

EXPLORING THE MOLECULAR EVOLUTION OF PROTEINS WITH DEEP  
MUTATIONAL SCANNING

By

Matthew Steven Faber

A DISSERTATION

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

Biochemistry and Molecular Biology – Doctor of Philosophy

2019

## ABSTRACT

### EXPLORING THE MOLECULAR EVOLUTION OF PROTEINS WITH DEEP MUTATIONAL SCANNING

By

Matthew Steven Faber

In this thesis, deep mutational scanning is expanded and applied to better understand the molecular evolution of proteins. Protein evolution is a complex process where subtle changes in molecular architecture can have massive impacts on biophysical properties, altering how well-adapted a protein is to a specific task or environment. Deep mutational scanning provides a finer level of understanding of molecular evolution by assessing the effect of every possible single-mutation on a protein's function. The technique combines site saturation mutant libraries, high throughput selections, and deep sequencing to tabulate the changes in mutant frequencies. From these changes the impacts of the mutations on protein function are characterized. This technology allows for efficient exploration of the local evolutionary landscape of a protein, making it a powerful tool for understanding evolution.

Here, I use deep mutational scanning to study how the initial likelihood of obtaining the native folded state of an enzyme *in vivo* constrains its evolution. We designed two unique single-point mutants of AmiE, an aliphatic amidase from *Pseudomonas aeruginosa*. These mutant enzymes are significantly less likely to reach the native folded state *in vivo* than the unmutated precursor and have catalytic efficiencies that are statistically indistinguishable from the initial unmutated enzyme. I tested the impacts of nearly all single-point mutations for the two impaired enzymes using high-throughput growth selections and compared them to the precursor enzyme. These comparisons provided insights into how evolutionary outcomes are changed following

decreases in the likelihood of native folding, and on how the impacts of single mutations combine to influence function.

The other primary goal of this thesis is the development of a new method that expands the utility of deep mutational scanning studies. This method assembles comprehensive single-site saturation, and large multi-point, mutant genome libraries of the bacteriophage  $\phi$ X174. To assemble the mutant genome libraries we combine nicking scanning mutagenesis and Golden Gate cloning. With these viral genome libraries, deep mutational scanning experiments can be performed *in situ*. These libraries are a valuable tool for studying the molecular determinants of viral host switching, the combination of inter- and intra-subunit mutations, and other aspects of the molecular evolution of viruses.

This thesis is dedicated to my family and friends who have inspired and encouraged me to pursue  
a life of exploration through science.

## **ACKNOWLEDGEMENTS**

I want to acknowledge my graduate advisor, Tim Whitehead, for his guidance, support, and patience in helping me mature my scientific practices. Our time working together has significantly impacted my character in many positive ways and imbued me with the skills necessary to succeed. I also want to acknowledge my labmates - Angelica Medina-Cucurella, Carolyn Haarmeyer, Caitlin Stein, Emily Wrenbeck, James Stapleton, Justin Klesmith, Matilda Newton, Matt Bedewitz, Monica Kirby, and PJ Steiner - for their support and friendship. I want to thank my parents and siblings for their unending love, support, and for always encouraging me to be a curious and creative individual. I want to thank my friends for the great times spent together, and for the good times still to be had. Finally, I want to thank my wife Fabiola for her constant love, positivity, support, and inspiration.

## TABLE OF CONTENTS

<b>LIST OF TABLES .....</b>	viii
<b>LIST OF FIGURES .....</b>	x
<b>KEY TO ABBREVIATIONS .....</b>	xii
<b>CHAPTER 1 An introduction to deep mutational scanning, its limitations, and products ..</b>	1
Abstract .....	2
Introduction .....	3
The technology .....	3
Library preparation .....	4
High throughput selections .....	5
The limitations .....	6
Experimental .....	6
Evolutionary insights .....	7
The products .....	9
Growth selection products .....	9
YSD-FACS products .....	12
REFERENCES .....	14
<b>CHAPTER 2 Data-driven engineering of protein therapeutics .....</b>	18
Abstract .....	19
Introduction .....	20
Large-scale mutational analysis .....	22
Antibody deep mutational scanning .....	22
Enzyme deep mutational scanning .....	24
Writing libraries of synthetic genes .....	25
Nature-sourced data .....	27
Protein design by phylogeny .....	27
Engineering from human antibody repertoires .....	28
Perspective .....	29
Acknowledgements .....	30
REFERENCES .....	31
<b>CHAPTER 3 Impact of <i>in vivo</i> protein folding probability on local fitness landscapes ..</b>	37
Abstract .....	38
Introduction .....	39
Results .....	42
AmiE variants with lower <i>in vivo</i> folding probabilities and wild-type catalytic efficiencies engineered .....	43
Deep mutational scans for AmiE variants .....	47

Distribution of beneficial fitness effects are largely insensitive to initial <i>in vivo</i> protein folding probability .....	49
Moderate epistasis observed with decreasing enzyme <i>in vivo</i> folding probability...	51
.....	.....
Most beneficial mutations in the WT background are shared .....	52
Positive sign epistasis is overwhelmingly specific .....	55
A plurality of unique beneficial mutations is codon-dependent .....	55
Discussion .....	57
Materials and methods .....	59
Reagents .....	59
Computational design of folding impaired mutants .....	59
Plasmid construction .....	60
Near comprehensive single-site mutant library construction .....	61
Protein expression and purification .....	61
Biophysical analysis of proteins .....	62
Growth selections .....	65
Deep mutational scanning .....	66
Data availability .....	67
Acknowledgements .....	68
REFERENCES .....	69
<b>CHAPTER 4 Saturation mutagenesis genome engineering of infective ΦX174 bacteriophage via unamplified oligo pools and golden gate assembly .....</b>	<b>74</b>
Abstract .....	75
Introduction .....	76
Results .....	77
Discussion .....	81
Materials and methods .....	83
Reagents .....	83
Segmentation of the ΦX174 genome .....	83
Introduction of nicking site .....	83
Comprehensive single site mutant library construction .....	84
Illumina sequencing prep and analysis .....	85
Assembly of mutant genomes .....	85
REFERENCES .....	87
<b>CHAPTER 5 Conclusions and perspectives .....</b>	<b>91</b>
REFERENCES .....	95
<b>APPENDICES .....</b>	<b>97</b>
APPENDIX A Chapter 3 supporting information .....	98
APPENDIX B Chapter 4 supporting information .....	145
APPENDIX C Purification of the TROP2 extracellular domain from a stable insect cell line using ammonium sulfate precipitation .....	161
REFERENCES.....	179

## LIST OF TABLES

Table 4.1: Mutant library NGS statistics .....	80
Table 4.2: Mutant genome assembly statistics .....	81
Table A 1: Relative fitness for I38V and I22L synonymous codons in the AmiE WT background . .....	122
Table A 2: Biophysical analysis of the AmiE variants .....	123
Table A 3: Thermal shift analysis data and statistics.....	124
Table A 4: Circular dichroism analysis of thermal denaturation statistics .....	125
Table A 5: DNA sequences of transcriptional elements in plasmids used for deep mutational scans.....	126
Table A 6: Summary of MG1655 rph+ transformants obtained during selection strain mutant library preparation.....	127
Table A 7: Summary of mutational library statistics .....	128
Table A 8: Goodness-of-fit test statistics for distribution fittings .....	129
Table A 9: Deleterious empirical cumulative distribution function (ECDF) analysis statistics..	130
Table A 10: Inner and outer primers for PCR reactions for Illumina sequencing .....	131
Table A 11: Beneficial mutations shared by all enzymes .....	133
Table A 12: Beneficial mutations shared by only AmiE WT and AmiE I122L .....	136
Table A 13: Beneficial mutations shared by only AmiE WT and AmiE I38V .....	137
Table A 14: Beneficial mutations shared by only AmiE I38V and AmiE I122L .....	138
Table A 15: AmiE WT unique beneficial mutations .....	139
Table A 16: AmiE I122L unique beneficial mutations.....	140
Table A 17: AmiE I38V unique beneficial mutations .....	141

Table A 18: AmiE WT unique beneficial mutations with synonymous codon fitness disparities ....	142
Table A 19: AmiE I122L unique beneficial mutations with synonymous codon fitness disparities.	143
Table A 20: AmiE I38V unique beneficial mutations with synonymous codon fitness disparities ..	144
Table B 1: Mutant library preparation summary .....	158
Table B 2: Primers for incorporating BbvCI nicking sites into the pCR2.1-topo shuttle vector.....	159
Table B 3: Inner and outer primers for PCR reactions for Illumina sequencing .....	160

## LIST OF FIGURES

Figure 1.1: Deep mutational scanning overview .....	4
Figure 2.1: High throughput techniques used to overcome therapeutic protein engineering bottlenecks .....	21
Figure 2.2: Big data yields insights for efficient engineering of protein therapeutics.....	29
Figure 3.1: Design of deep mutational scanning experiment .....	44
Figure 3.2: Local fitness landscapes are nearly insensitive to initial protein folding probability <i>in vivo</i> .....	50
Figure 3.3: Positive sign epistatic mutations are spatially segregated and specific.....	54
Figure 3.4: Unique beneficial mutations have high percentages of synonymous codon fitness disparities .....	56
Figure 4.1: $\Phi$ X174 mutant library assembly .....	79
Figure A 1: Chromatogram of purified AmiE variants.....	102
Figure A 2: Representative far-UV spectra of the folded and unfolded AmiE variants .....	103
Figure A 3: Locations of I38V and I122L mutations in AmiE quaternary structure.....	104
Figure A 4: AmiE WT tryptophan emission spectra as a function of GDN-HCl concentration..	105
Figure A 5: Thermal shift analysis of the purified AmiE variants .....	106
Figure A 6: Activity loss of AmiE WT following dilution.....	107
Figure A 7: Thermal denaturation monitored by far-UV circular dichroism .....	108
Figure A 8: Frequency distribution of pre-selection read counts for AmiE libraries .....	109
Figure A 9: Deep mutational scanning replicates for AmiE WT compared with previous literature .....	110
Figure A 10: Relative fitness metrics for AmiE proteins as a function of pre-selection read counts .....	111

Figure A 11: Normalized distribution of fitness effects for WT, I122L, and I38V backgrounds .....	112
Figure A 12: Analysis of proportions of deleterious mutations.....	113
Figure A 13: Venn diagram of shared and unique beneficial mutations using the strict cutoff ..	114
Figure A 14: Correlation analysis of linear regression normalized shared beneficial mutations.	115
Figure A 15: Dot plots of fitness effect synonymous codon variances for beneficial mutations.....	116
Figure A 16: SDS-PAGE analysis of the purity of the purified AmiE variants .....	117
Figure A 17: Comparison of AmiE I122L outlier mutation technical replicates .....	118
Figure A 18: AmiE WT heatmap.....	119
Figure A 19: AmiE I122L heatmap .....	120
Figure A 20: AmiE I38V heatmap.....	121
Figure B 1: Introduction of nicking sites into shuttle vectors containing viral genes .....	149
Figure B 2: F1 heatmap of counts.....	150
Figure B 3: F2 tile 1 heatmap of counts.....	151
Figure B 4: F2 tile 2 heatmap of counts.....	152
Figure B 5: F3 tile 1 heatmap of counts.....	153
Figure B 6: F3 tile 2 heatmap of counts.....	154
Figure B 7: G1 tile 1 heatmap of counts .....	155
Figure B 8: G1 tile 2 heatmap of counts .....	156
Figure B 9: G2 heatmap of counts .....	157
Figure C 1: TROP2Ex and its expression .....	165
Figure C 2: Purification of TROP2Ex from insect cell cultures .....	167
Figure C 3: Verification of the proper assembly of the m7E6 scFv yeast display construct and of binding to TROP2Ex .....	168

