAGTGTAC--TATAAAGTGTAAACAAATAACATATATA gC  
  
 T E K L C E Y K  
 Y R K V M W I S T K  
 UACAGAAAAGUUUAUGUGAAUUAAG

B. EATRO 164 ND7

E1v1, D, C, B, A

Y K K T W L H D K Y H F M L F L V V F L H L Y R F T F G P Q H P A  
I Q K N M T T W S T V S F Y V I F G S F F T F V S F Y I W S T A S R  
GAUACAAAAAAACAUGACUACAUAGAUuAuCAuuuuGuuAuuuuGGuAGuuuuuuuACGuuuGuAuCGuuuuACGuuuG\*GUCCACAGCAuCCG

||: |||#|||:#|||#||  
gC 19TAAT-TAGATGTT-TAGAGC

||:|||:||:|||:||:|||:||:||:|||:||:||:|||:||:||:|||:||:||:|||:||:|||:  
gD 13TGTAAGTGTAGATATAGATAATGTAAGC-TGGGTGACGTAGATATATA

:||:|||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
09TATTATAGTAAGATGTAGTGAGAGCTATCAAAGAACATATAAA gE1v1

A H G V L C C L L Y F C G E F I V Y I D C I I G Y L H R G  
S T W C F M L F I V F L W W I Y C L Y W L Y Y R L F A S W

\*\*\*CAGCACAUuG\*\*GuGuuuuAuGuuGuuuAuGuGuuGuGuGA\*AuuuAuGuuGuuA\*\*UAUUGAuUGuAuAuA\*\*\*G\*GuuAUUUGCAUCGUGG

|||||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
gA 12TAAATAGTGAAT--ATAGTTAGTATGATAT---C-TAATAGACGTAGCACATATATA

:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
15TAATAAGTGTATGAAGATGCCATT-TAGATAGCAAAT--ATAACTACATA gB

#||#|||:||:||:||:||:||:||:  
----TCATGTAC--CGTGAAATACAACAAATAATATA gC

T E K L C E Y K  
Y R K V M W I S T K  
UACAGAAAAGUUUAUGUGAAUUAUAAAAG

### C. TREU 667 ND7 Variants

**E Variants**

**E1v2, D, C**

D T K K H D Y M I S T F M L F L V V F L H L Y R F T F G P Q H P A  
 Y K K T W L H D K Y I Y V I F G S F F T F V S F Y I W S T A S R  
 GAUACAAAAAAACAUGACUACAUAGUAACAUuGuuAuGuuGuuGGuAGuuuuuuuACGuuuGuAuCGuuuuACGuuuG\*GUCCACAGCAuCCCG\*\*  
 :|#||||#|||#: :|  
 gC <sub>17</sub>TAAGTGTATTAGAGT--  
 ||:|||:||:|||:||:|||:||:|||:||:|||:||:|||:  
 gD <sub>13</sub>TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA  
 ||:|||:||:|||:||:|||:||:|||:||:|||:  
<sub>12</sub>TAAAATGTAGTGAAGATTGTCGGAGAAATGTAAACATAGCATATACA gE1v2

**E2v1, D, C**

I Q K N M T T W S T V S F M L F L V V F L H L Y R F T F G P Q H P A  
 D T K K H D Y M I S I I Y V I F G S F F T F V S F Y I W S T A S R  
 GAUACAAAAAAACAUGACUACAUAGUAACAUuGuuAuGuuGuuGGuAGuuuuuuuACGuuuGuAuCGuuuuACGuuuG\*GUCCACAGCAuCCCG\*  
 :|#||||#|||#: :|  
 gC <sub>17</sub>TAAGTGTATTAGAGT--  
 ||:|||:||:|||:||:|||:||:|||:||:|||:  
 gD <sub>13</sub>TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA  
 :||:|||:||:|||:||:|||:||:|||:  
<sub>15</sub>TAATATAGTGAATATAATGGAGACTATCGAAGAAATGTAAACATATAAA gE2v1

**E2v2, D, C**

I Q K N M T T W S T V Q L L L V V F L H L Y R F T F G P Q H P A A  
 D T K K H D Y M I S T I V I G S F F T F V S F Y I W S T A S R S  
 GAUACAAAAAAACAUGACUACAUAGUAACAUuGuuAuGuuGGuAGuuuuuuuACGuuuGuAuCGuuuuACGuuuG\*GUCCACAGCAuCCCG\*\*\*CAGC  
 :|#||||#|||#: :##|  
 gC <sub>17</sub>TAAGTGTATTAGAGT---GAAG  
 ||:|||:||:|||:||:|||:||:|||:  
 gD <sub>13</sub>TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA  
 ||:|||:||:|||:||:|||:||:|||:  
<sub>14</sub>TAATAGTAATCATTAGAAGAATGTAGATATGGCAAATGTAAATATA gE2v2

E3<sup>t</sup>, D, C

I Q K N M T T W S T V S F M L F L V V F T F V S F Y I W S T A S R  
D T K K H D Y M I S I I Y V I F G S F Y I C I V L H L V H S I P  
GAUACAAAAAAACAUGACUACAUGAUAAAGUAuCAuuuAuGuuAuuuuuGGuAGuuuuuACAuGuAuCGuuuuuACAuGuG\*GUCCACAGCAuCCCG\*  
||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
<sup>12</sup>TATAGTAAGAGTCATTGAAGATGTGAGTATAGTAAAATGTAAAATATA gE4<sup>t</sup>  
:|:||:|:||:|:||:|:||:|:||:|:  
<sup>15</sup>TAATATAGTGAATATAATGGAGACTATCGAAGAAATGTAAACATATAAA gE2v1  
:|#|||:#|||:#|:  
gC <sup>17</sup>TAAGTGTATTAGAGT-  
||:||:||:|:||:||:||:||:||:||:||:  
gD <sup>13</sup>TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA

E4<sup>t</sup>, D, C

I Q K N M T T W S T  
D T K K H D Y M I S T S Y F W S T  
Y K K T W L H D K Y K L F L V V F T F V S F Y I W S T A S R  
GAUACAAAAAAACAUGACUACAUGAUAAAGUAuCAAGuuAuuuuuGGuAGuuuuuACAuGuAuCGuuuuuACAuGuG\*GUCCACAGCAuCCCG\*\*\*  
||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
<sup>12</sup>TATAGTAAGAGTCATTGAAGATGTGAGTATAGTAAAATGTAAAATATA gE4<sup>t</sup>  
:|#|||:#|||:#|:  
gC <sup>17</sup>TAAGTGTATTAGAGT---  
||:||:||:|:||:||:||:||:||:||:  
gD <sup>13</sup>TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA

C Variant

C<sup>Fst</sup>, B

D D I W S T A S R Y A H G V L C C L L Y F C G E F I V Y I D C  
R H L V H S I P L C T W C F M L F I V F L W W I Y C L Y W L  
T T F G P Q H P A M H M V F Y V V Y C I F V V N L L F I L I  
ACAGACGACAGUGUCCACAGCAuCCCG\*\*\*CuAuGCACAuG\*\*GuGuuuuAuGuuGuuGuAuGuGuGA\*GuuGuuA\*\*UAuuGuuA\*UAuuGAuU  
||:||:||:#|:  
gC<sup>Fst</sup> 57 Reads <sup>14</sup>TAATTGTAGAGT---GATATGTGTAC--TATAAGATAACAGCAAATAACATATATA

||||:||:||:||:||:||:||:||:||:||:||:  
gB 26624 Reads <sup>16</sup>TAATTAATAGTATGAGAGTGTCACT-TAAGTAGTAAAT--ATAACTAACATA

#### D. EATRO 164 ND7 Variants

**E Variants**

**E2v1, D, C**

I Q K N M T T W S T V S F M L F L V V F L H L Y R F T F G P Q H P A  
 D T K K H D Y M I S I I Y V I F G S F F T F V S F Y I W S T A S R  
 GAUACAAAAAAACAUGACUACAUAGUAuCAuuuAuGuuAuuuuGGuAGuuuuuuACAuuuGuAuCGuuuuACAuuuG\*GUCCACAGCAuCCCG\*  
 ||: ||| #||| :#||| #|||  
 gC 19TAAT-TAGATGTT-TAGAGC-  
 ||:||||:||:||||:||:||||||:|| :||||#|||||  
 gD 13TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA  
 ||:||||:||:||:||:||:||:||:||:||:||:||:||:  
 14TATTATAGTGAATACGGTGAGAGTTATCAGAGAAATGTAATAATATA gE2v1

**E4e, D**

M I S T Y C Y W STOP  
 M T T W STOP  
 Y K K T W L H D K Y I L L V V F T F V S F Y I W S T A S R S T  
 GAUACAAAAAAACAUGACUACAUAGUAuCAuGuuAuuGGuAGuuuuuACAuuuGuAuCGuuuuACAuuuG\*GUCCACAGCAuCCCG\*\*\*CAGCA  
 ||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
 13TATAATAGTAATCATCGAACGATGTAGACGTAGCAAATGTAACATATA gE4e  
 ||:||||:||:||:||:||:||:||:||:||:||:||:||:||:||:  
 gD 13TGTAAGTGTAGATATAGTAGAATGTAAGC-TGGGTGACGTAGATATATA

## BC Variants

C1<sup>ex</sup>, B, A

D T K K H D Y M I S T R G D R R Q C P Q H P F I V S F I G I C C L  
Y K K T W L H D K Y K R R Q T T V S T A P V H C F I H W D L L F

GAUACAAAAAAACAUGACUACAUGUAAGUACAAGAGGGAGACAGACGACAGUGUCCACAGCACCCG\*UUCAUuGuuuCAuuCAuuG\*\*GGAUuuGuuGuu

: | : |

**gC1<sup>ex</sup>** 11-T-AATTGATAGAGTAAGTGAC--CCTAAATGACAA

L Y F C G E F I V Y I D C I I G Y L H R G T E K L C E Y K  
I V F L W W I Y C L Y W L Y Y R L F A S W Y R K V M W I ST K

uAuuGuAuuuuuGuGGuGA\*AuuuAuuGuuA\*\*UAUUGAuUGuAuuAuA\*\*\*G\*GuuAUUUGCAUCGUUGGUACAGAAAGUUAUGUGAAUUAAG

**12**TAAATAGTGAAT--ATAGTTAGTATGAT  
:::|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
GTGATATGAAGATGCCATT-TAGATAGCAAAAT--ATAACTACATA gP

ATAAACAAAAATATA qC1<sup>ex</sup>

C2<sup>ex</sup>, B, A

D T K K H D Y M I S T R G D R R Q C P Q H P S L S Y S V F Y C C L  
Y K K T W L H D K Y K R R Q T T V S T A P V I V L Q C V L L L F

**GAUACAAAAAAACAUGACUACAUGUAAGUACAAGAGGAGACAGACGACAGUGUCCACAGCACCGG\***\*UCAuGuCuuACAG\*UGuGuuuuAuuGuuGuu

⋮ ⋮ ⋮ ⋮ ⋮

**gC2<sup>ex</sup>** ~TAATATAGAGTAT--ATTGATAGAATGTC-ACACAAGATAAAC-ACAA : |#|:|:|||||:|||||:#||| gB 15 TAATAAA

L Y F C G E F I V Y I D C I I G Y L H R G T E K L C E Y K  
I V F L W W I Y C L Y W L Y Y R L F A S W Y R K V M W I ST K

uuuuGuAuuuuuuGuGGuGA\*AuuuuAuuGuuuA\*\*\*UAUUGuAUGuAuuAuA\*\*\*G\*GuuAUUUGCAUCCGUGGUACAGAAAAGUUAUGUGAAUUAUAAAAG

12 TAAATAGTGAAT -- ATAGTTAGTATGAT  
:::|:::|:::|:::|:::|:::|:::|  
CTGTATGAGATGGATT-TAGATAGAANT-ATAACTAGATA G

GIGATATGAAGA.

B<sup>ex</sup>, A

I Q K N M T T W S T V Q E E T D D S V H S T R F S T V G Y L L S T I C  
D T K K H D Y M I S T R G D R R Q C P Q H P F Q H S W L F V V D L W  
GAUACAAAAAAACAUGACUACAUGAUAAAGUACAAGAGGAGACAGACGACAGUGGUCCACAGCACAGUUGGuuAuuuGuuGuAGAuuuGu

:||||:||:||||:|||:||:  
gB<sup>ex</sup> 13TAATAGATGACATTAGATA

G E F I V Y I D C I I G Y L H R G T E K L C E Y K  
W I Y C L Y W L Y Y R L F A S W Y R K V M W I ST K

GGuGA\*AuuuAuuGuuA\*\*UAUUGAuUGuAuuAuA\*\*\*G\*GuuAUUUGCAUCGUGGUACAGAAAAGUUAUGUGAAUAAAAG  
12TAAATAGTGAAT--ATAGTTAGTATGATAT---C-TAATAGACGTAGCACATATATA gA

||:|| |:||||| ||| |||||

CTGCT-TGAATAACAAAT--ATAACTATATA gB<sup>ex</sup>

**APPENDIX L. gRNAs identified to edit the ND7 5' mRNAs of found in both TREU 667 and EATRO 164 gRNA transcriptomes.**

Editing Region	Population	Sequences	Reads TREU 667	Reads EATRO 164
E	E1 version 1	AAAT ATACAAATGTAAGAAAACTATCGAGAGTGATGTAGAATGATATT ATTTTTTTTT	4,517	6
	E1 version 1	AAAT ATACAAGTGTAAAGAAAACTATCGAGAGTGATGTAGAATGATATT ATTTTTTTTTTT		1
	E1 version 1	AAAT ATACAATGTAAGAAAACTATCGAGAGTGATGTAGAATGATATT TTTTTTTTTTAT	803	
	E1 version 2	ATATA ATAATGTAAGAGACTATTGAGAGTGGCATAAGTGT TTTTATTTTTTTTTAT		193
	E1 version 2	ACATAT ACGATACAATGTAAGAGGCTTAGAAGTGATGTAAAT TTTTTTTTTTG	7,073	
	E1 version 2	ATACATAT ACGATACAGATGTGAAGAAACTATTAGAGATAATGTAAAT TTTTCTTTTT	213	
	E2 version 1	ATATA ATAATGTAAGAGACTATTGAGAGTGGCATAAGTGTATT ATTTTTTTTTTTTT		148,603
	E2 version 1	ATATA ATAATGTAAGAGACTATTGAGAGTGGCATAAGTGTATT TTTTTTTTTGT		6,174
	E2 version 1	ATATA ATAATGTAAGAGACTATTGAGAGTGGCATAAGTGTATT TTTTTTTTTCT		1,811
	E2 version 1	ATATA ATAATGTAAGAGACTATTGAGAGTGGCATAAGTGTATT ATTTTTTTTTGTTT		659
	E2 version 1	ATATAATAA TAAAGAGACTATTGAGAGTGGCATAAGTGTATT ATAATGATTTTTTTTTT		593
	E2 version 1	AAAT ATACAAATGTAAGAGACTATCAGAGGTAAATAGTGTATT ATTTTTTTGTTT	1,460	504
	E2 version 1	ATATAAT AATGTAAGAGACTATTGAGAGTGGCATAAGTGTATT ATAATGATATT		241
	E2 version 1	AT ATATAATGTAAGAGACTATTGAGAGTGGCATAAGTGTATT ATTTT		203
	E2 version 1	ATATAC ACAATGTAAGAGAGACTATCGAGAGTGACATAAGTGTATT ATTTTTTTATT	192	
	E2 version 2	ATATA ATAATGTAAGAGACTATTGATAGTGG CATAAGTGTATTATAATGATTTTTTTTT		34
	E2 version 2	ATATA TAAATGTAACCGGTATAGATGTAAGAAGATTACTAATGATAATT TTTTTTTTTT	259	
	E4 <sup>t</sup>	ATATA ATAATGTAAGAGACTATTGAGAGTGGCATAAGTGTATT ATAATTTTTTTTT		49
	E4 <sup>t</sup>	ATATA ATAATGTAAGAGACTATTGAGAGTGGCATAAGTGTATT ATAATTTTTTTTT		26
	E4 <sup>t</sup>	ATATA AAATGTAATGATATGACTGTAGAAGTTACTGAGAATGATAT TTTTTTTTT	10,349	
	E4 <sup>t</sup>	ATATA AAATGTAATGATATGACTGTAGAAGTTACTGAGAATGATATA TTTTTTTTT	2,680	
	E4 <sup>t</sup>	ATATA AAATGTAATGATATGACTGTAGAAGTTACTGAGAATGAT TTTTTTTTT	2,360	
	E4 <sup>e</sup>	ATATAC AATGTAACCGATGCAGATGTAGAAGCTACTAATGATAAT TTTTTTTTTTG		218,365
	E4 <sup>e</sup>	ATATAC AATGTAACCGATGCAGATGTAGAAGCTACTAATGATAAT TTTTTTTTTTC		18,562
	E4 <sup>e</sup>	ATATACAA TAAACGATGCAGATGTAGAAGCTACTAATGATAAT TTTTTTATTTTT		793
	E4 <sup>e</sup>	ATATACAATGTAACCGATT CAGATGTAGAAGCTACTAATGATAAT TTTCTTTTTT		777
	E4 <sup>e</sup>	ATATAC AATGTAACCGATGCAGATGTAGAAGCTACTAATAATAAT TTTTTTAATTTTTTTTT		712
	E4 <sup>e</sup>	ATATACAATGT AAACGATGCAGATGTAGAAGCTACTAATGATAAT TTTTTCTTTTT		439
	E4 <sup>e</sup>	ATATAC AATGTAACCGATACAGATGTAGAAGCTACTAATGATAAT TTTCTTTTT		203
	E4 <sup>e</sup>	ATATA AAATGTAATGATATGAGTGTAGAAGTTACTGATAATGATAT ATTTTTTTTT	16	
	E4 <sup>e</sup>	ATATA AAATGTAATGATATGAGTGTAGAAGTTACTGATAATGAT TCTTTTTTTTT	5	
D	D	ATATATAGATGCA GTGGGTCGAATGTAAGATGATATAGATGTGAATTTTTTT TTTTT	36,611	4,096
	D	ATATATAGATGCTGTA GATTAGATGTAGAGTGATATAAG CGTAAATTTTTTTTTTG		2,136
	D	ATATATAGATGCA GTGGGTCGAATGTAAGATGATATAGATGTGAATTTTTTT CTTTTT	3,590	797
	D	ATATATAA ATGCTGTGGATTAGATGTAGAATGATGATGAGTGTGAATTTTTTT TTTTT	15,885	24
	D	ATATA AAATGTAATGATGAGTGTAGAAG TTACTGAGAATGATTTTTTTTC	12,929	
	D	ATATATAATGCA GTGGGATCAGATGTAAAGATGGTATAAGTGTGAATTTTT TTTT	455	
	D	ATATATAGATGCA GTGGGTCGAATGTAAGATGATATAGATGTGAATGTTT GTTTTTTT	332	
	D	ATATATAA ATGCTGTGGATTAGATGTAGAATGATGATGAGTGTGAATTTTTTT TTTTG	332	
	D	GC GGGTCGAATGTAAGATGATAGATGTGAA TGTATTTTTTT	229	

D	D	AGATGCA GTGGGTGCGAATGTAAGATGATAGATGTGAATGTTCTTT TTTTTTT	211	
	D	ATATATAGATGCA GTGGGTGCGAATGTAAGATGATAGATGTGAATGTTTTT GTTTTTTT	203	
	D	TTATATAGATGCA GTGGGTGCGAATGTAAGATGATAGATGTGAATTTTTTGT TTT	195	
	D	ATATATAGATGCA GTGGGTCAAATGTAAGATGATAGATGTGAATTTTTTT	158	
	D	ATATATAGATGCA GTGGGTGCGAATGTAAGATGATAGATGTGAATGTTTTGT TTTTT	154	
	D	ATATATAGATGCAGT GGTCGAATGTAAGATGATAGATGTGAATGTTTTTT	153	
	D	ATATATAGATGCAT TGGGTGCGAATGTAAGATGATAGATGTGAATGTTTTTTTT	128	
	D	GTGGGCCGAATGTAAGATGATAGATGTGAATGTTTTTTTTTT	123	
	D	ATATATAGATGCA GTGGGTGCGAATGTAAGATGATAGATGTGAATGTTTTT GTTTTTT	112	
	D	TATAGATGCA GTGGGTGCGAATGTAAGATGATAGATGTGAATTCTTTTT TTTTTTT	109	
	D	ATATATAAATGC TGGATTAGATGAGATGATAGTGTGAAATTTTTTTTT	106	
	D	ATATATAGATGCA GTGGGTGCGAATGTAAGATGATAGATGTAAATGTTTTTTTT	101	
C	C	ATATAATAAACACATAAACAGTGCATGT ACTCGAGATTTGTAGATTAATTTTTTTTTTTTTTT		6,185
	C	ATAT ATACAATAAACACATGAAATATCATGTC AAGTGAGATATGTGAATTCTTTTTTTTT		114
	C <sup>1ex</sup>	ATATAAA ACAATAAACAGTAAATCCCAGTGAATGAGATAGT TAATTTTTTTTT		183
	C <sup>2ex</sup>	ATATATAAAC ACAATAGAACACTGTAAGATAGT TATATGAGATATAATTTTTT		22
	C <sup>2ex</sup>	ATATATAAAC ACAATAGAGCACACTGTAAGATAGT TATATGAGATATAATTTTTTT		14
	C <sup>2ex</sup>	ATATATATAAAC ACAATAGAACGACTGTAAGATAGT TATATGAGATATTTTTT		13
	C <sup>Fst</sup>	ATATA ACAATAAACACATAGGCACTGTGTATAGTG AGATGTTACTTTCTTATTTTTGTATTGTGTT		23
	C <sup>Fst</sup>	ATAT ATACAATAAACGACATAGAATATCATGTATAGTG AGATGTTAATTTTTTTGTTT		57
B	B	ATAC ATCAATATAAACGATAGATTACCGTAGAAGTATAGTGAATAAT TTTTTTTTTTATT		4,718
	B	ATATA CAATATAAACAGTAGATTCACTGCAGAAGTATGATAGATAAT TTTTTTTTTG		3,557
	B	ATACAT ATATAAACATGAATTCACTGTGAAGATACGATAGTGAAT TTTTTTTTTGTTT		390
	B	ATACA AATCAATATAATGATGAATTCACTGTGAGAGTATGATAATA TTAATTTTTTTTCT		26,624
	B	ATATA CAATATAAACAGTAGATTCACTGCAGAGATATGATAGATAATAA TTTTTTTTTT		3,050
	B	ATATA CAATATAAACAGTAGATTCACTGCAGAGATATGATAGATAATAA TTTTTTTTTCTA		1,675
	B	ATATA CAATATAAACAGTAGATTCACTGCAGAGATATGATAGATAATAAATTTT TTTTTT		798
	B	ATACAT ATATAAACAAATGAATTCACTGTGAAGATACAATAGTATA TTTTTTTTTTTT		579
	B	ATATA CAATATAAACAGTAGATTCACTGCAGAGATATGATAGATAAT TTTTTTTTT		327
	B	ATACA AATCAATATAATGATGAATTCACTGTGAGAGTATGAT TTTTTTTTTTT		110
	B	ATATAT AAAAATAATTCACTAGAGATATAGTAAAGTGTATGAGACATT TTTTTT		30
	B <sup>ex</sup>	ATAT ATCAATATAAACATAAGTCGTACATAGATTACAGTAGATAATT TTTTTTTTTT		20,696
	B <sup>ex</sup>	ATAT ATCAATATAAACATAAGTCGTACATAGATTACAGTAGATAATT TTTTTTTTTT		1,024
	B <sup>ex</sup>	AT ATCAATATAAACATAAGTCGTACATAGATTACAGTAGATAATTAA GTTTTTTTTTT		714
	B <sup>ex</sup>	AT ATCAATATAAACATAAGTCGTACATAGATTACAGTAGATAATT TTTTTTTTT		693
	B <sup>ex</sup>	ATAT ATCAATATAAACGATGAATTGTCATAGATTACAGTAGATAATT TTTTTTTT		57
	B <sup>ex</sup>	ATCAATATAAACGATGAATTGTCATAGATTACAGTAGATAATT TTTT		38
	B <sup>ex</sup>	ATAT ATCAATATAAACATGAGTTCATAGATTACAGTAGATAATT TTTTTTTT		25
A	A	ATATATA CACGATGCAGATAATCTATAGTATGATTGATATAAGTGATAAATT TTTTTTTT		993
	A	ATA TATGATGCAAATAACTTATGATACGATTGATGTAGATGATAAATT TTTTTTTTTG		4,939

#### APPENDIX M. Predicted ND7 protein sequences.

First start codons without premature termination codons were translated. Blue sequences are found in TREU cells only, orange sequences are found in EATRO cells only, and black sequences are found in both cell lines.

RF1

E1v1            MLFLVVFLHLYRFTFGPQHPAAHGVLCCLLYFCGEFIVYIDCIIGYLHRGTEKLCYK

E1v2            MISTFMLFLVVFLHLYRFTFGPQHPAAHGVLCCLLYFCGEFIVYIDCIIGYLHRGTEKLCYK

E2v1FS   MISIIYVIFGSFFTFSFYIWSTASRYAHGVLCCLLYFCGEFIVYIDCIIGYLHRGTEKLCYK

E2v2FS   MISTIVIGSFFTFSFYIWSTASRYAHGVLCCLLYFCGEFIVYIDCIIGYLHRGTEKLCYK

RF3

E2v1   MISIIYVIFGSFFTFSFYIWSTASRSTWCMLFIVFLWWIYCLYWLYYRLFASWYRKVMWI !

E2v2   MISTIVIGSFFTFSFYIWSTASRSTWCMLFIVFLWWIYCLYWLYYRLFASWYRKVMWI !

Minor RF2 variants

E3t            MISIIYVIFGSFYICIVLHLVHSIPQHMVFYVVYCIFVVNLLFILIVL !

E1v1FS        MLFLVVFLHLYRFTFGPQHPAMHMVFYVVYCIFVVNLLFILIVL !

E1v2FS   MISTFMLFLVVFLHLYRFTFGPQHPAMHMVFYVVYCIFVVNLLFILIVL !

No AUG

E4tRF3        LVVFTFSFYIWSTASRSTWCMLFIVFLWWIYCLYWLYYRLFASWYRKVMWI !

E4eRF3        LLLLVFTFSFYIWSTASRSTWCMLFIVFLWWIYCLYWLYYRLFASWYRKVMWI !

## APPENDIX N. CR3 gRNA Alignments for TREU 667 SDM79 (A), and all editing variants (B).

Amino acid translations are shown above the mRNA sequences. The cDNA sequence of the most abundant gRNA in its sequence class is shown aligned beneath the fully edited mRNA. Lowercase u's indicate uridines added by editing, asterisks indicate encoded uridines deleted during editing. Nucleotides and deletion sites in the fully edited mRNA were numbered starting from the 5' end (+1=0). gRNAs are colored based on transcript abundance as follows: Blue<100; Green<1,000; Purple<10,000; Orange<100,000; Red>100,000; Black=not quantified. Watson-Crick (|) and G:U base pairs (:) are indicated. Mismatches are indicated by an octothorpe (#). Highlighted sequence represents sequences where multiple CU configurations are possible.

#### A. TREU 667 CR3 gRNA alignment

$gG^t$ ,  $gF^t$ ,  $gE^{t'}$ ,  $gE$ ,  $gD$ ,  $gC$ ,  $gB1B2$ ,  $gA1$

M F D C L V L L F F Y C L F V H F  
AGAAAUAUAUAUAGUGUAUGAUUAuuAAuuAuuuuCAuuuuAuGuuuGA\*\*\*UUGGuuuGGGuuuuGuuGuuuuuUUUAuuGuuuGuuuGuACAUuu  
|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:  
G<sup>t</sup>, TTAATTAGTGAAGTGAGATGTAACT---AACAAAGTCAAAACAACATA

$\text{gF}^t$   ${}_{12}\text{TACATATTAGAGTGACAGAGAAGTGACGAGCAGACATGTAAA}$   
 $\text{gE}^{t'}$   ${}_{11}\text{TATATACTATGTAGA}$

L L C F I L Q I F S V I I I V Y K F S L L D STOP  
uUGuuGuGuuuuAuAuuACAGAuuuuuAGuGuuAuCAuUAuuAuuGuAuAuAAAGuuuuCGuuAUUAGAUAAAAAAAGUAUGCAAUAUUUUUGU

#### B. TREU 667 Cell line variants

## A Variants

A1

L L C F I L Q I F S V I I I I V Y K F S L L D S T K S M Q I I F  
C C V L Y Y R F L V L S L L L Y I S F R Y S T I K K V C K S T F L  
I V V F Y I T D F S T C Y H Y Y C I S T V F V I R L K K Y A N N F C

AuGUuuGuGuuuuAuAuuACAGAuuuuuAGuGuuAuCAuUAuuAuuGuAuAuAAGuUUUCGUUAAAAGAUAAAAAAAGUAUGCAAAUAUUUUUGU

For more information about the study, please contact Dr. John Smith at (555) 123-4567 or via email at [john.smith@researchinstitute.org](mailto:john.smith@researchinstitute.org).