## **KEY TO ABBREVIATIONS**

AmiE, Amidase E

CDR, complementarity determining regions

DFE, distribution of fitness effects

DMS, deep mutational scanning

FACS, fluorescence activated cell sorting

FPLC-SEC, fast protein liquid chromatography with a size exclusion chromatography gel

IMAC, immobilized metal affinity chromatography

KS, Kolmogorov-Smirnoff

NSM, nicking scanning mutagenesis

scFv, single chain fragment of variable regions

V<sub>H</sub>-V<sub>L</sub>, variable heavy and variable light chains

TROP2Ex, tumor associated calcium signal transducer 2 extracellular domain

YSD, yeast surface display

## **CHAPTER 1**

**An introduction to deep mutational scanning, its limitations, and products**

## **Abstract**

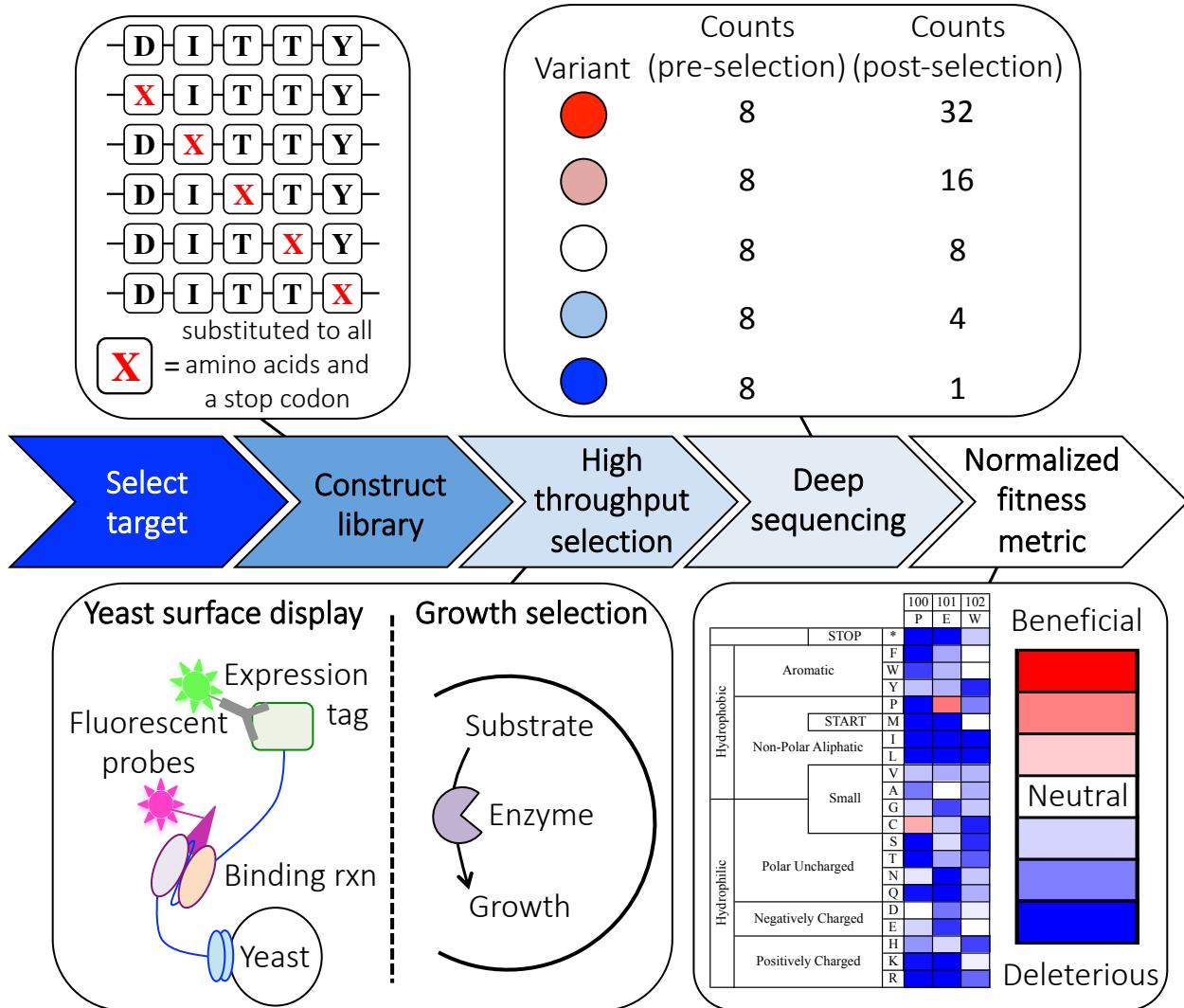
Deep mutational scanning combines saturation mutagenesis, high throughput selections, and deep sequencing to characterize the impacts of individual mutations on a respective protein. This technology allows for the characterization of thousands of mutations in parallel and as such has greatly expanded our abilities to understand protein evolution and perform rational design. This chapter will provide an introduction to deep mutational scanning, its methodological limitations, and what the outputs can and cannot tell us about the impacts of the tested mutations.

## **Introduction**

Deep mutational scanning (DMS) is a platform technology for efficiently assessing the impact of thousands of individual mutations on a protein of interest<sup>1,2</sup> (**Figure 1.1**). A DMS experiment requires a high throughput selection (HTS) for the protein of interest, the ability to generate comprehensive single-site saturation mutant libraries, and access to a deep sequencing platform. DMS is facilitated by clever cellular and molecular biology methods for HTS and for library preparations, as well as by the drop in deep sequencing costs to pennies per million base-pairs. The outputs of DMS are fitness metrics for each mutation that describe the magnitude by which a given mutation is beneficial or deleterious in the HTS performed. This chapter is not a comprehensive review of deep mutational scanning as there are already several excellent reviews that do so<sup>3,4</sup>. Additionally, this chapter will not focus on the application of the obtained datasets for studying evolution, or for forward engineering as these aspects have been reviewed in Wrenbeck et al.<sup>5</sup> and Faber and Whitehead<sup>6</sup> (Chapter 2). This chapter will provide a brief introduction to the technical aspects of DMS, what and where the limitations in the methods are, and what the obtained datasets can and cannot tell us about the mutations studied.

## **The technology**

As mentioned above DMS is the combination of three different processes: preparation of comprehensive site-saturation mutant libraries, high throughput selections, and deep sequencing (**Figure 1.1**). Here I will briefly describe advances in methods for preparing the mutant libraries, and the HTSs most used in our laboratory.



**Figure 1.1: Deep mutational scanning overview.** A DMS experiment begins with target selection, the selected target has to be compatible with a HTS. Next, comprehensive site-saturation mutant libraries are created. The HTS is applied to the mutant libraries to cause the weaker mutants to drop out, and beneficial mutations to proliferate. Counts of mutants in the pre- and post-selection libraries are tabulated using deep sequencing, and normalized fitness metrics are calculated based on the performance of the unmutated protein.

#### Library preparation

Single site saturation (SSM) mutagenesis libraries are constructed by template based mutagenesis using plasmid ssDNA and mixtures of mutagenic primers that encode the desired mutation(s). The mutant strand is ligated, the parent DNA strand digested, and a second

replication and ligation is performed<sup>7</sup>. Finally, the library is transformed into a cloning strain of *E. coli* to generate enough plasmid library for downstream applications. Initial library preparations using the Pfunkel method required preparation of ssDNA from bacteriophage and used uracil-containing DNA, which increased time and effort<sup>8</sup>. During my PhD studies my lab invented Nicking Mutagenesis, which allowed for the production of ssDNA by nicking a single strand of a dsDNA plasmid and digesting it with exonucleases<sup>8</sup>. This key advance allowed SSM library preparation from plasmid dsDNA in a single day. We have further advanced the production of these mutant libraries by replacing the degenerate mutagenic oligos with unamplified ink-jet printed oligo pools that contain user-defined mutations<sup>9</sup> (Chapter 5). Using unamplified oligo pools in place of degenerate oligos both simplifies the mutagenesis procedure and results in libraries with greater representations of all of the programmed mutations<sup>9</sup>.

### *High throughput selections*

The type of HTS used in a DMS experiment is dependent on the protein(s) being studied and the information one wants to obtain<sup>5,10</sup>. The two primary screens used in our laboratory are yeast surface display (YSD) paired with fluorescence activated cell sorting (FACS)<sup>9</sup>, and growth-based selections – where weaker variants become depleted and stronger variants enriched - for studying enzymes<sup>10-12</sup> (Chapter 3; **Figure 1.1**). Other selections exist, such as phage display in place of yeast display, lytic virus growth selections, or reporter based assays for enzyme function. Pairing YSD and FACS can be used to study protein-protein interactions, improve/engineer affinity and specificity<sup>13-16</sup>, and also to map paratopes and epitopes<sup>17-18</sup>. YSD and FACS can also be used to improve/study the stability of proteins and enzymes<sup>19-21</sup>. DMS experiments with growth-based selections can be used to inform the affinity/specificity

engineering of enzymes, and to improve enzymes for a respective reaction environment<sup>12</sup> (Chapter 3). Both YSD-FACS, and growth-based selections, can be used to answer molecular evolution questions because these methods map local fitness landscapes for a given environmental condition<sup>12,20</sup> (Chapter 3). Thus, DMS is a valuable tool for answering both applied and fundamental research questions.

## The limitations

### *Experimental*

DMS experiments require a suitable HTS, which significantly restricts the reach of these techniques. For YSD paired with FACS the desired target for scanning must be able to be displayed on the yeast cell surface, and the purified binding partners must be able to be produced<sup>10,17,18</sup>. Additionally, the reaction conditions have to be compatible with FACS and non-toxic to the displaying cells. Altered glycosylation patterns – inherent in YSD – can also impose significant issues. The mannose rich O- and N-linked glycosylation patterns in yeast are distinct from those in mammals, and cytosolic proteins are naturally unglycosylated. The attachment of non-native glycosylations can disrupt binding and alter other biophysical parameters for a protein. Also restricting the use of FACS is the oligomeric state of the protein being displayed and our group has successfully displayed up to homotrimers on surface of yeast<sup>17</sup>.

Growth-based HTSs require that cell growth be dependent on the reaction catalyzed by the enzyme of study. In Chapter 3 we study an amidase using DMS. When supplied with an amide the amidase releases ammonium, allowing for nitrogen restricted growth selections. Growth selections also require a growth rate ratio – growth rate in selective media over growth rate in unselective media – that is within a window of values<sup>10</sup>. Outside of the window the

growth of the cell is not proportional to the function of the enzyme. Tuning the growth rate ratio to be within this window requires certain biophysical parameters of enzymes like solubility, and catalytic efficiency to be met, as well as transcription and translational tuning with synthetic elements. Therefore, enzymes with very low catalytic efficiencies, and/or low solubilities, are unlikely to be compatible with DMS.

DMS is subject to the limits of the current mutant library preparation methods, cell sorting restrictions, and limitations in sequencing technologies<sup>9,10,18,22</sup>. The information obtained from DMS experiments is also restricted to a range of fitness metrics. DMS experiments performed using the pipelines developed in our lab can only scan an ~5-fold range of differences in fitness metrics<sup>20,23</sup>. DMS can discriminate sign differences with great accuracy and fidelity, but the ability to accurately quantify the impact of very deleterious mutations is dependent on sequencing depth<sup>10-12</sup> (Chapter 3). Often very deleterious mutations can only be qualitatively described.

When analyzing obtained datasets we are limited in the conclusion we can draw about the biophysical impacts of mutations. This limitation is typical in HTSs as these screens are often general competition assays that do not discriminate between the various biophysical properties that can be modified. Additional HTS screens<sup>19,20</sup> that select for a specific biophysical property, or experimental characterization of individual mutants<sup>12</sup> are required to determine how mutations impact the biophysical properties of the protein.

### *Evolutionary insights*

DMS datasets have to be approached with the understanding that these are limited test tube evolution experiments. Additionally, in any selection you get what you select for, meaning

that the conclusions and insights we obtain are biased for, and limited to, the selections performed<sup>10,12,20</sup> (Chapter 3). In Chapter 3, we perform our selections in a highly controlled environment at 37°C. However, our model enzyme - Amidase E from *P. aeruginosa* - has naturally evolved to function in a wide range of temperatures<sup>24</sup>. It is likely that mutations that are beneficial or deleterious at 37°C could have drastically different outcomes at different temperatures<sup>25</sup>. Additionally, our model enzyme is promiscuous and able to digest many short chain amides<sup>12,24</sup>, yet the selections we perform in Chapter 3 only use one substrate, acetamide. Therefore, we have selected for mutants that are only better at catalyzing the digestion of acetamide, or that increase the amount active enzyme *in vivo* at 37°C in the specific cellular host. It is possible that the beneficial mutations found in Chapter 3 will have decreased catalytic efficiencies with other substrates, and that some of the deleterious mutations found might be beneficial for digesting other aliphatic amides<sup>12</sup> (Chapter 3).

Selection conditions inherently constrain and bias the outcomes of test tube evolution experiments<sup>11,12,26,27</sup>. In laboratory evolution experiments the use of a single set of environmental conditions – which produce fitness values that can accurately describe mutational impacts only in the respective environment – is often required to limit the complexity of evolution such that it can be studied and understood<sup>27</sup>. For example, the Lenski evolution experiments have studied molecular and organismal evolution in several *E. coli* strains, individually, at a single temperature and constant chemical environment for over 60,000 generations<sup>28,29</sup>. The evolutionary insights obtained from this experiment are extremely valuable, yet they are biased for the organisms studied and the selection conditions used. By comparison, DMS approaches to studying molecular evolution are also subject to these limitations and biases. The environment used in a DMS experiment is user defined within certain parameters, and the selection

environment can be modified as long as the HTS is maintained<sup>10</sup>. However, altering the conditions of the HTS to test other environments – for example by performing growth-based selections at several different temperatures – can be extremely challenging and may not be feasible.

Finally, DMS is currently unable to model evolution over thousands of years and in highly complex environments. In nature there are many different environments a single organism can pass through, and many other organisms and pathogens to compete with, making the process of evolution massively more complex than in our laboratory experiments<sup>30-33</sup>. Therefore, while the data obtained from DMS is very informative, we must never lose sight of the fact that what is true in the test tube is an oversimplification of what is possibly occurring in nature.

## The products

Deep mutational scans provide local fitness landscapes typically defined as nearly all single mutational steps in the evolutionary space<sup>1,2</sup> (**Figure 1.1**). These can be used to gain insights into molecular evolution<sup>12,20</sup> (Chapter 3) and can be used in forward engineering proteins toward a desired purpose<sup>5,6</sup> (Chapter 2). As described previously, the definition of fitness is dependent on the HTS used and the environment of the assay. The obtained normalized fitness metrics describe how well a respective mutation is able to compete with the unmutated predecessor protein in the respective HTS<sup>10</sup>.

### *Growth selection products*

With growth based selections for studying enzymes, the normalized fitness metric describes how well a cell hosting a respective mutant is able to compete with – replicate faster

than – a cell hosting the unmutated enzyme in the assaying conditions at a respective temperature<sup>10</sup>. The mathematics for calculating normalized fitness metrics for growth based selections – as in Chapter 3 - from deep sequencing data was published by Kowalsky et al.<sup>10</sup> and is as follows. First, the frequency of a respective mutant ( $i$ ) in the pre- ( $f_{oi}$ ) and post-selection ( $f_{fi}$ ) populations is calculated with the tabulated counts of mutants from the deep sequencing data:

$$f_{oi} = \frac{x_{oi}}{\sum x_{oi}} \quad (1)$$

$$f_{fi} = \frac{x_{fi}}{\sum x_{fi}} \quad (2)$$

Where  $x_{oi}$  and  $x_{fi}$  are the number of counts of a specific variant in the pre- and post-selection populations respectively, and  $\sum x_{oi}$  and  $\sum x_{fi}$  are the total number of counts for all variants in the pre- and post-selection populations respectively. Next, the enrichment ratio,  $\varepsilon_i$ , for a mutant in the population is calculated from frequencies of the mutants obtained in (1) and (2):

$$\varepsilon_i = \log_2 \left( \frac{f_{fi}}{f_{oi}} \right) \quad (3)$$

The equation for calculating the enrichment ratios can be rewritten as:

$$\varepsilon_i = \log_2 \left( \frac{x_{fi}}{x_{oi}} \right) - \log_2 \left( \frac{\sum x_{fi}}{\sum x_{oi}} \right) \quad (4)$$

We also write the equation for the specific growth rate ( $\mu_i$ ) (5) of a cell as:

$$\mu_i = \ln \left( \frac{x_{fi}}{x_{oi}} \right) \frac{1}{t} \quad (5)$$

Where  $t$  is the difference in time between the end of the selection ( $f$ ) and initiation of the selection ( $o$ ). Combining the equations (4) and (5) gives:

$$\mu_i \log_2 e = \frac{1}{t} (\varepsilon_i + \log_2 \left( \frac{\sum x_{fi}}{\sum x_{oi}} \right)) \quad (6)$$

To simplify the equation we can first generate the average doubling period ( $g_p$ ) as a function of the change in the total number of counts for all variants in the pre- and post-selection populations:

$$g_p = \text{Number of doublings} = \log_2\left(\frac{\sum x_{fi}}{\sum x_{oi}}\right) \quad (7)$$

We can also remove  $t$  through redefining it:

$$t = \frac{\ln 2 * g_p}{\mu_p} \quad (8)$$

Where  $\mu_p$  is the bulk average growth rate of the population for the time between the initiation and end of the selection. Next we combine (7) and (8) into (6) to make the growth rate of a mutant a function of its respective enrichment ratio:

$$\mu_i = \mu_p \left( \frac{\varepsilon_i}{g_p} + 1 \right) \quad (9)$$

Finally, we can calculate our normalized fitness metric ( $\zeta_i$ ) where the growth rate of the respective mutant protein (i) is normalized to that of the unmutated predecessor protein (wt):

$$\zeta_i = \log_2 \left( \frac{\mu_i}{\mu_{wt}} \right) \quad (10)$$

We can also set the normalized fitness metric as a function of the enrichment ratios and average doubling period:

$$\zeta_i = \log_2 \left( \frac{\left( \frac{\varepsilon_i}{g_p} + 1 \right)}{\left( \frac{\varepsilon_{wt}}{g_p} + 1 \right)} \right) \quad (11)$$

The normalized fitness metric provides a quantitative description of how well a cell harboring a mutant  $i$  is able to compete with a cell hosting the predecessor enzyme in the respective environment. This data does not provide information on how the mutation impacts the biophysics of the enzyme.

For enzyme-based DMS, there are two general biophysical factors that are impacted by mutations and that determine fitness outcomes: specific velocity, and the amount of active enzyme being expressed within the cell. Both of these general properties are composed of a variety of different factors. Specific velocity is dependent on the standard Gibbs free energy of the reaction being catalyzed, substrate and product concentrations, flux of upstream and downstream reactions, the Michaelis constant ( $K_M$ ), and the maximum reaction velocity ( $k_{cat}$ )<sup>25,34</sup>. The probability of the enzyme being expressed and properly folded is a function of: the folding rate, the fidelity of the association process for oligomeric proteins, and the thermodynamic stability of the tertiary and quaternary structures and intermediates<sup>35,36</sup>. A growth selection is unable to discriminate between any of the above-mentioned factors, and additional biophysical analysis is required to understand why a mutation is deleterious or beneficial.

### *YSD-FACS products*

This chapter will not present the mathematics used in calculating the normalized fitness metrics for DMS experiments with YSD-FACS, those calculations can be found in the referenced work<sup>10,17</sup>. While I have performed DMS with YSD-FACS for epitope mapping<sup>17</sup>, paratope mapping (results incorporated into a patent in preparation), and affinity maturation (results incorporated into a patent in preparation), these projects are not included in this thesis. Normalized fitness metrics obtained from YSD-FACS experiments provide a quantitative measure of how well a mutant protein competes with its unmutated ancestor in either binding to a given target, or in displaying on the surface of the yeast in a respective environment<sup>10,17,20</sup>. Unlike the growth-based selections, YSD-FACS experiments can select for two different general biophysical properties: affinity/specificity, and surface display (**Figure 1.1**).

Comparing the obtained fitness landscapes from selections for binding/display with those from selections only for display allows for the discrimination of mutations that impact stability/solubility from those that impact affinity/specifcity<sup>17,19,20</sup>. Selections for display can be used to select for more stable variants, and to allow us to identify mutations that alter folding probabilities. To aid in discriminating the impacts of the mutations Kowalsky et al.<sup>17</sup> calculated Shannon entropy metrics for each position in the protein being studied for both the displaying sorted and binding/display sorted populations. These Shannon entropy metrics provide a quantitative measure of how well the respective position tolerates mutations. Positions with the lowest Shannon entropy scores are the least able to tolerate mutations. Comparison of Shannon entropy scores from selections for display with selections for binding/display allows for the identification of residues critical for the binding reaction. This comparative analysis is the foundation for mapping protein-protein interactions<sup>17,18</sup> with DMS. As described in Chapter 2, comparison of DMS datasets from different HTS with the same target protein allows for greater discrimination of the biophysical implications for respective mutations, and how this data can be used for forward engineering.