**TTAGTAGTATAGAGTGTAAATGTCTAACAGAGTCACAATAGTAATATATA gB1B2**

04 TATAGTAGTAATGATAGTATAATGTTTCAGAGGCCATAATCTAATTATATA gA1

A2

L L C L Y Y F R F Y G I I F I I V Y K F S L L D S T K S M Q I I F  
C C V Y I I S D F M V S F L L L Y I S F R Y S T I K K V C K S T F L  
I V V F I L F O I L W Y H F Y Y C I S T V F V I R L K K Y A N N F C

AuUGuuGuGuuuAuAuuAuuuCAGAuuuuAuGGuAuCAuuuuUAuuAuuGuAuAuaAGuUUUCGUUAUUAGAUUUAAAAGUAUGCAAAUAUUUUUGU

Digitized by srujanika@gmail.com

TAATGGTAGAAGTGTATGTTGAAAGCAGTAATCTAAAATATA gA2

Digitized by srujanika@gmail.com

TTAGTAGTATAGATATGATGAGGTTAACATACCATAGTAAAAATAATAA gB3<sup>t</sup>

## B Variants

### B1B2

D F C F L F N M G L L L C F I L Q I F S V I I I I V Y K F S  
I F V F Y L I W V Y C C V L Y Y R F L V L S L L L Y I S F R  
G F L F F I ST Y G F I V V F Y I T D F ST C Y H Y Y C I ST V F V  
GGAUUUUUGuuuuuuuAuuuAAuAuGGGUUUuAuUGuuGuGuuuuAuAuuACAGAUUUUUuAGuGuuAuCAuUAuuAuuGuAuAuAAGuUUUCG  
||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
<sub>12</sub>TAAGAGTGAAAGATAGATTATGTCCGAATAGTAACACAAATATAAA gC  
::||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
<sub>13</sub>TGAATAGTAGTATAGAGTGTAACTAGTCAGAGTCACAATAGTAATATATAA gB1B2  
:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
gA1 <sub>04</sub>TATAGTAGTAATGATAGTATATGTTCAGAGGC

### B4<sup>t</sup>

D F C F L F N M G L L L C L F F F I L S F D M L L L S F L L L Y I S F R  
I F V F Y L I W V Y C C V Y F F F L F Y H L I C C Y H F Y Y C I ST V F  
G F L F F I ST Y G F I V V F I F F F Y F I I W Y V V I I F I I V Y K F S  
GGAUUUUUGuuuuuuuAuuuAAuAuGGGUUUuAuUGuuGuGuuuuAuuuuuAuCAuuuGAuAuGuuAuCAuuuUauAuuGuAuAuAAGuUUUCG  
||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
<sub>12</sub>TAATGTAAGTGAGAGAAAAGAGTGAAATAGTAAACTATACAATATATAA gB4<sup>t</sup>  
||||:||||:||||:||||:||||:||||:||||:||||:||||:  
<sub>12</sub>TAAGAGTGAAAGATAGATTATGTCCGAATAGTAACACAAATATAAA gC gA2 <sub>14</sub>TAATGGTAGAAGTGTGGTATATGTTCGAAAG  
||||:||||:||||:||||:||||:||||:||||:||||:  
gB4 <sub>08</sub>TAATATAGTGAGTTATATAGTGATAGTAGAGATAATGACATATATTATATA

### B3<sup>t</sup>

D F C F L F N M G L L L C L Y Y F R F Y G I I F I I V Y K F S  
I F V F Y L I W V Y C C V Y I I S D F M V S F L L L Y I S F R  
G F L F F I ST Y G F I V V F I L F Q I L W Y H F Y Y C I ST V F V  
GGAUUUUUGuuuuuuuAuuuAAuAuGGGUUUuAuUGuuGuGuuuuAuAuuuCAAGAUUUuAuGGuAuCAuuuUauAuuGuAuAuAAGuUUUCG  
||||:||||:||||:||||:||||:||||:||||:||||:  
<sub>12</sub>TAAGAGTGAAAGATAGATTATGTCCGAATAGTAACACAAATATAAA gC gA2 <sub>14</sub>TAATGGTAGAAGTGTGGTATATGTTCGAAAGC  
||||:||||:||||:||||:||||:||||:  
<sub>12</sub>TATAGTAGTATAGATATGATGAGGTTAACGATACCATAGAAAAATATATAA gB3<sup>t</sup>

### C Variants

C

C L L F S F C F L L D F C F L F N M G L L L C F I L Q I F S V I  
V C Y L V F V F Y W I F V F Y L I W V Y C C V L Y Y R F L V L  
M F V I S T F L F F I G F L F F I S T Y G F I V V F Y I T D F S T C Y  
AuGuuuGuuA\*UUUAGuuuuuuGuuuuuuuAuGGGuuuGuuuuuGuuuuAAuAuGGGuuuAuUGuuGuGuuuAuAuuACAGAuuuuuAGuGuuA  
||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
<sup>12</sup>TAAGAGTGAAGATAGATTATGTCCGAATAGTAACACAAATATAAA gC  
||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:||||:  
TATGAGTAAT-GGATCAAAGACAGAGAATAACCTAAAGATATA gD      gB1B2 <sup>13</sup>TGAATAGTAGTATAGAGTGTAAAGAGTCACAAT

### D Variants

D

V S L Y F F V D F C L C L L F S F C F L L D F C F L F N M G L  
Y H C I F L W I F V Y V C Y L V F V F Y W I F V F Y L I W V Y  
I I V F F C G F L F M F V I S T F L F F I G F L F F I S T Y G F I  
GuAuCAuuGuAuuuuuuuGuGG\*\*\*AUUUUUGuuuAuGuuuGuuA\*UUUAGuuuuuuGuuuuuuuAuGGGuuuuuGuuuuuuAAuAuGGGuuuA  
:||||:||||:||||:||||:||| :||||:||||:||||:||||:||| :||||:||||:||||:||||:||||:||||:||||:  
TATGAGTAATATGAGAGAATATC---TAAAGACAGATACAAACAAT-ATATA gE<sup>t</sup>      gC <sup>12</sup>TAAGAGTGAAGATAGATTATGTCCGAAT  
:||||:||||:||| :||||:||||:||||:||||:||| :||||:||||:||||:||||:  
<sup>11</sup>TAGAATATGAGTAAT-GGATCAAAGACAGAGAATAACCTAAAGATATA gD

**E Variants**

**E**

C L V L L F F Y C L F V H F F C F L F V C D      L F L C L L F S F C  
 I V W F C C F F I V C L Y I F F V F Y L F V      I C F Y V C Y L V F V  
 L F G F V V F L L F V C T F F I C L W      F V F M F V I ST F L  
 A\*\*\*UUGuuuGGuuuuGuuGuuuuUUUAuuGuuuGuuGuACauuuuuuuGuuuuuuAuGuuGuG\*\*\*A\*\*UUUGuuuuuAuGuuGuuA\*UUUAGuuuuuG  
 | : | : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>12</sub>TACATATTAGAGTGCAGAGAACATGTAAATATA gF<sup>t</sup>      gD<sub>11</sub>TAGAATATGAGTAAT-GGATCAAAGAC  
 | : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>11</sub>TATATAGTATGTAGAGAACAGAACAC---T--AAATACATA gE<sup>t'</sup>  
 : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>13</sub>TATAGTTATAGATGGAGCAC---T--GAACGAGAACATAACAT-AGATA gE

**E<sup>t</sup>**

S F G G L L C V S L Y F F V D      F C L C L L F S F C F L L  
 V L V V Y C V Y H C I F L W      I F V Y V C Y L V F V F Y W  
 F W W F I V C I I V F F C G      F L F M F V I ST F L F F I G  
 A\*\*\*\*\*GuuuuGGuGGUUUAuuGuGuGuAuCAuuGuAuGG\*\*\*AUUUUUGuuuAuGuuGuuA\*UUUAGuuuuuGuuuuuuAuGuuG  
 | : | : | : | : | : | : | : | : | : | : | : | : |  
 T-----CGAGATTATTAGATGGCACATATAGTA-CATAAATTATA gFG<sup>t\*xp</sup>  
 : : : : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>13</sub>TGTGTATAGTAATATGAGAGAACATTC---TAAAGACAGAACATAACAT-ATATA gE<sup>t</sup>  
 : | : | : | : | : | : | : | : | : | : | : | : | : |  
 gD<sub>11</sub>TAGAATATGAGTAAT-GGATCAAAGACAGAGAACATAAC

**E<sup>t'</sup>**

C L V L L F F Y C L F V H F F C F L F V C D      L F L C L L F S F C  
 I V W F C C F F I V C L Y I F F V F Y L F V      I C F Y V C Y L V F V  
 L F G F V V F L L F V C T F F I C L W      F V F M F V I ST F L  
 A\*\*\*UUGuuuGGuuuuGuuGuuuuUUUAuuGuuuGuuGuACauuuuuuuGuuuuuuAuGuuGuG\*\*\*A\*\*UUUGuuuuuAuGuuGuuA\*UUUAGuuuuuG  
 | : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>12</sub>TACATATTAGAGTGCAGAGAACATGTAAATATA gF<sup>t</sup>      gD<sub>11</sub>TAGAATATGAGTAAT-GGATCAAAGAC  
 | : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>11</sub>TATATAGTATGTAGAGAACAGAACAC---T--AAATACATA gE<sup>t'</sup>  
 : | : | : | : | : | : | : | : | : | : | : | : |  
<sub>13</sub>TATAGTTATAGATGGAGCAC---T--GAACGAGAACATAACAT-AGATA gE

## FG Variants

FG<sup>tx</sup>

K T L V C S F G G L L C V S L Y F F V D F C L  
 K H ST F V V L V V Y C V Y H C I F L W I F V Y  
 K N I S L ST F W W F I V C I I V F F C G F L F M  
 AAAAAACAGuuuAGuuuGuA\*\*\*\*\*GuuuGGuGGUUUAuuGuGuGuAuCAuuGuAuuuuuuuuuGuGG\*\*\*AUUUUUJGuuuA

||||::|::||      ||:|::|::|:||::|||:|||:|||:#|||:||

**14** TAATTGAGTAT-----CGAGATTATTAGATGGCACATATAGTA-CATAAAATTATA gFG<sup>txp</sup>

**qE<sup>t</sup>** 13 TGTGTATAGTAATATGAGAGAATATC---TAAAGACAGAT

F<sup>t</sup>

I N Y F H F M F D C L V L L F F Y C L F V H F F C F L F V C D L F  
 L I I F I L C L I V W F C C F F I V C L Y I F F V F Y L F V I C  
 N S T L F S F Y V W L F G F V V F L L F V C T F F L F F I C L W F V  
 AA\*\*\*AA\*\*\*AA\*\*\*CA\*\*\*AA\*\*\*Gu\*\*\*GA\*\*\*GG\*\*\*GG\*\*\*Gu\*\*\*Gu\*\*\*Gu\*\*\*Gu\*\*\*Gu\*\*\*Gu\*\*\*Gu\*\*\*Gu\*\*\*Gu\*\*\*G\*\*\*A\*\*\*A\*\*\*A\*\*\*GG\*\*\*Gu

AuuuuAuuuAuuuCAuuuuAUGuuuuGA\*\*\*UUGuuuGGuuuGuuuuGuuuuGuuuuGuACAUuuuuuGuuuuGuuuuGuuuuGuuuuGuG\*\*\*A\*\*UUGu

— 1 —

**gE** <sub>13</sub>TATAGTTATAGATGGAGCAC---T--GAACG

||||||||||||||||||||||||||||||||||||||||

TATTAGAGTGACAGAGAAGTGACGAGCAGACATGTAAAAA

$QE^{t'}$  11 TATATAGTATGTAGAGAAGACAAGGAATAAGCAAACAC---T--AAATA

FG<sup>t</sup>

[View Details](#) [Edit Details](#) [Delete](#)

TATATAT----GTAGTAGTGAATAGAAT---TAAGACCAACATGTAATATATA gFGP

Digitized by srujanika@gmail.com

**11 TATATAGTATGTAGAGAACAGAAATAAGCAAACAC---T--AAATACATA** *get*

**gE** <sub>13</sub>TATAGTTATAGATGGAGCAC---T--GAACGAGAAT

G<sup>t</sup>

E I S T I C V W Y I I N Y F H F M F D C L V L L F F Y C L F V  
K Y K Y V Y D I S T L I I F I L C L I V W F C C F F I V C L  
R N I N M C M I Y N S T L F S F Y V W L F G F V V F L L F V C  
AGAAAUAUAAAUAUGUGUAUGAUAAAuuAuuAuuuuCAuuuuAuGuuuGA\*\*\*\*UUGuuuGGuuuuGuuGuuuuUUUAuuGuuuGuuuG  
|||||:|:|||||:|:||:|:||| |:|:|:|:|:|:|:|:  
<sub>13</sub>**TTAATTAGTGAAAGTGAGATGTAACT**----**AACAAGTCAAACAACAATATA** gG<sup>t</sup>  
| :|:|:|:|:|:|:|:|:|:|:|:|:|:  
gF<sup>t</sup> <sub>12</sub>**TACATATTAGTAGTGACAGAGAAGTGACGAGCAGAC**

## **APPENDIX O. CR3 gRNA Alignments for EATRO 164 (A), and all editing variants (B).**

Amino acid translations are shown above the mRNA sequences. The cDNA sequence of the most abundant gRNA in its sequence class is shown aligned beneath the fully edited mRNA. Lowercase u's indicate uridines added by editing, asterisks indicate encoded uridines deleted during editing. Nucleotides and deletion sites in the fully edited mRNA were numbered starting from the 5' end (+1=0). gRNAs are colored based on transcript abundance as follows: Blue<100; Green<1,000; Purple<10,000; Orange<100,000; Red>100,000; Black=not quantified. Watson-Crick (|) and G:U base pairs (:) are indicated. Mismatches are indicated by an octothorpe (#). Highlighted sequence represents sequences where multiple CU configurations are possible.

#### A. EATRO 164 CR3 SDM79/SDM80

gFG<sup>ep</sup>, gE<sup>ep</sup>, gD<sup>ep</sup>, gC<sup>e</sup>, gB5<sup>e</sup>, gB4, gA2  
 M C M I Y K L T I V L G G I L V I I V Y L V V M S  
 AGAAAUAUAAAUAUGUGUAUGAUAAUAAAuuACAAuuGuGuuA\*\*\*\*\*GGuGGG\*\*\*AuuuuGGuGAuCAuuGuuuAuuuGGuuG\*UUA\*\*\*\*UGAG  
 |||||:::|:||| |||::: |||:|||:|||:|||:  
 13TAATTGTTGGTATAAT----CCATTT---TAGAGTCACTAGTAGCAACATATA gFG<sup>ep</sup>  
 |||||:|||:|||:|||:|||:|||:  
 gE<sup>ep</sup> 12TATAGTAGTAAGTGAAATTGAC-GAT---GTTC  
 :||: |:| :||:  
 gD<sup>ep</sup> 12TAAT-AGT---GTAA  
  
 C I L C F V M V I V F Y L I W V Y C C V Y I T Y V F L L L L S F L  
 uuGuAUUUUAuGuuuuGuuAuGGuuAuuGuuuuuuAuuuAAuAuGGGuuuAuUGuuGuGuuuAuAuuGuGuuuuAuuGuuGuuAuCAuuuuUA  
 |:|||||:|||||:  
 AGTATAAAATACAAAATATATA gE<sup>ep</sup> |||:|||:|||:|||:|||:  
 :|||:|||:|||:|||:|||:  
 GACGTGAAATATAGAGTAGTACTAATAACAAAATATA gD<sup>ep</sup> 12TAATAGTAGTAAATGTAGTAGTAAATATAGAGATGACAACAATAGTATATA gB5<sup>e</sup>  
 ::|||:|||:|||:|||:|||:  
 :|||:|||:|||:|||:  
 13TGATAGTAGTGAAGATAGATTGTATCCAAATAGTAACACAAATATAAA gC<sup>e</sup> gB4 11TAATATAGTGAGTTATATAGTGATAGTAGAGAT  
 :|||:|||:|||:  
 ga2 13TAATAGTGGAAAGT

L L Y I S F R Y STOP  
uuAuuGuAuAuAAAGuUUUCGUUAUUAGAUUAAAAAGUAUGCAAAUAAUUUUUGU  
||||:|||||  
**AATGACATATATTATATA gB4**  
| : | : : | : | : | : | : | : |||||  
**AGTAGTATATGTTCAAGAGCAGTAATCTAAAATATA gA2**

## B. EATRO 164 CR3 Variants

### A Variants

#### A1

L L C F I L Q I F S V I I I I V Y K F S L L D STOP  
C C V L Y Y R F L V L S L L L Y I S F R Y STOP  
I V V F Y I T D F ST C Y H Y Y C I STOP  
AuUGuuGuGuuuuAuAuuACAGAuuuuuAGuGuuAuCAuUAuuAuuGuAuAuAAGuUUUCGUUAUUAGAUUAAAAAGUAUGCAAUAUUUUUGU  
:|:|:|:|||||:||:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
<sub>11</sub>TATAGTAGTAATGATACTATGTTCAAGAGCAATAATTATA gA1  
||:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
TAGTAGTATAGAATATGATGTCTAACAGACTACAATAGTAAATATA gB1B2

#### A2

L L C L Y Y L C I F I V V I I F I I V Y K F S L L D STOP  
C C V Y I T Y V F L L L L S F L L L Y I S F R Y STOP  
I V V F I L L M Y F Y C C Y H F Y Y C I STOP  
AuUGuuGuGuuuAuAuuACuuAuGuAuuuuuAuuGuGuuAuCAuuuuUAuuAuuGuAuAuAAGuUUUCGUUAUUAGAUUAAAAAGUAUGCAAUAUUUUUGU  
:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
<sub>13</sub>TAATAGTGGAAAGTAGTAGTATATGTTCAAGAGCAGTAATCTAAATATA gA2  
||:|:|:|:|:|:|:|:|:|:|:|:  
<sub>11</sub>TAATATAGTGAGTTATATAGTGATAGTAGAGATAATGACATATATTATA gB4  
||:|:|:|:|:|:|:|:|:|:|:  
TAATAGTGTAAATGTAGTGAATATATAGAGATGACAACAATAGTATATA gB5<sup>e</sup>

## B Variants

### B1B2

G Y C F L F N M G L L L C F I L Q I F S V I I I I I V Y K F S  
V I V F Y L I W V Y C C V L Y Y R F L V L S L L L Y I S F R  
L L F F I ST Y G F I V V F Y I T D F ST C Y H Y Y C I STOP  
GGuuAuuGuuuuuuAuuuAAuAuGGGuuAuUGuuGuGuuuuAuAuuACAGAuuuuuAGuGuuAuCAuUAuuAuuGuAuAuAAGuUUUCG

||||::|::|:|:|:|:|:|:|:|:|:|:|:|:|:  
**gB1B2** <sub>14</sub>TAATAGTAGTATAGAATATGATGTCTAAGAGTCACAATAGTAAATATA

:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:

**gA1** <sub>11</sub>TATAGTAGTAATGATAGTATATGTTCAGAGGC

:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:

<sub>13</sub>TGATAGTGAAAGATAGATTGTATCCAAATAGTAACACAAATATAAA gC<sup>e</sup>

### B5<sup>e</sup>

G Y C F L F N M G L L L C L Y Y L C I F I V V I I F I I V Y K F S  
V I V F Y L I W V Y C C V Y I T Y V F L L L S F L L L Y I S F R  
L L F F I ST Y G F I V V F I L L M Y F Y C C Y H F Y Y C I STOP  
GGuuAuuGuuuuuuAuuuAAuAuGGGuuAuUGuuGuGuuuuAuAuuGuGuuAuCAuuuuUAuuAuuGuAuAuAAGuUUUCG

||||:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
**gB5<sup>e</sup>** <sub>12</sub>TAATAGTGAAATGTAGTGAATATATAGAGATGACAACAATAGTATA

:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:

**gB4** <sub>11</sub>TAATATAGTGAGTTATATAGTGTAGTAGAGATAATGACATATATTATATA

:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:

<sub>13</sub>TGATAGTGAAAGATAGATTGTATCCAAATAGTAACACAAATATAAA gC<sup>e</sup>

**gA2** <sub>13</sub>TAATAGTGGAAGTAGTAGTATATGTTCAAGAGC

### B6<sup>e</sup>

G Y C F L F N M G L L L C L Y Y L M Y F Y C C Y H F Y Y C I STOP  
V I V F Y L I W V Y C C V Y I I L C I F I V V I I F I I V Y K F S  
L L F F I ST Y G F I V V F I L S Y V F L L L S F L L L Y I S F R  
GGuuAuuGuuuuuuAuuuAAuAuGGGuuAuUGuuGuGuuuuAuAuuGuGuuAuCAuuuuUAuuAuuGuAuAuAAGuUUUCG

||||:|:|:|:|:|:|:|:|:|:|:|:|:|:|:  
**gB6<sup>e</sup>** <sub>12</sub>TAGTAGTATAGATATGATAGAGTGCATAAGAATAACAACAATATATA

:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:

**gB4** <sub>11</sub>TAATATAGTGAGTTATATAGTGTAGTAGAGATAATGACATATATTATATA

:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:|:

<sub>13</sub>TGATAGTGAAAGATAGATTGTATCCAAATAGTAACACAAATATAAA gC<sup>e</sup>

**gA2** <sub>13</sub>TAATAGTGGAAGTAGTAGTATATGTTCAAGAGC

B7e

S N L ST C C Y L F G Y I I W Y V V I I F I I V Y K F S  
G V I Y S V V I C L D I S F D M L L S F L L L Y I S F R  
E ST F I V L L F V W I Y H L I C C Y H F Y Y C I STOP

GGAGuAAuuuAuAGuGuuGuuAuuUGuuuGGAuAuAuCAuuuGAuAuGuuAuCAuuuuUAuuAuuGuAuAuAAAGuUUUCG

:|||||:||:|||:||:||:||:||:|||||:|||||

<sup>09</sup>TATTAAGTGTACAGTAGTGAATGAGTCTATATAG gB7e

||||:||:|||:||:||:||:|||||:|||||

<sup>14</sup>TAAGTGTACAGTAGTGAATGAGTCTATATAGTAAACGAATATAAA gB7e

|||||||:||:|||:||:||:|||||:||:|||||:|||||

gB4 <sup>11</sup>TAATATAGTGAGTTATATAGTGATAGTAGAGATAATGACATATATTATATA

:|||||:||:||:||:||:|||||:||:|||||:|||

gA2 <sup>15</sup>TAATAGTGGAAGTAGTAGTATATGTTCAAGAGC

C Variants

C

C L L F S F C F L L D F C F L F N M G L L L C L Y Y L C I F I V V I  
V C Y L V F V F Y W I F V F Y L I W V Y C C V Y I T Y V F L L L L  
M F V I ST F L F F I G F L F F I ST Y G F I V V F I L L L M Y F Y C C Y  
AuGuuuGuuA\*UUUAGuuuuuuGuuuuuuuuAuuGGAuuuuuGuuuuuuAuuuAAuAuGGGuuuAuUGuuGuGuuAuAuuACuuAuGuAuuuuuAuuGuuGuuA  
||||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:  
gC<sub>14</sub> TAATTTAGAAGTAGAGAGTGAATTATGCTCAAATAACAACATATATA  
||||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||:  
TATGAGTAAT-AGATCGAACAGACAGAATAACCTAAAGATA gD  
gB5<sup>e</sup> 12 TAATAGTGTAAATGTAGTGAATATATAGAGATGACAACAAT

C<sup>e</sup>

L Y F M F C Y G Y C F L F N M G L L L C L Y Y L C I F I V V I  
S C I L C F V M V I V F Y L I W V Y C C V Y I T Y V F L L L L  
V V F Y V L L W L F F I ST Y G F I V V F I L L L M Y F Y C C Y  
AGuuGuAUUUUAuGuuuuGuuAuGGuuAuuGuuuuuuAuuuAAuAuGGGuuuAuUGuuGuGuuAuAuuACuuAuGuAuuuuuAuuGuuGuuA  
::|:|||:|||:|||:|||:|||:|||:|||:|||:  
13 TGATAGTGAAAGATAGATTGTATCCAAATAGTAACACAAATATAAA gC<sup>e</sup>  
||||:|||:|||:|||:|||:|||:|||:|||:  
gB5<sup>e</sup> 12 TAATAGTGTAAATGTAGTGAATATATAGAGATGACAACAAT  
|:|||:|||:|||:|||:  
TTGACGTGAAATATAGAGTAGTACTAATAACAAAATATA gD<sup>eP</sup>

C<sup>e80</sup>

S C I L C F V M F C M I I L ST C D L L C L Y Y L C I F I V V I  
V V F Y V L L C F V W L F Y S V I C C V Y I T Y V F L L L L  
L Y F M F C Y V L Y D Y F I V W F V V F I L L L M Y F Y C C Y  
AGuuGuAUUUUAuGuuuuGuuAuGuAuGAuuAuuuuAuAGuGuGuGAuuUGuuGuGuuAuAuuACuuAuGuAuuuuuAuuGuuGuuA  
:|||:|||:|||:|||:|||:|||:|||:  
14 TAATATAGAGTATATTGATAGAATGTCACACTAGATAACACAAATATATAA gC<sup>e80</sup>  
||||:|||:|||:|||:|||:|||:  
gB5<sup>e</sup> 12 TAATAGTGTAAATGTAGTGAATATATAGAGATGACAACAAT  
|:|||:|||:|||:  
TTGACGTGAAATATAGAGTAGTACTAATAACAAAATATA gD<sup>eP</sup>

D Variants

D

Y Q Y L F C D      L F L C L L F S F C F L L D F C F L F N M G L  
 I S I C F V      I C F Y V C Y L V F V F Y W I F V F Y L I W V Y  
 V S V F V L W      F V F M F V I ST F L F F I G F L F F I ST Y G F I  
 GuAuCAGuAuuuGuuuuGuG\*\*\*A\*\*UUUGuuuuuAuGuuuGuuA\*UUUAGuuuuuGuuuuuuAuGGGuuuuGuuuuuuAuuuAAuAuGGGuuuA

:||||::|:||:| :||:| :||:|:||:|:||:|:||:|:|

**gD<sub>13</sub>TATAGAATATGAGTAAT-AGATCGAACAGAGAATAACCTAAAGATA**

:||||:||:|:||:| :||:| :||:|:||:|:||:|

**TATAGTTATAGATGGAGCAC---T--GAACGAGAACAAACAAT-AGATA gE**

:||:|:||:||:|:||:|:||:|:||:|:||:|:||:|:||:|

**gC<sub>14</sub>TAATTTAGAAGTAGAGAGTGATTATGCTCAAAT**

D<sup>e</sup>

D H C L F **G** C Y E      L Y F M F C Y **G** Y C F L F N M G L  
 I I V Y L V V M S      C I L C F V M V I V F Y L I W V Y  
 S L F I W L L W      V V F Y V L L W L L F F I ST Y G F I  
 D C R **L** F S C Y E      L Y F M F C Y D Y C F C F I G D A ND7 protein seq

GAuuGUCGuuuAuuuAGuuG-uuA\*\*\*\*UGA\*\*\*\*\*GuUGuAuuuuuAuGuuuuGuuAuGAuuAuuGuuuuuGuuuuAuAGGuGAuGCAuuu ND7 777-866  
 GAuCAuuGuuuAuuuGGuuG\*UUA\*\*\*\*UGA---GuuGuAUUUUAuGuuuuGuuAuGGuuAuuGuuuuuuAuuuAAuAuGGGuuuA

:||:|:||:|:||:|:||:|:||:|:||:|:||:|:||:|:||:|

**gD<sup>ep</sup><sub>12</sub>TAAT-AGT---GTT---TGACGTGAAATATAGAGTAGTACTAATAACAAAATATA**

:||||:||:||:||:|:||:|:||:|:||:|:||:|:||:|

**ATAGTAGTAAGTGAATTGAC-GAT---GTT---CAGTATAAAATACAAAATATA gE<sup>ep</sup>**

:||:|:||:||:||:||:||:||:||:||:||:||:||:||

**gC<sup>e</sup><sub>13</sub>TGATAGTGAAGATAGATTGTATCCAAAT**

## E Variants

E<sup>e</sup>

R	W D	F G D H C L F G C Y	E L Y F M F C Y G
	G G	I L V I I V Y L V V M	S C I L C F V M
	V G	F W W S L F I W L L	W V V F Y V L L W

A\*\*\*\*\*GGuGGG\*\*\*AuuuuGGuGAuCAuuGuuAuuuGGuuG\*UUA\*\*\*\*UGAGuuGuAUUUUUauGuuuuGuuAuG

| :||::: ||:|::|||||:|

T-----CCATTT---TAGAGTCACTAGTAGCAACATATA gFG<sup>ep</sup>

|||||:|||:|||::||| :||| ::|||:|||:|||:|||:|||

**gE<sup>ep</sup>** 12 TATAGTAGTAAGTGAATTGAC-GAT-----GTTCAAGTATAAAAATACAAAAATATATA

gDep 12 TAAT-AGT----GTTTGACGTGAAATATAGAGTAGTAC

E

AUUUUUUUGGGGGGUUUAGGGuAuCAGGuAuuuGuuuuGuG\*\*\*A\*\*UUUGuuuuuuuAuGuuuGuuA\*UUUAGGuuuuuGuuuuuuAuuGu

1 : 1 | 1 : 1 | 1 : 1 | 1 : 1 | 1 : 1 | 1 : 1 | 1 : 1 | 1 : 1 |

gD<sub>12</sub>-TATAGAATATGAGTAAT-AGAT

9-13-2011-11:00 AM - 11:30 AM

### FG Variants

F<sup>e</sup>

K Y Y H I C V R	W D F G D H C L F G C Y E
N I I I F V L	G G I L V I I V Y L V V M S
I L S Y L C ST	V G F W W S L F I W L L W
AAuAuuAuCAuAuuuGuGuuA*****GGuGGG***AuuuuGGuGAuCAuuGuuuAuuuGGuuG*UUA****UGA	
:  :  :  :  :	:  :  :  :  :  :  :
<sup>11</sup> TATAATAGTGTAAAGTATAGT-----CTATCT---TAAAACCATCAGTAGCATATATA gFG <sup>eP</sup>	
:  :  :  :  :  :  :  :  :	
gE <sup>eP</sup> <sup>12</sup> TATAGTAGTAAGTGAATTGAC-GAT---GTT	

F<sup>e\*</sup>

K H I C V R	W D F G D H C L F G C Y E
K N I F V L	G G I L V I I V Y L V V M S
K T Y L C ST	V G F W W S L F I W L L W
AAAAACAuAuuuGuGuuA*****GGuGGG***AuuuuGGuGAuCAuuGuuuAuuuGGuuG*UUA****UGA	
:  :  :  :	:  :  :  :  :  :  :  :
<sup>13</sup> TAATTGTTAGGTATAAT-----CCATTT---TAGAGTCACTAGTAGAACATATA gFG <sup>eP</sup>	
:  :  :  :  :  :  :  :  :	
gE <sup>eP</sup> <sup>12</sup> TATAGTAGTAAGTGAATTGAC-GAT---GTT	

F<sup>ex</sup>

I I I K F V I	W C F V F D L F C V F H C L F G C Y E
I L S T S S L L	F G V L F L I C F V Y F I V Y L V V M S
Y Y N Q V C Y	L V F C F W F V L C I S L F I W L L W
AuAuuAuAAuCAAuGuuGuuA***UUUGGuGuuuuGuuuuG***AuuuGuuuuGuGuAuuuCAuuGuuuAuuuGGuuG*UUA****UGA	
:  :  :  :  :  :  :#	:  :  :  :  :  :
TATGATGTTAGTCAGTAGC---AAACCAAAAA gG <sup>eP</sup>	
:  :  :	:  :  :  :  :  :  :
TATAATAGTGACTTAGACAGT---AAACCATAAAACAAAAACATATA gF <sup>exp</sup>	
:  :  :  :  :  :	
gF <sup>exp</sup> <sup>12</sup> TATAAAGTGAAGC---TAAGTGAATATAGAGTAACAAATAACATA	
:  :  :  :  :  :  :  :	
gE <sup>eP</sup> <sup>12</sup> TATAGTAGTAAGTGAATTGAC-GAT---GTT	

Ge

K Y K Y V Y D I Y I I I K F V I W C F V F D L F C V  
R N I N M C M I Y I L S T S S L L F G V L F L I C F V  
E I S T I C V W Y I Y Y N Q V C Y L V F C F W F V L C  
AGAAAUAUAAAUAUGUGUAUGAUAAuAuuAuAAuCAAGuuuGuuA\*\*\*UUUGGuGuuuuGuuuuuG\*\*\*AuuuGuuuuGuG  
||||:||:|||||:||:# |||||  
<sub>12</sub>TAATAGAATATGATGTTAGTTCAAGTAGC---AAACCAAAAA gGe<sup>p</sup>  
|:|:|:||:| ||||||:|||||:|||||  
<sub>04</sub>TAATAGAGTATAATAGTGACTIONAGT---AAACCATAAAACAAAAACATATA gF<sup>exp</sup>  
:|:|||:|||:| |||:|||:|||:  
gF<sup>exp</sup> <sub>12</sub>TATAAAGTGAAAGC---TAAGTGGAAATAT