## **REFERENCES**

## REFERENCES

1. Fowler DM, Araya CL, Fleishman SJ, Kellogg EH, Stephany JJ, Baker D, Fields S: **High-resolution mapping of protein sequence–function relationships.** *Nat Methods* 2010, 7:741-746.
2. Hietpas RT, Jensen JD, Bolon DNA: **Experimental illumination of a fitness landscape.** *Proc Natl Acad Sci U S A* 2011, 108:7896-7901.
3. Araya CL, Fowler DM (2011) **Deep mutational scanning: Assessing protein function on a massive scale.** *Trends Biotechnol* 29(9):435–442.
4. Fowler DM, Fields S: **Deep mutational scanning: a new style of protein science.** *Nat Methods* 2014, 11:801–807.
5. Wrenbeck EE, Faber MS, Whitehead TA: **Deep sequencing methods for protein engineering and design.** *Curr Opin Struct Biol* 2017, 45:36-44. Wrenbeck 2019
6. Faber MS, Whitehead TA: **Data-driven engineering of protein therapeutics.** *Curr Opin in Biotech* 2019, 60:104-110.
7. Firnberg E, Ostermeier M: **PFunkel: efficient, expansive, user- defined mutagenesis.** *PLoS ONE* 2012, 7:e52031.
8. Wrenbeck EE, Klesmith JR, Stapleton JA, Adeniran A, Tyo KEJ, Whitehead TA: **Plasmid-based one-pot saturation mutagenesis.** *Nat Methods* 2016, 13:928-930.
9. Medina-Cucurella AV, Steiner PJ, Faber MS, Beltrán J, Borelli A, Kirby MB, Cutler SR, Whitehead TA: **User-defined single pot mutagenesis using unpurified oligo pools.** *Prot Eng. Des. Sel.* 2019, TBD:TBD.
10. Kowalsky CA, Klesmith JR, Stapleton JA, Kelly V, Reichkitzer N, Whitehead TA: **High-Resolution Sequence-Function Mapping of Full-Length Proteins.** *PLoS One* 2015, 10:e0118193.
11. Klesmith JR, Bacik JP, Michalczyk R, Whitehead TA: **Comprehensive Sequence-Flux Mapping of a Levoglucosan Utilization Pathway in *E. coli*.** *ACS Synth Biol* 2015, 4:1235–1243.
12. Wrenbeck EE, Azouz LR, Whitehead TA: **Single-mutation fitness landscapes for an enzyme on multiple substrates reveal specificity is globally encoded.** *Nat Comm* 2017, 8:15695.

13. Whitehead TA, Chevalier A, Song Y, Dreyfus C, Fleishman SJ, Mattos CD, Myers CA, Kamisetty H, Blair P, Wilson IA, Baker D: **Optimization of affinity, specificity and function of designed influenza inhibitors using deep sequencing.** *Nat Biotechnol* 2012, 30:543-548.
14. Strauch E-M, Fleishman SJ, Baker D: **Computational design of a pH-sensitive IgG binding protein.** *Proc Natl Acad Sci U S A* 2014, 111:675-680.
15. Procko E, Berguig GY, Shen BW, Song Y, Frayo S, Convertine AJ, Margineantu D, Booth G, Correia BE, Cheng Y et al.: **A computationally designed inhibitor of an Epstein-Barr viral Bcl-2 protein induces apoptosis in infected cells.** *Cell* 2014, 157:1644-1656.
16. Koenig P, Lee CV, Walters BT, Janakiraman V, Stinson J, Patapoff TW, Fuh G: **Mutational landscape of antibody variable domains reveals a switch modulating the interdomain conformational dynamics and antigen binding.** *Proc Natl Acad Sci U S A* 2017, 114:E486-E495.
17. Kowalsky CA, Faber MS, Nath A, Dann HE, Kelly VW, Liu L, Shanker P, Wagner EK, Maynard JA, Chan C et al.: **Rapid fine conformational epitope mapping using comprehensive mutagenesis and deep sequencing.** *J Biol Chem* 2015, 290:26457-26470.
18. Medina-Cucurella AV, Whitehead TA: **Characterizing Protein-Protein Interactions Using Deep Sequencing Coupled to Yeast Surface Display.** *Met in Mol Biol* 2018, 1764:101-121.
19. Julian MC, Li L, Garde S, Wilen R, Tessier PM: **Efficient affinity maturation of antibody variable domains requires co-selection of compensatory mutations to maintain thermodynamic stability.** *Sci Rep* 2017, 7:45259.
20. Klesmith JR, Bacik JP, Wrenbeck EE, Michalczyk R, Whitehead TA: **Trade-offs between enzyme fitness and solubility illuminated by deep mutational scanning.** *Proc Natl Acad Sci U S A* 2017, 114:2265-2270.
21. Wrenbeck EE, Bedewitz MA, Klesmith JR, Noshin S, Barry CS, Whitehead TA: **An automated data-driven pipeline for improving heterologous enzyme expression.** *ACS Synth Biol* 2019, 15:474-481.
22. Glanville J, D'Angelo S, Khan TA, Reddy ST, Naranjo L, Ferrara F, Bradbury ARM: **Deep sequencing in library selection projects: what insight does it bring?** *Curr Opin Struct Biol* 2015, 33:146-160.
23. Tokuriki N, Stricher F, Schymkowitz J, Serrano L, Tawfik DS: **The stability effects of protein mutations appear to be universally distributed.** *J mol biol* 2007, 369:1318-1332.
24. Kelly M, Clarke PH: **An inducible amidase produced by a strain of *Pseudomonas aeruginosa*.** *J gen Microbiol* 1962, 27:305-316.

25. Somero GN: **Adaptation of enzymes to temperature: searching for basic “strategies”**. *Comp Biochem and Physio* 2004, 139:321-333.
26. Bloom JD, Labthavikul ST, Otey CR, Arnold FH: **Protein stability promotes evolvability**. *Proc Natl Acad Sci U S A* 2006, 103:5869–5874.
27. Lenski RE: **Experimental evolution and the dynamics of adaptation and genome evolution in microbial populations**. *The ISME J* 2017, 11:2181-2194.
28. Wiser MJ, Ribeck N, Lenski RE: **Long-term dynamics of adaptation in asexual populations**. *Science* 2013, 342:1364-1367.
29. Good BH, McDonald MJ, Barrick JE, Lenski RE, Desai MM: **The dynamics of molecular evolution over 60,000 generations**. *Nature* 2017, 551:45-50.
30. Holmes EC: **What can we predict about viral evolution and emergence?** *Curr Opin in Vir* 2013, 3:180-184.
31. Neher RA: **Genetic draft, selective interference, and population genetics of rapid adaptation**. *Annu Rev Ecol Evol Syst* 2013, 44:195-215.
32. Lieberman TD, Flett KB, Yelin I, Martin TR, McAdam AJ, Priebe GP, Kishony R: **Genetic variation of a bacterial pathogen within individuals with cystic fibrosis provides a record of selective pressures**. *Nature Gen* 2014, 46:82-87.
33. Liao M, Somero GN, Dong Y: **Comparing mutagenesis and simulations as tools for identifying functionally important sequence changes for protein thermal adaptation**. *Proc Natl Acad Sci U S A* 2019, 116:679-688.
34. Siddiqui KS: **Defying the activity–stability trade-off in enzymes: taking advantage of entropy to enhance activity and thermostability**. *Critical Reviews in Biotechnology* 2017, 37:309-322.
35. Baker D, Agard DA: **Kinetics versus thermodynamics in protein folding**. *Biochemistry* 1994, 33(24):7505-7509.
36. Shakhnovich EI: **Theoretical studies of protein-folding thermodynamics and kinetics**. *Curr Opin in Struc Biol* 1997, 7:29-40.

## CHAPTER 2

### **Data-driven engineering of protein therapeutics**

This chapter is adapted with permission from the article “Data-driven engineering of protein therapeutics” in *Current Opinion in Biotechnology* 60:104–110 by Matthew S. Faber and Timothy A. Whitehead. Copyright 2019 Elsevier Ltd.

## **Abstract**

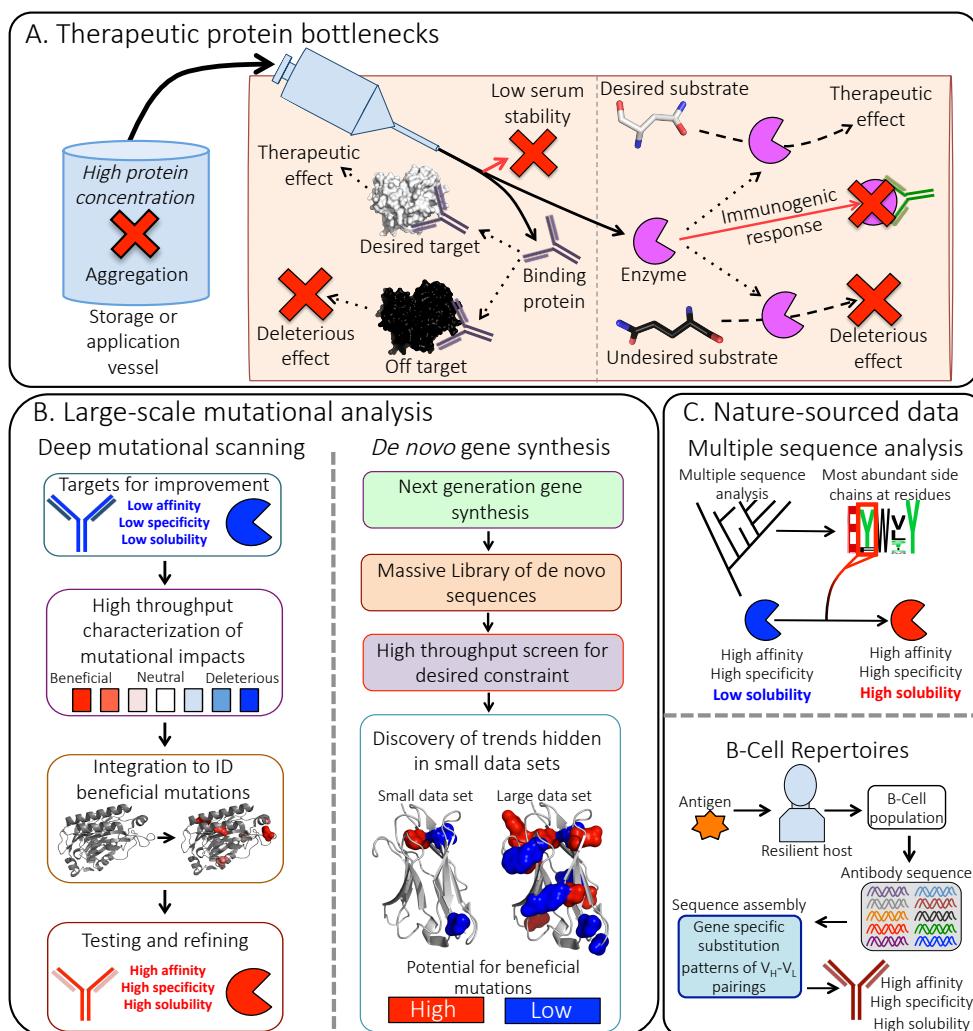
Protein therapeutics requires a series of properties beyond biochemical activity, including serum stability, low immunogenicity, and manufacturability. Mutations that improve one property often decrease one or more of the other essential requirements for therapeutic efficacy, making the protein engineering challenge difficult. The past decade has seen an explosion of new techniques centered around cheaply reading and writing DNA. This review highlights the recent use of such high throughput technologies for engineering protein therapeutics. Examples include the use of human antibody repertoire sequence data to pair antibody heavy and light chains, comprehensive mutational analysis for engineering antibody specificity, and the use of ancestral and inter-species sequence data to engineer simultaneous improvements in enzyme catalytic efficiency and stability. We conclude with a perspective on further ways to integrate mature protein engineering pipelines with the exponential increases in the volume of sequencing data expected in the forthcoming decade.