SDM80 only

R N I N M C M I Y K N N G S C G F V G W F R L G Y C Y C E  
AGAAAUAUAAAUAUGUGUAUGAUAAUAAAACAAuGGuA\*\*\*\*\*GuuGuGGuuuGu\*\*UAGGuuGAuuCAGAuuGGG\*UUA\*\*\*UUGGuuAuuGuGA\*\*  
||:||| :|:|||:||:|:| :|||:|||:|||:#|||:| ||| |||||  
gFGe<sup>80</sup> <sub>14</sub>TATCAT-----TAGTATCAGAGT--GTCTAATTAAAGTGTAACTC-AAT---AACATATA  
|||:||:|:| |::|:|||:|||:  
<sub>11</sub>TAATTT-AGT---AGTAGTGACATT--  
  
C C S F C M I I L S T C D L L C L Y Y L C I F I V V I I F I I V Y K  
\*\*AuGuuGuAGuuuuuGuAuGAuuAuuuuAuAGuGuGAuuUGuuGuGuuAuAuuACuuAuGuAuuuuuAuuGuuGuAuCAuuuuUAuuAuuGuAuAuA  
||:|||:|||:|||:|||:#|||:  
--TATAACGTCAAAAACATACGAATATATA gDE<sup>e80</sup>

F S L L D ST

AGuUUUCGUUAUUAGAUAAAAGUAUGCAAAUAAUUUUUGU

**APPENDIX P. gRNAs identified to edit the CR3 mRNAs of found in both TREU 667 and EATRO 164 gRNA transcriptomes.**

Editing Region	Population	Sequence	Reads TREU 667	Reads EATRO 164
G	gG <sup>t</sup>	ATATGT ACAACAAAACCGAGCAATCAGATATAGAGTGAAGTGATTAAATT TTTTTTTTTTC		21,297
	gG <sup>t</sup>	ATAT AACACAAAACGTGAAACAATCAAATGTAGAGTGAAGTGATTAAATT TTTTTTTTTTT	493	
	gG <sup>ep</sup>	AAAAACCAAAC GATGAACCTGATTGTAGTATA AGATAATTTTTTTTTTT	8,171	2,736
	gG <sup>ep</sup>	AATACCGAGC GACAGATTGATTGTAGTATA AGATAATATTTTTTTTTTT		94
FG	gFG <sup>tp</sup>	ATATAT AAAATGTACAACCAGAATTAAGATAAAGTGTGATGTATATAATT TT	21	
	gFG <sup>bp</sup>	ATATTAATAC ATGATATACCGCGATAGATTATTAGAGTTATGAGTTAAT TTTTTTTTTGGTT		2,772
	gFG <sup>bp</sup>	ATAC ATGATATATAACAGTGAACATTAGAAATTATAGGT AATGAGATTATTTTTTTTTTT		354
	gFG <sup>bp</sup>	ATATTAATAC ATGATATGCGCAGTAGACTATTAAAGTTATGAGTTAAT TTTTTTTTTTT		159
	gFG <sup>bp</sup>	ATATTAATAC ATGATATACACGGTAGATTATTAGAGCTATGAGTTAAT TTTTTTTTTTT	44,119	
	gFG <sup>bp</sup>	ATATTAATAC ATGATATACACGGTAGATTATTAGAGCTATGAGTTA TTTTTTTTTTTCT		2,106
	gFG <sup>bp</sup>	ATATTAATAC ATGATATACACGGTAGATTATTAGAGCTATGAGTT TTTTTTTTTTT		936
	gFG <sup>bp</sup>	TATTAATAC ATGATATACACGGTAGATTATTAGAGCTATGAGTTAA CTTTTTTTTTTTT		404
	gFG <sup>bp</sup>	ATATTAATAC ATATACACGGTAGATTATTAGAGCTATGAGTTAAT TTTTT		133
	gFG <sup>bp</sup>	ATATTAATAC ATGATATACACGGTAGATTATTAGAGCTATGA TTAAAATTTTTTTTTTT		115
	gFG <sup>e+p</sup>	ATATAC AACGATGATCACTGAGATTTACCTAATATGG ATGTTAATTTTTTTGGGAACTGAA	53,389	53,459
	gFG <sup>e+p</sup>	ATACAA AACGATGATCACCGAGATTCA GTTAATATGATTGCTAATTTTTTTTT		717
	gFG <sup>e+p</sup>	ATATACAAC ATGATCACTGAGATTTACCTAATATGG ATGTTAATTTTTTTTT		133
	gFG <sup>e+p</sup>	ATATATACAA AAACGATGATACCGAGATTTCATTAAATATGATTGT CTAATTT		143
	gFG <sup>e+p</sup>	ATATAC AACGATGATCACTGAGATTTACCTAATATGG TATGTTAATTTTTCTTTT		97
	gFG <sup>e+p</sup>	ATATAC AACGATGATCACTGAGATTTACCTAATATGGTTAATTT TAGATTTTT		42
	gFG <sup>e+p</sup>	ATATAC AACAAATGATCACTGAGATTTACCTAATATGGTTAATTT TTTTTTTTT		1
	gFG <sup>e+p</sup>	ATATAGAACGATGAC TGCTAGAATTCTGCTTGATATGGATG TAATTTTTTTTTTT		12,250
	gFG <sup>e+p</sup>	ATATAC AACGATGATCACTGAGATTTACCTAATACGA TGTTAATTTTTTTTTTAT		1,418
	gFG <sup>e+p</sup>	ATATAC AACGATGATCACTGAGATTTACCTAATATGGAT TTTTTTTTTAAAGTGCGGCCATAGGGTG		534
	gFG <sup>ep</sup>	ATATATACGATGAC TACAAAATCTATCTGATATGAAATGTGATAATTTT TTTTTTT		1,854
	gFG <sup>ep</sup>	ATATATACGATGAC TGCCAAAATCTATCTGATATGAAATGTGATAATTTT TTTTTTT		173
	gFG <sup>ep</sup>	ATATAC AACGATGATCACTGAGATTTACCTAAT TTTTTTTTTTT		140
	gFG <sup>e80</sup>	ATATACAATAACTCAATG TGAATTAATCTGTGAGACTATGATTACT TTTTTTTGTTTT		146
	gFG <sup>e80</sup>	ATATACAATAACTCAATG TGAATTAATCTGTGAGACTATGATTACTATT TTTTTTTATTT		112
F	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAAATTAGATGAT ATTTTTTTTTTTTC		138,710
	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAAATTAGATGATT TTTTTTTTTTT		33,909
	gF <sup>t</sup>	ATATAATT AAATGTACAGACAAATGATAGAGAGAGACGTGAGATTAAAGT TATATTTTTTTTT		16,892
	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAAATTAGAT TTTTTTTTTTT		941
	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAAATTAGATG TTTTTTTTTTT		932
	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTAAATAGATGAT ATTTTTTTTTTT		679
	gF <sup>t</sup>	ATATATAAAATA TACAAACGGACAATGAGAGAACAGTGAAATTAGATGAT AATATATTTT		603
	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAAATTAGATGA AATATTTTTTTTT		459
	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAAATTAGATA TAATATTTTTTTTTA		263
	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGTGAAATTAGATGAT ATTTTATATTTTTTT		234

	gF <sup>t</sup>	ATATAT AAAATGTACAAACGGACAATGAGAGAACAGCGAAATTAGATGAT ATCTTTTTT		208
	gF <sup>t</sup>	ATAT AAAATGTACAGACGAGCAGTGAGAGACAGTGAGATTA TACATTTTTTTTT	19,078	
	gF <sup>t</sup>	ATATAA GTACAAACAGACAATGAAAAGATGGTAGACTGAGTA TACATTTTTTTT	592	
	gF <sup>t</sup>	AAATT AATGTACAAATAAACGATAGAGAGACAGTGAGATTA TGATTGTAATTCTTTTTT	243	
	gF <sup>t</sup>	ATAT AAAATGTACAGACAAGCAATGAAGAGACAGTGAGATAGTT TTTTTTTTTC	210	
	gF <sup>exp</sup>	ATACAT AATAAACAAATGAGATATATAAGGTGAATCGAAAGTGAATATT TTTTTTTTT		360
	gF <sup>exp</sup>	ATATA CAAAAACAAAATACCAATGACAGATT CAGTGATAATATGAGATAATT		235
	gF <sup>exp</sup>	ATATAC AAAACAAAATACCAATGACAGATT CAGTGATAATATGAGATAATT		197
	gF <sup>exp</sup>	ATACAT AATAAACAAATGAGATATATAAGGTGAATCGAAAGTGAATATT ATT		179
E	gE	A TAGATAACAAACATAAGAGCAAGTCACCGAGGTAGATATTGATATT TTTTTTTTC	10,598	15,900
	gE	A TAGATAACAAACATAAGAGCAAGTCACCGAGGTAGATATTGATATTAA TTTTTTTATT	867	1,342
	gE	A TAGATAACAAACATAAGAGCAAGTCACCGAGGTAGATATTGATATT TTTTTTTTTTT	100	199
	gE	TTAGATAACAAACATAAGAGCAAGTCACCGAGGTAGATATTGATATTAA TTTTTT		102
	gE <sup>t</sup>	ATATAT AATCACAAACAAATAGAAAATGAGAGAGGTGTATGA TACATTTTTTTCTTTT		584
	gE <sup>t</sup>	ATAC ATAAATCACAAACGAATAAGGAACAGAACAGATGTATGA TATATT		388
	gE <sup>t</sup>	AT AATAAATCACAAACAGATAAGAGCAAGAAAGGTGTATGA TTATATT		86
	gE <sup>t</sup>	ATAAAC AACACAAACAGATAGAGACAGAACAGAGATGTATAGATA TTAAAATT		80
	gE <sup>t</sup>	ATAT ATAACAAACATAGACAGAACTATAAGAGATATAATGATATGTGT TTTTTTTTT		29
	gE <sup>ep</sup>	ATAT ATAAAACATAAAATGACTGTAGCACTTAAGTGAATGTATGA TATTTTTTTT		2,493
	gE <sup>ep</sup>	AAATA ATAACAAACATGAGATATAACTGTAGTGTAGATGAATGTATGA TTTTTTTTT		12,105
	gE <sup>ep</sup>	ATATAT TAAAATACAACCTATGATGACTAAGTGAATGTGATTGTCAA TTTTTTTCTTTT		5,559
	gE <sup>ep</sup>	ATAT ATAAAACATAAAATGACTGTAGCACTTAAGTAAATGTATGA TTTTTTTT		131
	gE <sup>ep</sup>	ATATAT TAAAATACAACCTATGATGGCTAAGTGAATGTGATTGTCAA TTTTTTTTTT		2,491
	gE <sup>ep</sup>	ATAT ATAAAACATAAAATGACTGTAGCACTTAAGTGAATGTGATT TTTTTTTATT		428
	gE <sup>ep</sup>	ATATT ATAAAACATAAGATATAACTCATAGTGTATGAATAAGTGTATTTTTT		157
DE	gDE <sup>e80</sup>	ATATATAAG CATAACAAAACGTCAATATTACAGTGATGATTTAATT TTTTTT	525	149
	gDE <sup>e80</sup>	ATATATAAG CATAACAAAACGTCACTATTACAGTGATGATTTAATT T		10
D	gD	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAAG ATATTTTTTTTTT		27,320
	gD	ATATAT AAATCCAATAAGAAATGAAAGCTAGATAGTGAGTATAAGT TTTTTTTTTGTT	535	331
	gD	GTAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAAG ATATTTTTTTGTTTTTTT		315
	gD	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAAG TTTTTTTTTT		206
	gD	ATAT ACAAAAATCCAATGAAAATAAGACTGAGTGATGGATG CAATTTCTTTTTT		142
	gD	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATA TTTTTCTTTTTT		137
	gD	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAAG ATATATT		90
	gD	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAAG ATATATT		74
	gD	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAAG ATATATT		23
	gD	ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATTAGTATAAG ATATATATT		19
	gD	ATATAT AAACAAAATTCAATAAGAAAAGAGACTGAGTAATTGATAAA TTTT		17
	gD	AT ATAGAAATCCAATAAGAGACAGAACAGCTAGATAATGAGTATAAG ATTTTTTTTTT		718
	gD	ATAT AACAAAATCCAATGAGAAATAGAGACTGAGTAATTGATATA TATTTTTATT		217
	gD <sup>ep</sup>	ATAT AAAACAATAATCATGATGAGATATAAGTGCAGTTGTGATAATT TTTTGT		12
	gD <sup>ep</sup>	ATAT ATAATCATAACAAAGATGTAGAGTACGATTAGTAGTTAA TTTTTTTTTT		2
	gD <sup>ep</sup>	ATAT AAAACAATAATCATGATGAGATATAAGTGCAGTTGTGATAATT TTTTTGTTT		31
	gD <sup>ep</sup>	ATAT AAAACAATAATCATATAAGATGTAGAGTACGATTAGTAGATAATT TTTTTGTTT		15
	gD <sup>ep</sup>	ATAT ATAATCATAACAGAGCATAGAATAACAGTTAGTAGTGTAA TTTTTT		8

	gD <sup>ep</sup>	ATAT AAAACAATAATCATAATAAGATGTAAGGTACGATTATGAT TTTTTTTTATGCAACGGTACTGGAG	1	
	gD <sup>ep</sup>	ATAT ATAATCATAACAAGATATAGAATACGATTAGTAA TTTTTTTTTTTT	1	
C	gC	ATAT ATACAACAATAACTCGTATTAAGTGAGAGATGAAGATTAAT TTTTTTTTATTT		1,509
	gC	ATAT TAAACACAACGATAGATCTATTAAGTAGAAGATAGAAATTAA TTTTTTTTTTTTT		88
	gC	A AATATAAACACAATGATAAGCCTGTATTAGATAGAAGTGAGAATT TTTTTTTTG		10,402
	gC	TA AAACACAACAATAGATTCTGTATTAAGTAGAGATAGAGATTAA TTTTTTTTTT		209
	gC	ATAT TAAACACAACAGTAGATCTATTAAGTAGAAGATAGAAATT TTTTTTTT		170
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATATA TTTTTTTTTTTTT		214,918
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATAT TTTTTTTTTTN		180,589
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATATTTT TTTTTTTTT		67,000
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGAT TTTTTTTTTGTT		2,828
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATATTTGT TTTTTT		553
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATATTTT GTTTTTTTT		621
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATATTT TTTTTTTTT		1,014
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATA TTTTTTTTTTATTTTTT		330
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATAT TTTTTTTTT		358
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCGACTGTAAGATAGTTATGAGATATA TTTTTTTTTTG		33,434
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCGACTGTAAGATAGTTATGAGATATA TTTTTTTTTTC		14,634
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCGACTGTAAGATAGTTATGAGATATTT TTTTTTTTT		3,887
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCGACTGTAAGATAGTTATGAGAT TTTTTTTTTT		828
	gC <sup>e80</sup>	ATAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGGGATATA TTTTTTTCTTTT		387
	gC <sup>e80</sup>	NTAT ATATAAACACAATAGATCACACTGTAAGATAGTTATGAGATATACTTT TTTTT		374
	gC <sup>e</sup>	A AATATAAACACAATGATAAACCTATGTAGATAGAAGTGATAGTT TTTTTTTTT		32
	gC <sup>e</sup>	AAAA AATATAAACACAATGATAAGCCTGTATTAGATAGAAGTGATAATT TTTT		3
B	gB1B2	ATATATA AATGATAACACTGAGAATCTGTAGTATAAGATATGATGATAA TTTTTTTTTTTT		37,769
	gB1B2	ATATATA AATGATAACACTGAGAATCTGTAGTATAAGATATGATGATA TTTTTTTATTT		330
	gB1B2	ATATATAATA ATAACACTGAGAATCTGTAGTATAAGATATGATGATAA TTTTTTTTTTTTTTT		180
	gB1B2	ATATATA AATGATAACACTGAGAATCTGTAGTATAAGATATGATAAT TTTTTTTTTT		98
	gB1B2	ATAT ATAATGATAACACTGAGAATCTGTAATGTGAGATATGATGATAAGTT TTTTTTTTT		10,346
	gB1B2	ATAT ATAATGATAACACTGAGAATCTGTAATGTGAGATATGATGATAA TTTTTTTTTTT		2,337
	gB1B2	ATATATA AATGATAACACTGAGAATCTGTAGTATAAGATATGATGATAA TTTTTTTCTTTTT		139
	gB1B2	ATAT ATAATGATAACACTGAGAATCTGTAATGTGAGATATGATGATA TTTTTTTTATTT		124
	gB3 <sup>t</sup>	ATAT ATAAAAATGATACCATAGAATTGGAGTAGTATAGATATGATGATA TTTTTTTTTT		1,183
	gB3 <sup>t</sup>	ATAT ATAAAAATGATACCATAGAGTTGGAGGTGATGTAGATATGATGGTA TTTT		474
	gB3 <sup>t</sup>	ATAT ATAAAAATGATACCATAGAATTGGAGTAGTATAGATATGATGAT TTTTTTTTTT		383
	gB4	ATATA TTATATACAGTAATAGAGATGATAGTGATATTGAGTGATATA ATTNTTTTTTT		260
	gB4	ATATACAGTAATAGAGATGATAGTGATATTGAGTGATATAATTAATATGT GATTTTAATTTTTTTCTTT		8
	gB4	ATACAGTAATAGAGATGATAGTGATATTGAGTGATATAATTAATATGTTT T		6
	gB4	ATA TATATACAATAAAAGAGTGATAGCAGTG TATTAGTGATGATTAATTTTTTTTTT		5
	gB4	ATA TATATACAATAATGAGAATGATGACAGTG TATTAGTGATGTAATATTTTTTTT		367
	gB4	ATATA TTATATACAGTAATAGAGATGATAGTGATATTGAGTGATATAATTAATATGT GATTTTAATATAATA		18
	gB4 <sup>t</sup>	ATAT ATAACATATCAAATGATAAAAGTGAGAAAAGAGAGTGATGTAAT TTTTTCTTTCTATCTATTACAT		13
	gB5 <sup>e</sup>	ATAT ATGATAACACAGTAGAGATATAAAGTGATGTAATGTGATAAT TTTTTTTTT		3,708
	gB6 <sup>e</sup>	ATAT ATAACACAATAAGAATACGTGAGATAGTATAGATATGATGAT TTTTTTTTTG		2,300
	gB6 <sup>e</sup>	ATATATAAT ACAATAAGAATACGTGAGATAGTATAGATATGATGAT TTTTTTTTTA		222

	gB7 <sup>e</sup>	GATATATCTGAGTAAGTGATGACATTGTGAATTATTTTTT T		489
	gB7 <sup>e</sup>	AAATATAAG CAAATGATATATCTGAGTAAGTGATGACATTGTGAATT TTTTTTTTTTT	3,042	250
	gB7 <sup>e</sup>	ATGATATATCTGAGTAAGTGATGACATTGTGAATTAT GGTATATTTTG		117
	gB7 <sup>e</sup>	TGATATATCTGAGTAAGTGATGACATTGTGAATTACTTTTT T		49
	gB7 <sup>e</sup>	AATGATATATCTGAGTAAGTGATGACATTGTGAATTATTTTTT TTTTT		14
	gB7 <sup>e</sup>	AATATAAG CAAATGATATATCTGAGTAAGTGATGACATTGTGAATT TTTTTTTTTTTTTTTTTA AAAAAAAA	884	
	gB7 <sup>e</sup>	AAATATAAG CAAATGATATATCTGAGTAAGTGATGACATTGTGAATTAT GGTATATAAGTTAAATAATTATTC	522	
A	gA1	ATA TTAATCTAATAACGGAGACTTGTATATGATAGTAATGATGATATT TTTTTTTT		2,697
	gA1	ATATA TTAATCTAATAGCGGAGACTTGTATATGATAGTAATGATGATATT TT	104	
	gA2	ATATAA AATCTAATGACGAGAACTTGTATATGATGATGAAGGTGATA ATT TTTTTTTTTTTTTG		16,984
	gA2	ATA AATCTAATAACGAGAATTATGTACGATAATGAAAGTGATAT ATTTTTTTCTTTT		540
	gA2	ATATAA AATCTAATGACGAAAGCTTGTATATGGTAATGAAAGATGGTATT TTTTTTTTT		540
	gA2	ATATAA AATCTAATGACGAAAGCTTGTATATGGTAATGAAAGATGGTA ATT TTTTTTTCTTTT		419
	gA2	ATATAA AATCTAATGACGAGAACTTGTATATGATGATGAAGGTGATATT TCTTTTTTTT		331
	gA2	ATATA AATCTAATAACGAGAATTATGTACGATAATGAAAGTGATATT TTTTTCTTTTTTT		106
	gA2	ATACAT ATCTAATAACGGAAGCCTTGTATGGTAGTAGTGAAAGATGGTA ATT TTTTTTTTTTTT		103
	gA2	ATATAA AATCTAATGACGAAAGCTTGTATATGGTAGTGAAAGATGGTA ATT TTTTTTTTTTTT	1,234	
	gA2	ATATAA AATCTAATGACGAAAGCTTGTATATGGTAGTGAAAGATGGTATT TTTTTTTTT	382	
	gA2	ATATA AATCTAATAACGGAAATTGTATATGATGATAGAAGTGATAGT TTTTTTTTTTT	40	

## **REFERENCES**

## REFERENCES

1. Shapiro TA, Englund PT. The structure and replication of kinetoplast DNA. *Annu Rev Microbiol.* 1995;49: 117+.
2. Vickerman K. The evolutionary expansion of the trypanosomatid flagellates. *Int J Parasitol.* 1994;24: 1317–1331. doi:10.1016/0020-7519(94)90198-8
3. Simpson AGB, Stevens JR, Lukeš J. The evolution and diversity of kinetoplastid flagellates. *Trends Parasitol.* 2006;22: 168–174. doi:10.1016/j.pt.2006.02.006
4. Read LK, Lukeš J, Hashimi H. Trypanosome RNA editing: the complexity of getting U in and taking U out. *Wiley Interdiscip Rev RNA.* 2016;7: 33–51. doi:10.1002/wrna.1313
5. Priest JW, Hajduk SL. Developmental regulation of mitochondrial biogenesis in *Trypanosoma brucei*. *J Bioenerg Biomembr.* 1994;26: 179–191. doi:10.1007/BF00763067
6. Aphasizhev R, Aphasizheva I. Uridine insertion/deletion editing in trypanosomes: a playground for RNA-guided information transfer. *Wiley Interdiscip Rev RNA.* 2011;2: 669–685. doi:10.1002/wrna.82
7. Aphasizhev R, Aphasizheva I. Mitochondrial RNA editing in trypanosomes: Small RNAs in control. *Biochimie.* 2014;100: 125–131. doi:10.1016/j.biochi.2014.01.003
8. Hong M i. n., Simpson L. Genomic Organization of *Trypanosoma brucei* Kinetoplast DNA Minicircles. *Protist.* 2003;154: 265–279. doi:10.1078/143446103322166554
9. Koslowsky D, Sun Y, Hindenach J, Theisen T, Lucas J. The insect-phase gRNA transcriptome in *Trypanosoma brucei*. *Nucleic Acids Res.* 2014;42: 1873–1886. doi:10.1093/nar/gkt973
10. CDC - African Trypanosomiasis [Internet]. 2 May 2017 [cited 25 Jan 2019]. Available: <https://www.cdc.gov/parasites/sleepingsickness/index.html>
11. Programme Against African Trypanosomosis (PAAT) | Food and Agriculture Organization of the United Nations [Internet]. [cited 25 Jan 2019]. Available: <http://www.fao.org/paat/en/>
12. Horn D. Antigenic variation in African trypanosomes. *Mol Biochem Parasitol.* 2014;195: 123–129. doi:10.1016/j.molbiopara.2014.05.001
13. Sudarshi D, Lawrence S, Pickrell WO, Eligar V, Walters R, Quaderi S, et al. Human African Trypanosomiasis Presenting at Least 29 Years after Infection—What Can This Teach Us

about the Pathogenesis and Control of This Neglected Tropical Disease? PLOS Negl Trop Dis. 2014;8: e3349. doi:10.1371/journal.pntd.0003349

14. Matthews KR. Developments in the Differentiation of *Trypanosoma brucei*. Parasitol Today. 1999;15: 76–80. doi:10.1016/S0169-4758(98)01381-7
15. van Hellemond JJ, Bakker BM, Tielens AGM. Energy Metabolism and Its Compartmentation in *Trypanosoma brucei*. In: Poole RK, editor. Advances in Microbial Physiology. Academic Press; 2005. pp. 199–226. doi:10.1016/S0065-2911(05)50005-5
16. Hannaert V, Bringaud F, Opperdoes FR, Michels PA. Evolution of energy metabolism and its compartmentation in Kinetoplastida. Kinetoplastid Biol Dis. 2003;2: 11. doi:10.1186/1475-9292-2-11
17. Nolan DP, Voorheis HP. The mitochondrion in bloodstream forms of *Trypanosoma brucei* is energized by the electrogenic pumping of protons catalysed by the F1F0-ATPase. Eur J Biochem. 1992;209: 207–216. doi:10.1111/j.1432-1033.1992.tb17278.x
18. Vertommen D, Van Roy J, Szikora J-P, Rider MH, Michels PAM, Opperdoes FR. Differential expression of glycosomal and mitochondrial proteins in the two major life-cycle stages of *Trypanosoma brucei*. Mol Biochem Parasitol. 2008;158: 189–201. doi:10.1016/j.molbiopara.2007.12.008
19. Weelden SWH van, Fast B, Vogt A, Meer P van der, Saas J, Hellemond JJ van, et al. Procytic *Trypanosoma brucei* Do Not Use Krebs Cycle Activity for Energy Generation. J Biol Chem. 2003;278: 12854–12863. doi:10.1074/jbc.M213190200
20. Weelden SWH van, Hellemond JJ van, Opperdoes FR, Tielens AGM. New Functions for Parts of the Krebs Cycle in Procytic *Trypanosoma brucei*, a Cycle Not Operating as a Cycle. J Biol Chem. 2005;280: 12451–12460. doi:10.1074/jbc.M412447200
21. Bringaud F, Rivière L, Coustou V. Energy metabolism of trypanosomatids: Adaptation to available carbon sources. Mol Biochem Parasitol. 2006;149: 1–9. doi:10.1016/j.molbiopara.2006.03.017
22. Oberle M, Balmer O, Brun R, Roditi I. Bottlenecks and the Maintenance of Minor Genotypes during the Life Cycle of *Trypanosoma brucei*. PLOS Pathog. 2010;6: e1001023. doi:10.1371/journal.ppat.1001023
23. Abbeele JVD, Claes Y, Bockstaele DV, Ray DL, Coosemans M. *Trypanosoma brucei* spp. development in the tsetse fly: characterization of the post-mesocyclic stages in the foregut and proboscis. Parasitology. 1999;118: 469–478.
24. Michelotti EF, Hajduk SL. Developmental regulation of trypanosome mitochondrial gene expression. J Biol Chem. 1987;262: 927–932.

25. Feagin JE, Jasmer DP, Stuart K. Developmentally regulated addition of nucleotides within apocytochrome b transcripts in *Trypanosoma brucei*. *Cell*. 1987;49: 337–345. doi:10.1016/0092-8674(87)90286-8
26. Read LK, Wilson KD, Myler PJ, Stuart K. Editing of *Trypanosoma brucei* maxicircle CR5 mRNA generates variable carboxy terminal predicted protein sequences. *Nucleic Acids Res*. 1994;22: 1489–1495. doi:10.1093/nar/22.8.1489
27. Koslowsky DJ, Bhat GJ, Perrollaz AL, Feagin JE, Stuart K. The MURF3 gene of *T. brucei* contains multiple domains of extensive editing and is homologous to a subunit of NADH dehydrogenase. *Cell*. 1990;62: 901–911. doi:10.1016/0092-8674(90)90265-G
28. Souza AE, Myler PJ, Stuart K. Maxicircle CR1 transcripts of *Trypanosoma brucei* are edited and developmentally regulated and encode a putative iron-sulfur protein homologous to an NADH dehydrogenase subunit. *Mol Cell Biol*. 1992;12: 2100–2107. doi:10.1128/MCB.12.5.2100
29. Souza AE, Shu HH, Read LK, Myler PJ, Stuart KD. Extensive editing of CR2 maxicircle transcripts of *Trypanosoma brucei* predicts a protein with homology to a subunit of NADH dehydrogenase. *Mol Cell Biol*. 1993;13: 6832–6840. doi:10.1128/MCB.13.11.6832
30. Read LK, Myler PJ, Stuart K. Extensive editing of both processed and preprocessed maxicircle CR6 transcripts in *Trypanosoma brucei*. *J Biol Chem*. 1992;267: 1123–1128.
31. Bhat GJ, Koslowsky D, Feagin J, Smiley B, Stuart K. An Extensively Edited Mitochondrial Transcript in Kinetoplastids Encodes a Protein Homologous to ATPase Subunit 6. *Cell*. 1990;61: 885–894.
32. Feagin JE, Abraham JM, Stuart K. Extensive editing of the cytochrome c oxidase III transcript in *Trypanosoma brucei*. *Cell*. 1988;53: 413–422. doi:10.1016/0092-8674(88)90161-4
33. Blum B, Bakalara N, Simpson L. A model for RNA editing in kinetoplastid mitochondria: “Guide” RNA molecules transcribed from maxicircle DNA provide the edited information. *Cell*. 1990;60: 189–198. doi:10.1016/0092-8674(90)90735-W
34. Eperon IC, Janssen JWG, Hoeijmakers JHH, Borst P. The major transcripts of the kinetoplast DNA of *Trypanosoma brucei* are very small ribosomal RNAs. *Nucleic Acids Res*. 1983;11: 105–125. doi:10.1093/nar/11.1.105
35. Adler BK, Harris ME, Bertrand KI, Hajduk SL. Modification of *Trypanosoma brucei* mitochondrial rRNA by posttranscriptional 3' polyuridine tail formation. *Mol Cell Biol*. 1991;11: 5878–5884. doi:10.1128/MCB.11.12.5878

36. Aphasizheva I, Maslov D, Wang X, Huang L, Aphasizhev R. Pentatricopeptide Repeat Proteins Stimulate mRNA Adenylation/Uridylation to Activate Mitochondrial Translation in Trypanosomes. *Mol Cell*. 2011;42: 106–117. doi:10.1016/j.molcel.2011.02.021
37. Bhat GJ, Souza AE, Feagin JE, Stuart K. Transcript-specific developmental regulation of polyadenylation in *Trypanosoma brucei* mitochondria. *Mol Biochem Parasitol*. 1992;52: 231–240. doi:10.1016/0166-6851(92)90055-O
38. Hensgens LA, Brakenhoff J, De Vries BF, Sloof P, Tromp MC, Van Boom JH, et al. The sequence of the gene for cytochrome c oxidase subunit I, a frameshift containing gene for cytochrome c oxidase subunit II and seven unassigned reading frames in *Trypanosoma brucei* mitochondrial maxi-circle DNA. *Nucleic Acids Res*. 1984;12: 7327–7344.
39. Benne R, Van Den Burg J, Brakenhoff JPJ, Sloof P, Van Boom JH, Tromp MC. Major transcript of the frameshifted coxII gene from trypanosome mitochondria contains four nucleotides that are not encoded in the DNA. *Cell*. 1986;46: 819–826. doi:10.1016/0092-8674(86)90063-2
40. Feagin JE, Stuart K. Differential expression of mitochondrial genes between life cycle stages of *Trypanosoma brucei*. *Proc Natl Acad Sci*. 1985;82: 3380–3384.
41. Payne M, Rothwell V, Jasmer DP, Feagin JE, Stuart K. Identification of mitochondrial genes in *Trypanosoma brucei* and homology to cytochrome c oxidase II in two different reading frames. *Mol Biochem Parasitol*. 1985;15: 159–170. doi:10.1016/0166-6851(85)90117-3
42. Kannan S, Burger G. Unassigned MURF1 of kinetoplastids codes for NADH dehydrogenase subunit 2. *BMC Genomics*. 2008;9: 455. doi:10.1186/1471-2164-9-455
43. Feagin JE, Jasmer DP, Stuart K. Apocytochrome b and other mitochondrial DNA sequences are differentially expressed during the life cycle of *Trypanosoma brucei*. *Nucleic Acids Res*. 1985;13: 4577–4596. doi:10.1093/nar/13.12.4577
44. Stuart K, Feagin JE, Jasmer DP. Regulation of Mitochondrial Gene Expression in *Trypanosoma brucei*. Sequence Specificity in Transcription and Translation. Alan R. Liss; 1985. pp. 621–631.
45. Jasmer DP, Feagin JE, Stuart K. Diverse patterns of expression of the cytochrome c oxidase subunit I gene and unassigned reading frames 4 and 5 during the life cycle of *Trypanosoma brucei*. *Mol Cell Biol*. 1985;5: 3041–3047. doi:10.1128/MCB.5.11.3041
46. Feagin JE, Stuart K. Developmental aspects of uridine addition within mitochondrial transcripts of *Trypanosoma brucei*. *Mol Cell Biol*. 1988;8: 1259–1265. doi:10.1128/MCB.8.3.1259