## **Introduction**

Over 100 therapeutic proteins are currently FDA-approved for use as drugs, with functions as variegated as enzymes that rob tumors of necessary nutrients to antibodies that block signaling mechanisms important for clinical presentation of rheumatoid arthritis. Regardless of the exact category of protein, there are a large number of requirements necessary for therapeutic efficacy: the biologic must possess sufficient affinity or catalytic efficiency toward its intended target while minimizing non-specific interactions, harbor a sufficient in vivo half-life, and maintain low immunogenicity and aggregation propensity. Additionally, each protein must satisfy a number of manufacturability constraints such as the capability of expression and purification in sufficient quantities, sufficient storage stability, and so forth. Naturally occurring proteins usually are suboptimal in at least one of these properties. While alterations to these properties can be imparted by mutations (protein engineering), mutations are pleiotropic and often a property can only be improved at the expense of another important requirement. Such competing constraints on protein function render brute force screens – used extensively by contemporary protein engineers – inefficient; therefore, new insights into navigating sequence space efficiently are always of interest.

This past decade has seen the cost of reading and writing DNA drop precipitously<sup>1</sup>. These advances have led to completely new ways to interrogate biology, including construction of thousands of synthetic genes<sup>2</sup>, evaluation of functional effect of tens of thousands of mutations in a protein in massively parallel experiments<sup>3</sup>, sequencing of thousands of homologues across the tree of life for any given gene, and even ways to sequence human antibody repertoires<sup>4</sup>. Such data-rich reservoirs are beginning to be tapped by protein engineers.

Here, we review recent progress in the use of data-driven protein engineering of potential therapeutics in order to satisfy multiple competing constraints. We restrict this review to protein engineering by amino acid substitutions. There are topical reviews concerning the de-immunization of protein therapeutics<sup>5</sup>, engineering competing trade-offs in antibody function<sup>6</sup>, and a comprehensive overview of engineered therapeutic enzymes<sup>7</sup>. Thus, we will highlight representative and instructive examples utilizing de novo gene synthesis, deep mutational scanning, phylogenetic analysis, or antibody repertoire sequencing (**Figure 2.1**); we also offer



**Figure 2.1: High throughput techniques used to overcome therapeutic protein engineering bottlenecks.** A. Different bottlenecks encountered in developing therapeutics. B–C. Different high throughput techniques used for protein engineering.

perspective on future ways to harness the extraordinary torrents of data that we anticipate in the coming decade.

### **Large-scale mutational analysis**

Integrating deep sequencing with user-defined saturation mutagenesis<sup>8</sup> and screens such as yeast surface display sorting enables the determination of the functional effect of tens of thousands of mutations on a given protein sequence in a massively parallel fashion; this technique is known as deep mutational scanning<sup>3</sup>. The mutational effect on each functional property of the protein can be assessed separately (if a screen exists), and mutations can be incorporated only if all screens give a positive result. Here, we give the latest examples of leveraging such datasets to engineer proteins while considering multiple constraints.

### **Antibody deep mutational scanning**

Antibodies often bind two or more related proteins, and frequently the protein engineer is tasked with engineering specificity for one protein over the other. A recent notorious example involved the use of an antibody to monitor age-dependent levels of GDF11 in mice<sup>9</sup>. However, this antibody also bound the closely related homolog GDF8 (myostatin), which resulted in erroneous conclusions<sup>10</sup> in the original paper. This multi-specificity problem is general, as many potential targets have very similar homologs in humans with some differing by only 10% or so in pairwise sequence identity<sup>11-14</sup>; antibodies, therefore, must have exquisite specificity for therapeutic and diagnostic applications.

Deep sequencing combined with a suitable screen can be used to identify specificity-modulating mutations in protein binders<sup>15</sup>. A recent excellent example comes from Koenig et al.

from Genentech<sup>16</sup>, who sought to prevent binding to angiopoietin-1 (Ang-1) for a candidate antigen binding fragment (Fab) with desired binding to Ang-2 and vascular endothelial growth factor (VEGF). Using phage display, the authors determined the relative binding profile of many possible single point mutants in the six complementarity determining regions (CDRs) (483 light chain, 609 heavy chain) for each protein in parallel. From these 1092 mutations, 25 (2.3%) were shown in vitro to result in no or severely reduced binding to Ang-1. We note that all of these mutations slightly decreased binding affinities for Ang-2 and/or VEGF, illustrating the general difficulty – and in some cases, impossibility – of finding specificity-modulating mutations without a functional trade-off even when the local mutational space is comprehensively sampled.

Thermal stability is another parameter that is often negatively correlated with binding affinity<sup>17</sup>. Reduced Fab thermal stability results in lower expression titers<sup>18</sup>. The necessity of engineering specificity-affinity while maintaining stability is critical for therapeutic application of binding proteins<sup>19</sup>, and state of the art methods involve using co-screening for stability and affinity simultaneously<sup>20</sup>. Alternatively, one can screen for stability and affinity in parallel to uncover rare, globally optimal mutations. Recently, the same Genentech team used this approach to scan all possible Fab framework and CDR single point mutants for affinity and stable expression in a phage display context<sup>21</sup>. Affinity-enhancing mutations were identified using selections against VEGF, while stabilizing mutations were identified with selections against protein A or protein L. There were a handful of mutations improving both stability and affinity. In particular, mutation from Phe to Ala at a single framework mutation 25 Å from the binding site, light chain residue 83 (LC-F83A), was found to strongly improve thermostability and affinity. The authors attributed this improvement to alteration of the interface between the variable and constant regions on the light chain, which was confirmed by hydrogen–deuterium

exchange mass spectrometry. Intriguingly, LC-83 is one of the 5–10% highest somatically mutated LC positions as determined by deep sequencing of over a thousand human lymphoid tissues. This suggests that LC-83 mutations can generally improve thermostability in Fabs, which the authors confirmed by incorporating LC-F83A in several unrelated mAbs. This is a great example of integrating naturally sourced deep sequencing data on a single target to obtain general insights into the functional tradeoffs of mutations in a broader set of antibodies.

### **Enzyme deep mutational scanning**

Enzymes are applied in the treatment of diverse disorders such as cancer therapeutics to deplete essential amino acids or metabolic precursors needed by the cancer cells<sup>22,23</sup>, or as replacement therapies to return specific metabolites to healthier levels in a patient<sup>24</sup>. Engineering non-immunogenic, serum-stable, and active enzymes can be challenging. For example, human kynureneine-degrading enzyme has low serum stability, and engineering has proven difficult because mutations that increase catalytic efficiency have decreased serum stability (J. Blazeck, personal communication). Understanding the trade-offs between stability and catalytic efficiency is critical for simplifying enzyme engineering.

In contrast with protein binders such as antibodies, generalizable high-throughput screens for directly testing enzyme function do not exist. Thus, deep mutational scanning pipelines are not available for all enzyme classes. A creative solution for forward engineering stability and activity was presented by Klesmith et al., who presented the comparative analysis of deep mutational scanning datasets for solubility/stability and activity for two different model enzymes: levoglucosan kinase, and TEM-1 beta-lactamase<sup>25</sup>. From these datasets, Klesmith et al. were able to extract features common to mutations that improve stability while not hampering

catalytic activity. These mutations were sampled in the evolutionary history of the enzyme, were at least 15 angstroms from the active site, were in positions of the tertiary protein sequence without a large number of other residues in close contact, and were not mutations to or from proline. Combining these filtering metrics yields a greater than 90% probability of choosing a stabilizing, catalytically neutral mutation in any given enzyme. This method has been developed into a Rosetta-based script and applied to increase the thermal stability and *in vivo* expression yield of a Type III polyketide synthase (E. Wrenbeck and T. Whitehead, unpublished results).

### **Writing libraries of synthetic genes**

Compared with more limited datasets of proteins with 1 or 2 amino acid changes described above, datasets with wider tranches of sequence space may be more useful for identifying general engineering rules. For example, several enzyme discovery efforts evaluate candidates by synthesizing hundreds of genes encoding enzymes spread throughout the protein superfamily<sup>26,27</sup>. Over the past five years the cost of writing DNA for kb-size genes has plateaued at about \$0.10 per base pair<sup>28</sup>, which means that libraries of several hundred genes can be written and tested on a medium-sized lab's budget.

A stunning example of the use of large protein datasets for obtaining engineering insights was recently described by Adimab scientists<sup>18</sup>. Monoclonal antibodies represent the largest class of engineered therapeutic proteins by revenue, with \$98 billion in worldwide sales in 2017<sup>29</sup>. Even with dozens of FDA-approved antibodies, some still fail in late-stage clinical trials for reasons unrelated to their binding affinity to their respective target. Understanding how a primary sequence determines manufacturability could help guide antibody-drug development akin to the Lipinski ‘rule of five’ edict for small molecule drugs<sup>30</sup>. To that end, Adimab scientists produced

a set of nearly all mAbs commercialized or in advanced stage clinical trials described in the patent literature through mid-2009 (137 in total)<sup>18</sup>. They then assayed this set for a dozen biophysical properties including aggregation propensity, non-specific binding, and melting temperature. The first major insight was that – contrary to expectations – many of these antibodies had unfavorable scores in one or more measured biophysical property. The second major insight was the significantly better biophysical properties of the approved subset of mAbs compared to those in Phase 2 clinical trials, suggesting potential difficulties in translating some of these Phase 2 mAbs.

Although this mAb dataset was only published in 2017, researchers are already mining it for predicting the sequence determinants of poor biophysical properties. For example, the Tessier group has used the Adimab dataset to find that positive CDR net charge correlates with mAb self-association using simple sequence-based scoring methods<sup>31</sup>. This same group, following work on individual antibodies<sup>32</sup>, has shown that reducing CDR net charge also decreases antibody nonspecificity (P. Tessier, personal communication).

Going from hundreds to tens of thousands of proteins – with greater underlying predictive power – requires a continued decrease in the cost of writing DNA. To that end, recent advances in DNA synthesis using microarray-derived oligo pools<sup>1</sup> have been utilized to synthesize thousands of designed small proteins<sup>2,33</sup>. Producing longer genes of 0.8-2 kb in length from such oligo pools is difficult. The best approaches have a success rate of ~2%<sup>34</sup>, and increasing the fidelity rate is essential for these emerging techniques to reach wider application.