47. Stuart K. The RNA editing process in *Trypanosoma brucei*. *Semin Cell Biol.* 1993;4: 251–260. doi:10.1006/scel.1993.1030
48. Corell RA, Myler P, Stuart K. *Trypanosoma brucei* mitochondrial CR4 gene encodes an extensively edited mRNA with completely edited sequence only in bloodstream forms. *Mol Biochem Parasitol.* 1994;64: 65–74. doi:10.1016/0166-6851(94)90135-X
49. Hajduk SL, Adler BK, Madison S, McManus M, Sabatini R. Insertional and deletional RNA editing in trypanosome mitochondria. *Nucleic Acids Symp Ser.* 1996; 15–18.
50. Koslowsky DJ, Riley GR, Feagin JE, Stuart K. Guide RNAs for transcripts with developmentally regulated RNA editing are present in both life cycle stages of *Trypanosoma brucei*. *Mol Cell Biol.* 1992;12: 2043–2049. doi:10.1128/MCB.12.5.2043
51. Riley GR, Corell RA, Stuart K. Multiple guide RNAs for identical editing of *Trypanosoma brucei* apocytochrome b mRNA have an unusual minicircle location and are developmentally regulated. *J Biol Chem.* 1994;269: 6101–6108.
52. Greif G, Rodriguez M, Reyna-Bello A, Robello C, Alvarez-Valin F. Kinetoplast adaptations in American strains from *Trypanosoma vivax*. *Mutat Res Mol Mech Mutagen.* 2015;773: 69–82. doi:10.1016/j.mrfmmm.2015.01.008
53. Ooi C-P, Schuster S, Cren-Travaillé C, Bertiaux E, Cosson A, Goyard S, et al. The Cyclical Development of *Trypanosoma vivax* in the Tsetse Fly Involves an Asymmetric Division. *Front Cell Infect Microbiol.* 2016;6. doi:10.3389/fcimb.2016.00115
54. Ruvalcaba-Trejo LI, Sturm NR. The *Trypanosoma cruzi* Sylvio X10 strain maxicircle sequence: the third musketeer. *BMC Genomics.* 2011;12: 58. doi:10.1186/1471-2164-12-58
55. Cazzulo JJ. Protein and amino acid catabolism in *Trypanosoma cruzi*. *Comp Biochem Physiol Part B Comp Biochem.* 1984;79: 309–320. doi:10.1016/0305-0491(84)90381-X
56. Cazzulo JJ. Intermediate metabolism in *Trypanosoma cruzi*. *J Bioenerg Biomembr.* 1994;26: 157–165. doi:10.1007/BF00763064
57. Tyler KM, Engman DM. The life cycle of *Trypanosoma cruzi* revisited. *Int J Parasitol.* 2001;31: 472–481. doi:10.1016/S0020-7519(01)00153-9
58. Cannata JJB, Cazzulo JJ. The aerobic fermentation of glucose by *Trypanosoma cruzi*. *Comp Biochem Physiol Part B Comp Biochem.* 1984;79: 297–308. doi:10.1016/0305-0491(84)90380-8
59. Sanchez-moreno M, Fernandez-becerra MC, Castilla-calvente JJ, Osuna A. Metabolic studies by <sup>1</sup>H NMR of different forms of *Trypanosoma cruzi* as obtained by “in vitro”

- culture. *FEMS Microbiol Lett.* 1995;133: 119–125. doi:10.1111/j.1574-6968.1995.tb07871.x
60. Prevention C-C for DC and. CDC - Leishmaniasis [Internet]. 16 Oct 2018 [cited 25 Jan 2019]. Available: <https://www.cdc.gov/parasites/leishmaniasis/index.html>
61. McConville MJ, Naderer T. Metabolic Pathways Required for the Intracellular Survival of Leishmania. *Annu Rev Microbiol.* 2011;65: 543–561. doi:10.1146/annurev-micro-090110-102913
62. Saunders EC, Souza DPD, Naderer T, Sernee MF, Ralton JE, Doyle MA, et al. Central carbon metabolism of Leishmania parasites. *Parasitology.* 2010;137: 1303–1313. doi:10.1017/S0031182010000077
63. Simpson L, Thiemann OH, Savill NJ, Alfonzo JD, Maslov DA. Evolution of RNA editing in trypanosome mitochondria. *Proc Natl Acad Sci.* 2000;97: 6986–6993. doi:10.1073/pnas.97.13.6986
64. Porcel BM, Denoeud F, Opperdoes F, Noel B, Madoui M-A, Hammarton TC, et al. The Streamlined Genome of Phytomonas spp. Relative to Human Pathogenic Kinetoplastids Reveals a Parasite Tailored for Plants. *PLOS Genet.* 2014;10: e1004007. doi:10.1371/journal.pgen.1004007
65. Nawathean P, Maslov DA. The absence of genes for cytochrome c oxidase and reductase subunits in maxicircle kinetoplast DNA of the respiration-deficient plant trypanosomatid Phytomonas serpens. *Curr Genet.* 2000;38: 95–103. doi:10.1007/s002940000135
66. Blum B, Simpson L. Guide RNAs in kinetoplastid mitochondria have a nonencoded 3' oligo(U) tail involved in recognition of the preedited region. *Cell.* 1990;62: 391–397. doi:10.1016/0092-8674(90)90375-O
67. Ochsenreiter T, Hajduk SL. Alternative editing of cytochrome c oxidase III mRNA in trypanosome mitochondria generates protein diversity. *EMBO Rep.* 2006;7: 1128–1133. doi:10.1038/sj.emboj.7400817
68. Ochsenreiter T, Anderson S, Wood ZA, Hajduk SL. Alternative RNA Editing Produces a Novel Protein Involved in Mitochondrial DNA Maintenance in Trypanosomes. *Mol Cell Biol.* 2008;28: 5595–5604. doi:10.1128/MCB.00637-08
69. Covello P, Gray M. On the evolution of RNA editing. *Trends Genet.* 1993;9: 265–268. doi:10.1016/0168-9525(93)90011-6
70. Gray MW. Evolutionary Origin of RNA Editing. *Biochemistry (Mosc).* 2012;51: 5235–5242. doi:10.1021/bi300419r

71. Gray MW, Lukeš J, Archibald JM, Keeling PJ, Doolittle WF. Irremediable Complexity? *Science*. 2010;330: 920–921. doi:10.1126/science.1198594
72. Stoltzfus A. On the Possibility of Constructive Neutral Evolution. *J Mol Evol*. 1999;49: 169–181. doi:10.1007/PL00006540
73. Stoltzfus A. Constructive neutral evolution: exploring evolutionary theory's curious disconnect. *Biol Direct*. 2012;7: 35. doi:10.1186/1745-6150-7-35
74. Leeder W-M, Hummel NFC, Göringer HU. Multiple G-quartet structures in pre-edited mRNAs suggest evolutionary driving force for RNA editing in trypanosomes. *Sci Rep*. 2016;6. doi:10.1038/srep29810
75. Speijer D. Evolutionary Aspects of RNA Editing. In: Göringer HU, editor. *RNA Editing*. Berlin, Heidelberg: Springer Berlin Heidelberg; 2008. pp. 199–227. doi:10.1007/978-3-540-73787-2\_10
76. Buhrman H, van der Gulik P, Severini S, Speijer D. A mathematical model of kinetoplastid mitochondrial gene scrambling advantage. *ArXiv13071163 Q-Bio*. 2013; Available: <http://arxiv.org/abs/1307.1163>
77. Thiemann OH, Maslov DA, Simpson L. Disruption of RNA editing in *Leishmania tarentolae* by the loss of minicircle-encoded guide RNA genes. *EMBO J*. 1994;13: 5689–5700.
78. Savill Nicholas J., Higgs Paul G. A theoretical study of random segregation of minicircles in trypanosomatids. *Proc R Soc Lond B Biol Sci*. 1999;266: 611–620. doi:10.1098/rspb.1999.0680
79. Lynch M, Bürger R, Butcher D, Gabriel W. The Mutational Meltdown in Asexual Populations. *J Hered*. 1993;84: 339–344. doi:10.1093/oxfordjournals.jhered.a111354
80. LaBar T, Adami C. Evolution of drift robustness in small populations. *Nat Commun*. 2017;8: 1012. doi:10.1038/s41467-017-01003-7
81. Muller HJ. The relation of recombination to mutational advance. *Mutat Res Mol Mech Mutagen*. 1964;1: 2–9. doi:10.1016/0027-5107(64)90047-8
82. Haigh J. The accumulation of deleterious genes in a population—Muller's Ratchet. *Theor Popul Biol*. 1978;14: 251–267. doi:10.1016/0040-5809(78)90027-8
83. Poon A, Otto SP. Compensating for Our Load of Mutations: Freezing the Meltdown of Small Populations. *Evolution*. 2000;54: 1467–1479. doi:10.1111/j.0014-3820.2000.tb00693.x

84. Whitlock MC. Fixation of New Alleles and the Extinction of Small Populations: Drift Load, Beneficial Alleles, and Sexual Selection. *Evolution*. 2000;54: 1855–1861. doi:10.1111/j.0014-3820.2000.tb01232.x
85. Normark S, Bergström S, Edlund T, Grundström T, Jaurin B, Lindberg FP, et al. Overlapping genes. *Annu Rev Genet*. 1983;17: 499–525.
86. Liang H. Decoding the dual-coding region: key factors influencing the translational potential of a two-ORF-containing transcript. *Cell Res*. 2010;20: 508–509. doi:10.1038/cr.2010.62
87. Belshaw R, Pybus OG, Rambaut A. The evolution of genome compression and genomic novelty in RNA viruses. *Genome Res*. 2007;17: 000–000. doi:10.1101/gr.6305707
88. Brandes N, Linial M. Gene overlapping and size constraints in the viral world. *Biol Direct*. 2016;11: 26. doi:10.1186/s13062-016-0128-3
89. Mouilleron H, Delcourt V, Roucou X. Death of a dogma: eukaryotic mRNAs can code for more than one protein. *Nucleic Acids Res*. 2016;44: 14–23. doi:10.1093/nar/gkv1218
90. Liang H, Landweber LF. A genome-wide study of dual coding regions in human alternatively spliced genes. *Genome Res*. 2006;16: 190–196. doi:10.1101/gr.4246506
91. Ribrioux S, Brüniger A, Baumgarten B, Seuwen K, John MR. Bioinformatics prediction of overlapping frameshifted translation products in mammalian transcripts. *BMC Genomics*. 2008;9: 122. doi:10.1186/1471-2164-9-122
92. Pallejà A, Harrington ED, Bork P. Large gene overlaps in prokaryotic genomes: result of functional constraints or mispredictions? *BMC Genomics*. 2008;9: 335. doi:10.1186/1471-2164-9-335
93. Chung W-Y, Wadhawan S, Szklarczyk R, Pond SK, Nekrutenko A. A First Look at ARFome: Dual-Coding Genes in Mammalian Genomes. *PLOS Comput Biol*. 2007;3: e91. doi:10.1371/journal.pcbi.0030091
94. Klemke M. Two overlapping reading frames in a single exon encode interacting proteins-- a novel way of gene usage. *EMBO J*. 2001;20: 3849–3860. doi:10.1093/emboj/20.14.3849
95. Yoshida H, Oku M, Suzuki M, Mori K. pXBP1(U) encoded in XBP1 pre-mRNA negatively regulates unfolded protein response activator pXBP1(S) in mammalian ER stress response. *J Cell Biol*. 2006;172: 565–575. doi:10.1083/jcb.200508145
96. Peleg O, Kirzhner V, Trifonov E, Bolshoy A. Overlapping Messages and Survivability. *J Mol Evol*. 2004;59: 520–527. doi:10.1007/s00239-004-2644-5

97. Sykes SE, Hajduk SL. Dual Functions of  $\alpha$ -Ketoglutarate Dehydrogenase E2 in the Krebs Cycle and Mitochondrial DNA Inheritance in *Trypanosoma brucei*. *Eukaryot Cell*. 2013;12: 78–90. doi:10.1128/EC.00269-12
98. Sykes S, Szempruch A, Hajduk S. The Krebs Cycle Enzyme  $\alpha$ -Ketoglutarate Decarboxylase Is an Essential Glycosomal Protein in Bloodstream African Trypanosomes. *Eukaryot Cell*. 2015;14: 206–215. doi:10.1128/EC.00214-14
99. Ochsenreiter T, Cipriano M, Hajduk SL. Alternative mRNA Editing in Trypanosomes Is Extensive and May Contribute to Mitochondrial Protein Diversity. *PLoS ONE*. 2008;3: e1566. doi:10.1371/journal.pone.0001566
100. Aphasizheva I, Maslov DA, Aphasizhev R. Kinetoplast DNA-encoded ribosomal protein S12. *RNA Biol*. 2013;10: 1679–1688. doi:10.4161/rna.26733
101. Stuart K, Gobright E, Jenni L, Milhausen M, Thomashow L, Agabian N. The Istar 1 Serodeme of *Trypanosoma brucei*: Development of a New Serodeme. *J Parasitol*. 1984;70: 747–754. doi:10.2307/3281757
102. Chomczynski P, Sacchi N. Single-step method of RNA isolation by acid guanidinium thiocyanate-phenol-chloroform extraction. *Anal Biochem*. 1987;162: 156–159. doi:10.1016/0003-2697(87)90021-2
103. Agabian N, Thomashow L, Milhausen M, Stuart K. Structural Analysis of Variant and Invariant Genes in Trypanosomes. *Am J Trop Med Hyg*. 1980;29: 1043–1049.
104. Friedman RC, Farh KK-H, Burge CB, Bartel DP. Most mammalian mRNAs are conserved targets of microRNAs. *Genome Res*. 2009;19: 92–105. doi:10.1101/gr.082701.108
105. Madina BR, Kumar V, Metz R, Mooers BHM, Bundschuh R, Cruz-Reyes J. Native mitochondrial RNA-binding complexes in kinetoplastid RNA editing differ in guide RNA composition. *RNA*. 2014; doi:10.1261/rna.044495.114
106. Clement SL, Mingler MK, Koslowsky DJ. An Intragenic Guide RNA Location Suggests a Complex Mechanism for Mitochondrial Gene Expression in *Trypanosoma brucei*. *Eukaryot Cell*. 2004;3: 862–869. doi:10.1128/EC.3.4.862-869.2004
107. Cristodero M, Seebeck T, Schneider A. Mitochondrial translation is essential in bloodstream forms of *Trypanosoma brucei*. *Mol Microbiol*. 2010;78: 757–769. doi:10.1111/j.1365-2958.2010.07368.x
108. MacLeod A, Turner CMR, Tait A. A high level of mixed *Trypanosoma brucei* infections in tsetse flies detected by three hypervariable minisatellites. *Mol Biochem Parasitol*. 1999;102: 237–248. doi:10.1016/S0166-6851(99)00101-2

109. Balmer O, Caccone A. Multiple-strain infections of *Trypanosoma brucei* across Africa. *Acta Trop.* 2008;107: 275–279. doi:10.1016/j.actatropica.2008.06.006
110. Szempruch AJ, Choudhury R, Wang Z, Hajduk SL. In vivo analysis of trypanosome mitochondrial RNA function by artificial site-specific RNA endonuclease-mediated knockdown. *RNA.* 2015; doi:10.1261/rna.052084.115
111. Surve S, Heestand M, Panicucci B, Schnaufer A, Parsons M. Enigmatic Presence of Mitochondrial Complex I in *Trypanosoma brucei* Bloodstream Forms. *Eukaryot Cell.* 2012;11: 183–193. doi:10.1128/EC.05282-11
112. Verner Z, Čermáková P, Škodová I, Kriegová E, Horváth A, Lukeš J. Complex I (NADH:ubiquinone oxidoreductase) is active in but non-essential for procyclic *Trypanosoma brucei*. *Mol Biochem Parasitol.* 2011;175: 196–200. doi:10.1016/j.molbiopara.2010.11.003
113. Speijer D. Is kinetoplastid pan-editing the result of an evolutionary balancing act? *IUBMB Life.* 2006;58: 91–96. doi:10.1080/15216540600551355
114. Hudson K m., Taylor AE r., Elce B j. Antigenic changes in *Trypanosoma brucei* on transmission by tsetse fly. *Parasite Immunol.* 1980;2: 57–69. doi:10.1111/j.1365-3024.1980.tb00043.x
115. Gibson W. The origins of the trypanosome genome strains *Trypanosoma brucei brucei* TREU 927, *T. b. gambiense* DAL 972, *T. vivax* Y486 and *T. congolense* IL3000. *Parasit Vectors.* 2012;5: 71.
116. Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, et al. Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Mol Syst Biol.* 2014;7: 539–539. doi:10.1038/msb.2011.75
117. Gonnet GH, Cohen MA, Benner SA. Exhaustive matching of the entire protein sequence database. *Science.* 1992;256: 1443–1445.
118. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in Python. *J Mach Learn Res.* 2011;12: 2825–2830.
119. Kirby LE, Sun Y, Judah D, Nowak S, Koslowsky D. Analysis of the *Trypanosoma brucei* EATRO 164 Bloodstream Guide RNA Transcriptome. *PLOS Negl Trop Dis.* 2016;10: e0004793. doi:10.1371/journal.pntd.0004793
120. Duarte M, Tomás AM. The mitochondrial complex I of trypanosomatids - an overview of current knowledge. *J Bioenerg Biomembr.* 2014;46: 299–311. doi:10.1007/s10863-014-9556-x

121. Simpson L, Neckelmann N, Cruz VF de la, Simpson AM, Feagin JE, Jasmer DP, et al. Comparison of the maxicircle (mitochondrial) genomes of *Leishmania tarentolae* and *Trypanosoma brucei* at the level of nucleotide sequence. *J Biol Chem.* 1987;262: 6182–6196.
122. Hanada K, Shiu S-H, Li W-H. The Nonsynonymous/Synonymous Substitution Rate Ratio versus the Radical/Conservative Replacement Rate Ratio in the Evolution of Mammalian Genes. *Mol Biol Evol.* 2007;24: 2235–2241. doi:10.1093/molbev/msm152
123. Firth AE, Brown CM. Detecting overlapping coding sequences with pairwise alignments. *Bioinformatics.* 2005;21: 282–292. doi:10.1093/bioinformatics/bti007
124. Firth AE, Brown CM. Detecting overlapping coding sequences in virus genomes. *BMC Bioinformatics.* 2006;7: 75. doi:10.1186/1471-2105-7-75
125. Landweber LF, Gilbert W. RNA editing as a source of genetic variation. *Nat Lond.* 1993;363: 179.
126. Tielens AGM, van Hellemond JJ. Surprising variety in energy metabolism within Trypanosomatidae. *Trends Parasitol.* 2009;25: 482–490. doi:10.1016/j.pt.2009.07.007
127. Verner Z, Čermáková P, Škodová I, Kováčová B, Lukeš J, Horváth A. Comparative analysis of respiratory chain and oxidative phosphorylation in *Leishmania tarentolae*, *Crithidia fasciculata*, *Phytomonas serpens* and procyclic stage of *Trypanosoma brucei*. *Mol Biochem Parasitol.* 2014;193: 55–65. doi:10.1016/j.molbiopara.2014.02.003
128. Jackson AP, Berry A, Aslett M, Allison HC, Burton P, Vavrova-Anderson J, et al. Antigenic diversity is generated by distinct evolutionary mechanisms in African trypanosome species. *Proc Natl Acad Sci U S A.* 2012;109: 3416–3421. doi:10.1073/pnas.1117313109
129. Morrison LJ, Vezza L, Rowan T, Hope JC. Animal African Trypanosomiasis: Time to Increase Focus on Clinically Relevant Parasite and Host Species. *Trends Parasitol.* 2016;32: 599–607. doi:10.1016/j.pt.2016.04.012
130. Maslov DA, Hollar L, Haghigat P, Nawathean P. Demonstration of mRNA editing and localization of guide RNA genes in kinetoplast–mitochondria of the plant trypanosomatid *Phytomonas serpens* Note: Nucleotide sequences from *P. serpens* 1G reported in this work were deposited in GenBank™ database with the following accession numbers: AF034624 (Sau3AI-cut minicircle), AF034625 (HindIII-cut minicircle), AF034626 (fully edited sequence of RPS12 mRNA), AF034627 (genomic sequence of RPS12 cryptogene).1. *Mol Biochem Parasitol.* 1998;93: 225–236. doi:10.1016/S0166-6851(98)00028-0
131. David V, Flegontov P, Gerasimov E, Tanifuji G, Hashimi H, Logacheva MD, et al. Gene Loss and Error-Prone RNA Editing in the Mitochondrion of Perkinsela, an Endosymbiotic Kinetoplastid. *mBio.* 2015;6: e01498-15. doi:10.1128/mBio.01498-15

132. Maslov DA. Complete set of mitochondrial pan-edited mRNAs in *Leishmania mexicana amazonensis* LV78. *Mol Biochem Parasitol.* 2010;173: 107–114.  
doi:10.1016/j.molbiopara.2010.05.013
133. Käll L, Krogh A, Sonnhammer ELL. A Combined Transmembrane Topology and Signal Peptide Prediction Method. *J Mol Biol.* 2004;338: 1027–1036.  
doi:10.1016/j.jmb.2004.03.016
134. Käll L, Krogh A, Sonnhammer ELL. Advantages of combined transmembrane topology and signal peptide prediction—the Phobius web server. *Nucleic Acids Res.* 2007;35: W429–W432. doi:10.1093/nar/gkm256
135. Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJE. The Phyre2 web portal for protein modeling, prediction and analysis. *Nat Protoc.* 2015;10: 845–858.  
doi:10.1038/nprot.2015.053
136. Xu Y, Tao Y, Cheung LS, Fan C, Chen L-Q, Xu S, et al. Structures of bacterial homologues of SWEET transporters in two distinct conformations. *Nature.* 2014;515: 448–452.  
doi:10.1038/nature13670
137. Bender T, Pena G, Martinou J-C. Regulation of mitochondrial pyruvate uptake by alternative pyruvate carrier complexes. *EMBO J.* 2015;34: 911–924.  
doi:10.15252/embj.201490197
138. Štáfková J, Mach J, Biran M, Verner Z, Bringaud F, Tachezy J. Mitochondrial pyruvate carrier in *Trypanosoma brucei*. *Mol Microbiol.* 2016;100: 442–456.  
doi:10.1111/mmi.13325
139. Saunders EC, Ng WW, Kloehn J, Chambers JM, Ng M, McConville MJ. Induction of a Stringent Metabolic Response in Intracellular Stages of *Leishmania mexicana* Leads to Increased Dependence on Mitochondrial Metabolism. *PLOS Pathog.* 2014;10: e1003888.  
doi:10.1371/journal.ppat.1003888
140. Aphasizheva I, Aphasizhev R. U-Insertion/Deletion mRNA-Editing Holoenzyme: Definition in Sight. *Trends Parasitol.* 2016;32: 144–156. doi:10.1016/j.pt.2015.10.004
141. Kirby LE, Koslowsky D. Mitochondrial dual-coding genes in *Trypanosoma brucei*. *PLoS Negl Trop Dis.* 2017;11: e0005989. doi:10.1371/journal.pntd.0005989
142. Lamour N, Rivière L, Coustou V, Coombs GH, Barrett MP, Bringaud F. Proline Metabolism in Procyclic *Trypanosoma brucei* Is Down-regulated in the Presence of Glucose. *J Biol Chem.* 2005;280: 11902–11910. doi:10.1074/jbc.M414274200
143. Coustou V, Biran M, Breton M, Guegan F, Rivière L, Plazolles N, et al. Glucose-induced Remodeling of Intermediary and Energy Metabolism in Procyclic *Trypanosoma brucei*. *J Biol Chem.* 2008;283: 16342–16354. doi:10.1074/jbc.M709592200

144. Bochud-Allemann N, Schneider A. Mitochondrial Substrate Level Phosphorylation Is Essential for Growth of Procyclic *Trypanosoma brucei*. *J Biol Chem.* 2002;277: 32849–32854. doi:10.1074/jbc.M205776200
145. Horváth A, Horáková E, Dunajčíková P, Verner Z, Pravdová E, Šlapetová I, et al. Downregulation of the nuclear-encoded subunits of the complexes III and IV disrupts their respective complexes but not complex I in procyclic *Trypanosoma brucei*. *Mol Microbiol.* 2005;58: 116–130. doi:10.1111/j.1365-2958.2005.04813.x
146. Gnipová A, Panicucci B, Paris Z, Verner Z, Horváth A, Lukeš J, et al. Disparate phenotypic effects from the knockdown of various *Trypanosoma brucei* cytochrome c oxidase subunits. *Mol Biochem Parasitol.* 2012;184: 90–98. doi:10.1016/j.molbiopara.2012.04.013
147. Kuile BH ter. Adaptation of metabolic enzyme activities of *Trypanosoma brucei* promastigotes to growth rate and carbon regimen. *J Bacteriol.* 1997;179: 4699–4705. doi:10.1128/jb.179.15.4699-4705.1997
148. Simpson RM, Bruno AE, Bard JE, Buck MJ, Read LK. High-throughput sequencing of partially edited trypanosome mRNAs reveals barriers to editing progression and evidence for alternative editing. *RNA.* 2016;22: 677–695. doi:10.1261/rna.055160.115
149. Carnes J, McDermott S, Anupama A, Oliver BG, Sather DN, Stuart K. In vivo cleavage specificity of *Trypanosoma brucei* editosome endonucleases. *Nucleic Acids Res.* 2017;45: 4667–4686. doi:10.1093/nar/gkx116
150. Otaka E, Hashimoto T, Mizuta K. The ribosomal proteins. I: An introduction to a compilation of the protein species equivalents from various organisms by a universal code system. *Protein Seq Data Anal.* 1993;5: 285–300.
151. Lawson SD, Igo RP, Salavati R, Stuart KD. The specificity of nucleotide removal during RNA editing in *Trypanosoma brucei*. *RNA.* 2001;7: 1793–1802.
152. Baradaran R, Berrisford JM, Minhas GS, Sazanov LA. Crystal structure of the entire respiratory complex I. *Nature.* 2013;494: 443–448. doi:10.1038/nature11871
153. Källberg M, Wang H, Wang S, Peng J, Wang Z, Lu H, et al. Template-based protein structure modeling using the RaptorX web server. *Nat Protoc.* 2012;7: 1511–1522. doi:10.1038/nprot.2012.085
154. Peng J, Xu J. A multiple-template approach to protein threading. *Proteins Struct Funct Bioinforma.* 2011;79: 1930–1939. doi:10.1002/prot.23016
155. Peng J, Xu J. Raptord: Exploiting structure information for protein alignment by statistical inference. *Proteins Struct Funct Bioinforma.* 2011;79: 161–171. doi:10.1002/prot.23175

156. Sturm NR, Maslov DA, Blum B, Simpson L. Generation of unexpected editing patterns in *Leishmania tarentolae* mitochondrial mRNAs: Misediting produced by misguiding. *Cell*. 1992;70: 469–476. doi:10.1016/0092-8674(92)90171-8
157. Maslov DA, Thiemann O, Simpson L. Editing and misediting of transcripts of the kinetoplast maxicircle G5 (ND3) cryptogene in an old laboratory strain of *Leishmania tarentolae*. *Mol Biochem Parasitol*. 1994;68: 155–159. doi:10.1016/0166-6851(94)00160-X
158. Alatortsev VS, Cruz-Reyes J, Zhelonkina AG, Sollner-Webb B. Trypanosoma brucei RNA Editing: Coupled Cycles of U Deletion Reveal Processive Activity of the Editing Complex. *Mol Cell Biol*. 2008;28: 2437–2445. doi:10.1128/MCB.01886-07
159. Necas D, Ohtamaa M, Määttä E, Haapala A. python-Levenshtein 0.12.0.
160. Zimmer SL, Simpson RM, Read LK. High throughput sequencing revolution reveals conserved fundamentals of U-indel editing. *Wiley Interdiscip Rev RNA*. 2018;9: e1487. doi:10.1002/wrna.1487
161. Simpson RM, Bruno AE, Chen R, Lott K, Tylec BL, Bard JE, et al. Trypanosome RNA Editing Mediator Complex proteins have distinct functions in gRNA utilization. *Nucleic Acids Res*. 2017;45: 7965–7983. doi:10.1093/nar/gkx458
162. Koslowsky DJ, Jayarama Bhat G, Read LK, Stuart K. Cycles of progressive realignment of gRNA with mRNA in RNA editing. *Cell*. 1991;67: 537–546. doi:10.1016/0092-8674(91)90528-7
163. Abraham JM, Feagin JE, Stuart K. Characterization of cytochrome c oxidase III transcripts that are edited only in the 3' region. *Cell*. 1988;55: 267–272. doi:10.1016/0092-8674(88)90049-9
164. Sturm NR, Simpson L. Partially edited mRNAs for cytochrome b and subunit III of cytochrome oxidase from leishmania tarentolae mitochondria: RNA editing intermediates. *Cell*. 1990;61: 871–878. doi:10.1016/0092-8674(90)90197-M
165. Decker CJ, Sollner-Webb B. RNA editing involves indiscriminate U changes throughout precisely defined editing domains. *Cell*. 1990;61: 1001–1011. doi:10.1016/0092-8674(90)90065-M
166. Ammerman ML, Presnyak V, Fisk JC, Foda BM, Read LK. TbRGG2 facilitates kinetoplastid RNA editing initiation and progression past intrinsic pause sites. *RNA*. 2010;16: 2239–2251. doi:10.1261/rna.2285510
167. Wang Z, Drew ME, Morris JC, Englund PT. Asymmetrical division of the kinetoplast DNA network of the trypanosome. *EMBO J*. 2002;21: 4998–5005. doi:10.1093/emboj/cdf482

168. Lukeš J, Skalický T, Týč J, Votýpká J, Yurchenko V. Evolution of parasitism in kinetoplastid flagellates. *Mol Biochem Parasitol*. 2014;195: 115–122.  
doi:10.1016/j.molbiopara.2014.05.007