## Nature-sourced data

### *Protein design by phylogeny*

The cumulative sequencing data of thousands of protein homologs across the tree of life, for nearly all known therapeutically relevant proteins, are a rich resource for the protein engineer. It has been known for decades that, given any position in a protein family, the consensus residue is likely to be stabilizing<sup>35</sup>. This insight enables the engineering of stable enzymes by simultaneously incorporating mutations at consensus positions in a protein sequence. An interesting variation on this idea comes from Nguyen et al., who used consensus sequence information to refine the specificity of an L-asparaginase from *Erwinia chrysanthemi* (ErA)<sup>36</sup>. The therapeutic effect of wild type ErA is diminished by its dual activity as an L-glutaminase. While asparagine depletion starves cancer cells, glutamine depletion is associated with many negative side effects<sup>37</sup>. To reduce the L-glutaminase activity of ErA, conserved stretches of residues were identified within the active site. Nonconserved adjacent residues were targeted for saturation mutagenesis with the assumption that such mutations could modulate specificity without disrupting desired catalytic function. This strategy resulted in a multi-point mutant with conserved L-asparaginase activity and a ~25-fold reduction of L-glutaminase activity. In vivo experiments using the engineered ErA have determined it is an effective therapy for both T and B-cell acute lymphoblastic leukemia without the side effects of glutamine depletion<sup>38</sup>.

Consensus mutations informed by phylogeny are increasingly incorporated into structure-based computational design algorithms to engineer stabilized proteins<sup>39</sup>. To cite as the best example, Goldenzweig et al. developed a Rosetta-based method to stabilize human enzymes<sup>40</sup>. Evolutionary conservation data were converted to a position specific scoring matrix, and only mutations above a certain conservation threshold were considered. Mutations passing additional

structure-based filters were combined into new designs. They used this approach to increase the expression yield and stability of a human acetylcholinesterase variant, which can be used for organophosphate detoxification following nerve agent exposure. Use of this method led to an acetylcholinesterase variant with 51 mutations, \$2000-fold increased bacterial expression compared to the original enzyme, and with nearly the same specific activity.

The consensus mutation approach has reached its logical end-state with the use of ancestral sequence reconstruction to improve protein function<sup>41-43</sup>. Resurrecting ancient proteins often leads to enzymes with incredible stabilities<sup>44</sup>. Zakas et al. used this general understanding to resurrect several ancestral Factor VIII proteins<sup>45</sup>. By doing so, they engineered a therapeutically effective ancestral Factor VIII variant with higher stability than extant forms.

Immunogenicity is always a major concern for candidate therapeutic proteins. For example, a highly engineered Factor VII variant failed late stage clinical trials because of immunogenicity<sup>46</sup>. Interestingly, the reconstructed ancestral Factor VIII variants share 95% sequence identity with human Factor VIII, but have reduced cross-reactivity with known anti-human Factor VIII antibodies<sup>45</sup>. This reduced cross-reactivity was verified in in vitro analyses. It will be interesting to see if these in vitro experiments translate to lower immunogenicity in animal models, and whether ancestral proteins have lower immunogenicity than extant proteins in general.

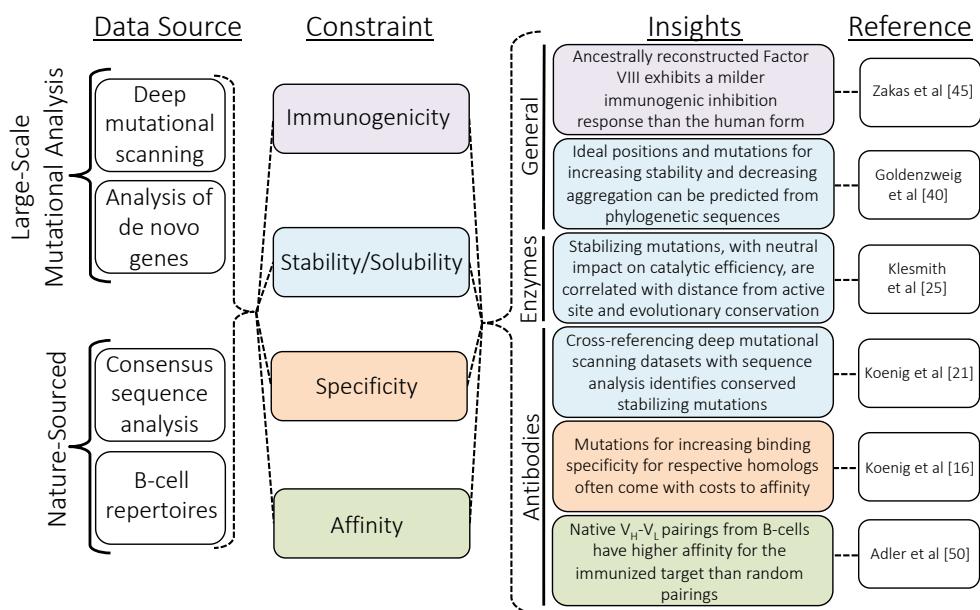
#### *Engineering from human antibody repertoires*

We are in an exciting age of molecular serology, moving on from bulk measurements of humoral responses towards descriptions of full sequences and functions of individual antibodies in human repertoires<sup>47-49</sup>. These repertoire sets are a critical resource for understanding affinity

maturity, CDR development, and pairing between heavy and light variable chains ( $V_H$ - $V_L$ ). While there has not been a great deal of engineering work in the open literature exploiting these repertoire datasets, we expect that to change in the next few years. The most recent use of information from these repertoires was described from Adler et al.<sup>50</sup>, who showed that incorporating native  $V_H$ - $V_L$  pairing improved the quality and quantity of candidate anti-interleukin 21 receptor antibodies isolated from a yeast display screen. Anticipated use of these datasets include designing more efficient antibody libraries by incorporating information from CDR lengths and net charges, and position specific substitution patterns. To facilitate this goal, Sheng et al. constructed gene-specific substitution patterns for 69 common V genes<sup>51</sup>, while Kovaltsuk et al. have collected nearly all repertoire datasets into a single data-base called the Observed Antibody Space<sup>52</sup> (antibodymap.org).

## Perspective

The ability to write and read DNA at scale has transformed protein engineering into a big data



**Figure 2.2: Big data yields insights for efficient engineering of protein therapeutics.**

field. Here, we have discussed a diversity of data-rich methods for improving protein therapeutics, with each of these methods providing unique engineering insights (**Figure 2.2**). In the next few years, we anticipate further ways to leverage and combine data from multiple sources to improve protein engineering. One of the key emerging ideas is to enhance protein design by incorporating constraints from the evolutionary history of the protein. These constraints are mostly identification of conservation at a single position, but for human antibodies with millions of sequences we imagine more sophisticated algorithms that take into account correlation between residues within or across CDRs, bulk biophysical properties like CDR net charge, and sequence properties as a function of CDR loop length. We also anticipate a marriage of data analysis from a post-hoc view by combining nature-sourced data of existing protein sequences with the forward-evolutionary view of deep mutational scanning to observe the range of possible mutations for a given candidate therapeutic. This integration will allow a comprehensive look at ‘what is’ with ‘what ought to be’, enabling the protein engineer to design functional proteins on demand.

## Acknowledgements

This work was supported by the National Science Foundation [Career Award #1254238 CBET to T.A.W] and Michigan State University [Johansen Crosby endowed chair to T.A.W].

## **REFERENCES**

## REFERENCES

1. Hughes RA, Ellington AD: **Synthetic DNA synthesis and assembly: putting the synthetic in synthetic biology.** *Cold Spring Harb Perspect Biol* 2017, 9:a023812.
2. Rocklin GJ, Chidyausiku TM, Goreshnik I, Ford A, Houliston S, Lemak A, Carter L, Ravichandran R, Mulligan VK, Chevalier A et al.: **Global analysis of protein folding using massively parallel design, synthesis, and testing.** *Science* 2017, 357:168-175.
3. Wrenbeck EE, Faber MS, Whitehead TA: **Deep sequencing methods for protein engineering and design.** *Curr Opin Struct Biol* 2017, 45:36-44.
4. Friedensohn S, Khan TA, Reddy ST: **Advanced methodologies in high-throughput sequencing of immune repertoires.** *Trends Biotechnol* 2017, 35:203-214.
5. Sauna ZE, Lagassé D, Pedras-Vasconcelos J, Golding B, Rosenberg AS: **Evaluating and mitigating the immunogenicity of therapeutic proteins.** *Trends Biotechnol* 2018, 36:1068-1084.
6. Rabia LA, Desai AA, Jhajj HS, Tessier PM: **Understanding and overcoming trade-offs between antibody affinity, specificity, stability and solubility.** *Biochem Eng J* 2018, 137:365-374.
7. Lutz S, Williams E, Muthu P: **Engineering therapeutic enzymes.** *Directed Enzyme Evolution: Advances and Applications.* Springer International Publishing; 2017:17-67.
8. Wrenbeck EE, Klesmith JR, Stapleton JA, Adeniran A, Tyo KEJ, Whitehead TA: **Plasmid-based one-pot saturation mutagenesis.** *Nat Methods* 2016, 13:928-930.
9. Sinha M, Jang YC, Oh J, Khong D, Wu EY, Manohar R, Miller C, Regalado SG, Francesco SL, Pancoast JR et al.: **Restoring systemic GDF11 levels reverses age-related dysfunction in mouse skeletal muscle.** *Science* 2014, 344:649-652.
10. Egerman MA, Cadena SM, Gilbert JA, Meyer A, Nelson HN, Swalley SE, Mallozzi C, Jacobi C, Jennings LJ, Clay I et al.: **GDF11 increases with age and inhibits skeletal muscle regeneration.** *Cell Metab* 2015, 22:164-174.
11. McPherron AC, Huynh TV, Lee SJ: **Redundancy of myostatin and growth/differentiation factor 11 function.** *BMC Dev Bio* 2009, 9:24-33.
12. Fredriksson R, Lagerström MC, Lundin LG, Schiöth HB: **The G-protein-coupled receptors in the human genome form five main families. Phylogenetic analysis, paragon groups, and fingerprints.** *Mol Pharmacol* 2003, 63:1256-1272.

13. Achen MG, Jeltsch M, Kukk E, Mäkinen T, Vitali A, Wilks AF, Alitalo K, Stacker SA: **Vascular endothelial growth factor D (VEGF-D) is a ligand for the tyrosine kinases VEGF receptor 2 (Flk1) and VEGF receptor 3 (Flt4).** *Proc Natl Acad Sci U S A* 1998, 95:548-553.
14. Fredriksson L, Li H, Eriksson U: **The PDGF family: four gene products form five dimeric isoforms.** *Cytokine Growth Factor Rev* 2004, 15:197-204.
15. Whitehead TA, Chevalier A, Song Y, Dreyfus C, Fleishman SJ, Mattos CD, Myers CA, Kamisetty H, Blair P, Wilson IA, Baker D: **Optimization of affinity, specificity and function of designed influenza inhibitors using deep sequencing.** *Nat Biotechnol* 2012, 30:543-548.
16. Koenig P, Sanowar S, Lee CV, Fuh G: **Tuning the specificity of a two-in-one fab against three angiogenic antigens by fully utilizing the information of deep mutational scanning.** *MAbs* 2017, 6:959-967.
17. Houlihan G, Gatti-Lafranconi P, Lowe D, Hollfelder F: **Directed evolution of anti-HER2 DARPinS by SNAP display reveals stability/function trade-offs in the selection process.** *Protein Eng Des Sel* 2015, 28:269-279.
18. Jain T, Sun T, Durand S, Hall A, Houston NR, Nett JH, Sharkey B, Bobrowicz B, Caffry I, Yu Y et al.: **Biophysical properties of the clinical-stage antibody landscape.** *Proc Natl Acad Sci U S A* 2017, 114:944-949.
19. Koday MT, Nelson J, Chevalier A, Koday M, Kalinoski H, Stewart L, Carter L, Nieusma T, Lee PS, Ward AB et al.: **A computationally designed hemagglutinin stem-binding protein provides in vivo protection from influenza independent of a host immune response.** *PLoS Pathog* 2016, 12:e1005409.
20. Julian MC, Li L, Garde S, Wilen R, Tessier PM: **Efficient affinity maturation of antibody variable domains requires co-selection of compensatory mutations to maintain thermodynamic stability.** *Sci Rep* 2017, 7:45259.
21. Koenig P, Lee CV, Walters BT, Janakiraman V, Stinson J, Patapoff TW, Fuh G: **Mutational landscape of antibody variable domains reveals a switch modulating the interdomain conformational dynamics and antigen binding.** *Proc Natl Acad Sci U S A* 2017, 114:E486-E495.
22. Triplett TA, Garrison KC, Marshall N, Donkor M, Blazeck J, Lamb C, Qerqez A, Dekker JD, Tanno Y, Lu WC et al.: **Reversal of indoleamine 2,3-dioxygenase-mediated cancer immune suppression by systemic kynurenone depletion with a therapeutic enzyme.** *Nat Biotechnol* 2018, 36:758-764.
23. Knott SRV, Wagenblast E, Khan S, Kim SY, Soto M, Wagner M, Turgeon MO, Fish L, Erard N, Gable AL et al.: **Asparagine bioavailability governs metastasis in a model of breast cancer.** *Nature* 2018, 554:378-381.

24. Germain DP, Charrow J, Desnick RJ, Guffon N, Kempf J, Lachmann RH, Lemay R, Linthorst GE, Packman S, Scott CR et al.: **Ten-year outcome of enzyme replacement therapy with agalsidase beta in patients with Fabry disease.** *J Med Genet* 2015, 52:353-358.
25. Klesmith JR, Bacik JP, Wrenbeck EE, Michalczyk R, Whitehead TA: **Trade-offs between enzyme fitness and solubility illuminated by deep mutational scanning.** *Proc Natl Acad Sci U S A* 2017, 114:2265-2270.
26. Carlin DA, Caster RW, Wang X, Betzenderfer SA, Chen CX, Duong VM, Ryklansky CV, Alpekin A, Beaumont N, Kapoor H et al.: **Kinetic characterization of 100 glycoside hydrolase mutants enables the discovery of structural features correlated with kinetic constants.** *PLoS One* 2016, 11:e0147596.
27. Macdonald SS, Patel A, Larmour VLC, Morgan-Lang C, Hallam SJ, Mark BL, Withers SG: **Structural and mechanistic analysis of a b-glycoside phosphorylase identified by screening a metagenomic library.** *J Biol Chem* 2018, 293:3451-3467.
28. Servick K: **Genome writing project confronts technology hurdles.** *Science* 2017, 356:673-674.
29. Grilo AL, Mantalaris A: **The increasingly human and profitable monoclonal antibody market.** *Trends Biotechnol* 2018, 37:9-16 <http://dx.doi.org/10.1016/j.tibtech.2018.05.014>.
30. Lipinski CA: **Lead-and drug-like compounds: the rule-of-five revolution.** *Drug Discov Today Technol* 2004, 1:337-341.
31. Alam ME, Geng SB, Bender C, Ludwig SD, Linden L, Hoet R, Tessier PM: **Biophysical and sequence-based methods for identifying monovalent and bivalent antibodies with high colloidal stability.** *Mol Pharm* 2018, 15:150-163.
32. Datta-Mannan A, Thangaraju A, Leung D, Tang Y, Witcher DR, Lu J, Wroblewski VJ: **Balancing charge in the complementarity-determining regions of humanized mAbs without affecting pI reduces non-specific binding and improves the pharmacokinetics.** *MAbs* 2015, 7:483-493.
33. Chevalier A, Silva DA, Rocklin GJ, Hicks DR, Vergara R, Murapa P, Bernard SM, Zhang L, Lam KH, Yao G et al.: **Massively parallel de novo protein design for targeted therapeutics.** *Nature* 2017, 550:74-79.
34. Plesa C, Sidore AM, Lubock NB, Zhang D, Kosuri S: **Multiplexed gene synthesis in emulsions for exploring protein functional landscapes.** *Science* 2018, 359:343-347.
35. Lehmann M, Pasamontes L, Lassen S, Wyss M: **The consensus concept for thermostability engineering of proteins.** *Biochim Biophys Acta* 2000, 1543:408-415.

36. Nguyen HA, Su Y, Lavie A: **Design and characterization of *Erwinia Chrysanthemi* L-asparaginase variants with diminished L-glutaminase activity.** *J Biol Chem* 2016, 291:17664-17676.
37. Vidya J, Sajitha S, Ushasree MV, Sindhu R, Binod P, Madhavan A, Pandey A: **Genetic and metabolic engineering approaches for the production and delivery of L-asparaginases: an overview.** *Bioresour Technol* 2017, 245:1775-1781.
38. Nguyen HA, Su Y, Zhang JY, Antanasićević A, Caffrey M, Schalk AM, Liu L, Rondelli D, Oh A, Mahmud DL et al.: **A novel L-asparaginase with low L-glutaminase coactivity is highly efficacious against both T and B cell acute lymphoblastic leukemias in vivo.** *Cancer Res* 2018, 6:1549-1560.
39. Bednar D, Beerens K, Sebestova E, Bendl J, Khare S, Chaloupkova R, Prokop Z, Brezovsky J, Baker D, Damborsky J: **FireProt: energy-and evolution-based computational design of thermostable multiple-point mutants.** *PLoS Comput Biol* 2015, 11:e1004556.
40. Goldenzweig A, Goldsmith M, Hill SE, Gertman O, Laurino P, Ashani Y, Dym O, Unger T, Albeck S, Prilusky J et al.: **Automated structure-and sequence-based design of proteins for high bacterial expression and stability.** *Mol Cell* 2016, 63:337-346.
41. Trudeau DL, Kaltenbach M, Tawfik DS: **On the potential origins of the high stability of reconstructed ancestral proteins.** *Mol Biol Evol* 2016, 33:2633-2641.
42. Lazarus RA, Scheiflinger F: **Mining ancient proteins for next-generation drugs.** *Nat Biotechnol* 2017, 35:28-29.
43. Risso VA, Sanchez-Ruiz JM, Ozkan SB: **Biotechnological and protein-engineering implications of ancestral protein resurrection.** *Curr Opin Struct Biol* 2018, 51:106-115.
44. Romero-Romero ML, Risso VA, Martinez-Rodriguez S, Ibarra-Molero B, Sanchez-Ruiz JM: **Engineering ancestral protein hyperstability.** *Biochem J* 2016, 473:3611-3620.
45. Zakas PM, Brown HC, Knight K, Meeks SL, Spencer HT, Gaucher EA, Doering CD: **Enhancing the pharmaceutical properties of protein drugs by ancestral sequence reconstruction.** *Nat Biotechnol* 2017, 35:35-37.
46. for the adept<sup>TM</sup> 2 investigators: Lentz S, Ehrenforth S, Karim FA, Matsushita T, Welding KN, Windyga J, Mahlangu JN: **Recombinant factor VIIa analog in the management of hemophilia with inhibitors: results from a multicenter, randomized, controlled trial of vatreptacog alfa.** *J Thromb Haemost* 2014, 12:1244-1253.
47. Hwang JK, Wang C, Du Z, Meyers RM, Kepler TB, Neuberg D, Kwong PD, Mascola JR, Joyce MG, Bonsignori M et al.: **Sequence intrinsic somatic mutation mechanisms contribute to affinity maturation of VRC01-class HIV-1 broadly neutralizing antibodies.** *Proc Natl Acad Sci U S A* 2017, 114:8614-8619.

48. DeKosky BJ, Kojima T, Rodin A, Charab W, Ippolito GC, Ellington AD, Georgiou G: **In-depth determination and analysis of the human paired heavy-and light-chain antibody repertoire.** *Nat Med* 2015, 21:86-91.
49. DeKosky BJ, Lungu OI, Park D, Johnson EL, Charab W, Chrysostomou C, Kuroda D, Ellington AD, Ippolito GC, Gray JJ, Georgiou G: **Large-scale sequence and structural comparisons of human naive and antigen-experienced antibody repertoires.** *Proc Natl Acad Sci U S A* 2016, 113:E2636- E2645.
50. Adler AS, Bedinger D, Adams MS, Asensio MA, Edgar RC, Leong R, Leong J, Mizrahi RA, Spindler MJ, Bandi SR et al.: **A natively paired antibody library yields drug leads with higher sensitivity and specificity than a randomly paired antibody library.** *MAbs* 2018, 10:431-443.
51. Sheng Z, Schramm CA, Kong R, NISC Comparative Sequencing Program, Mullikin JC, Mascola JR, Kwong PD, Shapiro L: **Gene-specific substitution profiles describe the types and frequencies of amino acid changes during antibody somatic hypermutation.** *Front Immunol* 2017, 8:537.
52. Kovaltsuk A, Leem J, Kelm S, Snowden J, Deane CM, Krawczyk K: **Observed antibody space: a resource for data mining next generation sequencing antibody repertoires.** *J Immunol* 2018, 201:2502-2509.

## CHAPTER 3

### **Impact of in vivo protein folding probability on local fitness landscapes**

This chapter is adapted with permission from the publication “Impact of in vivo protein folding probability on local fitness landscapes” in *Molecular Biology and Evolution* by Matthew S. Faber, Emily E. Wrenbeck, Laura R. Azouz, Paul J. Steiner, and Timothy A. Whitehead.

Copyright 2019 Oxford University Press.

## Abstract

It is incompletely understood how biophysical properties like protein stability impact molecular evolution and epistasis. Epistasis is defined as specific when a mutation exclusively influences the phenotypic effect of another mutation, often at physically interacting residues. In contrast, nonspecific epistasis results when a mutation is influenced by a large number of non-local mutations. As most mutations are pleiotropic, the *in vivo* folding probability - governed by basal protein stability - is thought to determine activity-enhancing mutational tolerance, implying that nonspecific epistasis is dominant. However, evidence exists for both specific and nonspecific epistasis as the prevalent factor, with limited comprehensive datasets to support either claim. Here we use deep mutational scanning to probe how *in vivo* enzyme folding probability impacts local fitness landscapes. We computationally designed two different variants of the amidase AmiE with statistically indistinguishable catalytic efficiencies but lower probabilities of folding *in vivo* compared to wild-type. Local fitness landscapes show slight alterations among variants, with essentially the same global distribution of fitness effects. However, specific epistasis was predominant for the subset of mutations exhibiting positive sign epistasis. These mutations mapped to spatially distinct locations on AmiE near the initial mutation or proximal to the active site. Intriguingly, the majority of specific epistatic mutations were codon-dependent, with different synonymous codons resulting in fitness sign reversals. Together, these results offer a nuanced view of how protein folding probability impacts local fitness landscapes, and suggest that transcriptional-translational effects are as important as stability in determining evolutionary outcomes.

## Introduction

Understanding the mechanisms of molecular evolution is important to molecular biology, virology, evolutionary biology, and protein engineering. Researchers interested in evolving natural proteins, designing proteins *de novo*, or understanding the extent of contingency on extant proteins must contend with the implicit evolutionary limitations set forth by nature. The challenge, then, is to understand what constrains protein evolution and by what mechanisms. How do these factors interact with one another to alter the frequency of mutations with increased fitness in a given environment, and how do they govern evolvability for new functions?

A particularly important component of evolution is epistasis, or the non-additive combination of mutations<sup>1</sup>. Epistasis impacts the rate of evolution and the spectrum of possible evolutionary pathways available to a protein<sup>2</sup>. Epistasis is said to be specific when a mutation exclusively influences the phenotypic effect of only a few select mutations, usually at physically interacting residues<sup>3</sup>. In contrast, epistasis is said to be nonspecific when a mutation impacts a global property like stability that can be rescued by large numbers of non-local mutations. Of the two classes, specific epistatic effects exert the greatest influence on the possible evolutionary outcomes<sup>3</sup>. This is the result of the precise and long-lasting amino acid constraints imposed by specific epistatic mutations, which decrease the evolutionary reversibility of a protein sequence in a given evolutionary trajectory. In turn, the delocalized and unconstrained effects resulting from non-specific epistasis often breakdown over evolutionary time, largely restricting its influence to short-term evolution. Non-specific epistasis can temporarily alter the mutational robustness of a protein by changing the permissivity to additional mutations at many more positions than does specific epistasis. By modulating the evolutionary trajectories available, epistatic phenomena exert immense influence on the short and long-term evolution of proteins<sup>4</sup>.

What remains incompletely understood is how biophysical parameters like protein stability constrain epistasis. Protein stability as defined here is the cumulative balance of the thermodynamic stability, the folding rate, and the fidelity of the association process for oligomeric proteins; these parameters combine to determine the likelihood that an enzyme will assume its native state when expressed *in vivo*: the *in vivo* folding probability of a given protein<sup>5,6</sup>. For enzymes, fitness is often a function of flux through a pathway, which is a product of steady state enzyme concentration and specific velocity. The steady state enzyme concentration, in turn, is a function of the translation rate, the probability of the nascent peptide folding into the native state, and the protein degradation rate. Typical evolutionary models control for translation rate and protein degradation rate and then assume that (i.) thermodynamics of protein folding can be described by a 2-state model; and (ii.) this single Gibbs free energy term can account for the probability of folding. However, these assumptions fail for much of a typical proteome. Many proteins are oligomeric, multi-domain proteins have more complicated folding trajectories, and there is increasing evidence that formation of secondary structure and partial hydrophobic collapse before ribosomal release is important for on-target folding<sup>7</sup>. The *in vivo* folding probability encompasses all of these biophysical terms into the probability of reaching the folding state.

Analyses of the impacts of stability in evolution at the genomic<sup>8</sup>, protein<sup>9-11</sup>, and organismal<sup>12</sup> levels have uncovered a complex and dynamic equilibrium between stabilizing and destabilizing mutations. For enzymes in particular, previous studies have shown that missense mutations often act pleiotropically where catalytically enhancing mutations are, on average, moderately destabilizing<sup>9,13,14</sup>. Consequently, high basal stability can buffer catalytically beneficial but destabilizing mutations<sup>15,16</sup>, allowing fixation. Deleterious destabilizing mutations

can be repaired by reversion mutations<sup>17</sup>, or by specific and non-specific epistatic mutations that rescue stability<sup>8,18</sup>. These epistatic mutations are a central phenomenon in the stabilizing-destabilizing equilibrium, with significant consequences in long-term evolution<sup>19,20</sup>. It is uncertain whether specific or non-specific epistatic mutations are more likely to rescue a destabilized protein, with evidence existing for both arguments<sup>17-21</sup>.

Deep mutational scanning experiments provide a wealth of mutational data that can be used to address questions in molecular evolution<sup>22</sup>. This technology comprises the use of large mutational libraries with selections coupled to deep sequencing to evaluate relative fitness of thousands of variants in a massively parallel fashion<sup>14,23-25</sup>. We previously used deep mutational scanning on the homohexameric aliphatic amidase AmiE from *Pseudomonas aeruginosa* to understand how local fitness landscapes, defined here as the set of all possible single-point amino acid substitutions from wild-type (WT), change with different substrates<sup>25</sup>. In this original study, AmiE was chosen as a model as it is stable in its genetic background and has a high probability of folding upon translation. To comprehensively assess how the initial probability of folding *in vivo* constrains mutational outcomes, we designed two variants of AmiE in which catalytic activity is unperturbed but the proteins have different *in vivo* folding probabilities. We then used deep mutational scanning to probe the local fitness landscapes of these variants. While we found moderate epistasis, local fitness landscapes are largely insensitive to the initial *in vivo* folding probability of the enzyme variant. In particular, the great majority of beneficial mutations were shared between all three starting points: WT AmiE, and the two disrupted single point-mutant enzymes. However, positive sign epistasis was present and was dominated by specific epistasis. Remarkably, we found that the sign of the fitness metric for many mutations depends on the codon used to encode the mutation, suggesting more complicated fitness landscapes than

predicted from intrinsic protein biophysics. Together, these results provide a nuanced view of how local fitness landscapes are perturbed under slightly different initial *in vivo* folding probabilities.

## Results

The experimental pipeline used in this study is shown in **Figure 3.1A**. First, we designed variants of AmiE that possess wild-type catalytic activity but with a reduced probability of folding *in vivo*. Second, we developed selection conditions for the variants under which cell growth is proportional to enzyme activity using a growth selection with acetamide as the sole nitrogen source. Third, near-comprehensive single-site saturation mutant libraries for our variants were prepared<sup>25</sup> and growth selections performed. Fourth, pre- and post-selection populations were deep sequenced to extract mutant frequencies in the selected and reference populations. These frequencies were converted into a relative fitness metric ( $\zeta_i$ ) for each mutant *i* defined as

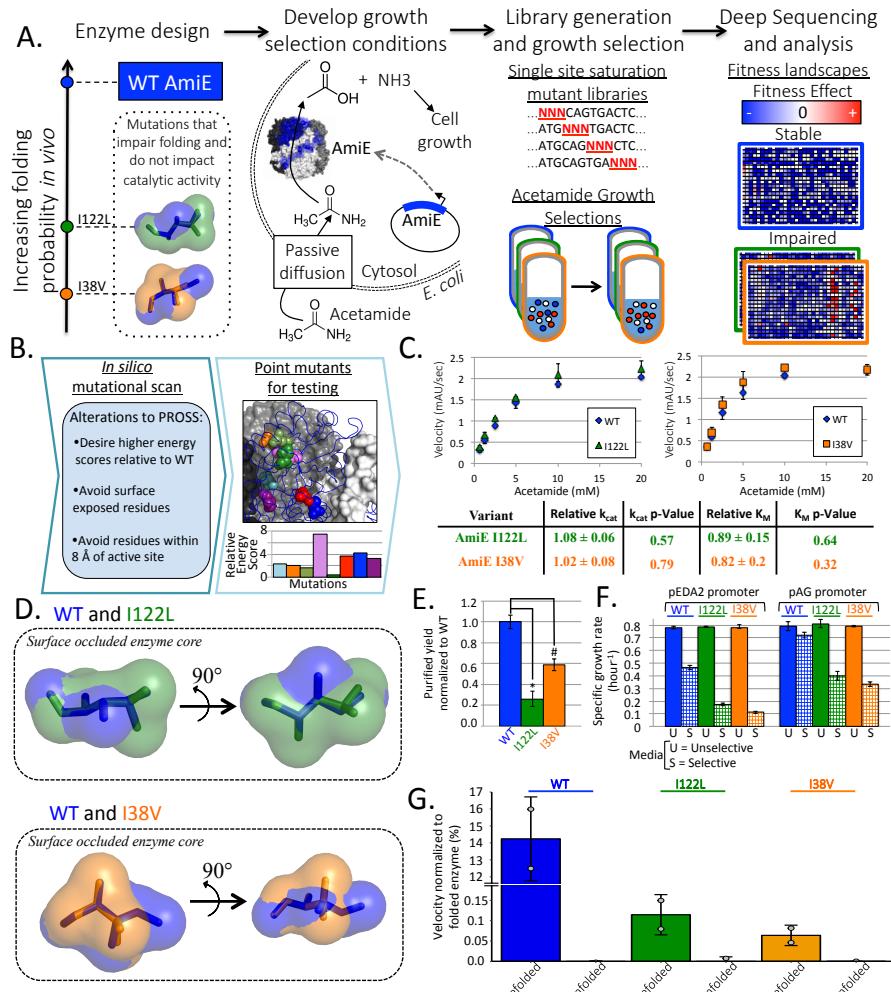
$$\zeta_i = \log_2\left(\frac{\mu_i}{\mu_{REF}}\right) \quad (1)$$

where  $\mu_i$  and  $\mu_{REF}$  represent the specific growth rates in selection media for the mutant ( $\mu_i$ ) and unmutated AmiE variant ( $\mu_{REF}$ ), respectively. A relative fitness score above zero means that a strain harboring a given mutant has higher fitness than those carrying the unmutated variant. It is important to understand the limitations of deep mutational scanning experiments. While these experiments provide quantitative measurements of fitness relative to the respective backgrounds, deep mutational scanning cannot provide information on why a given mutation is beneficial or deleterious. Thus, specifying whether a mutation impacts stability, catalytic efficiency, on-path folding rate, etc. is left to reasoned speculation or further biophysical analysis.

*AmiE variants with lower in vivo folding probabilities and wild-type catalytic efficiencies engineered*

We first sought to identify mutations to AmiE that, under the selection conditions, would decrease the in vivo folding probability of the protein while maintaining wild-type catalytic efficiency. To identify such mutants, we chose to use a computational approach by modifying PROSS<sup>26</sup>. Briefly, PROSS designs a protein sequence that will have an improved probability of reaching the folded state in vivo relative to its input. This improved folding probability correlates with biophysical properties like improved protein stability, faster on-target folding rate, or reduced aggregation propensity. As our experimental objective is essentially the inverse problem, we modified the Rosetta FilterScan protocol undergirding PROSS and then selected point-mutations with higher energy scores relative to AmiE wild-type (WT) (**Figure 3.1B, Table A 1**). For each mutant these scores were then cross-referenced with experimental relative fitness scores previously determined for AmiE<sup>25</sup> to ensure that their relative fitness was below zero (**Table A 1**).

Of thirteen variants with 1-3 mutations from WT selected for experimental characterization, nine expressed as soluble proteins in E. coli BL21\* (DE3). We purified a subset of these nine variants and assessed their catalytic efficiency with the substrate acetamide. While most mutants showed reduced enzymatic activity, both AmiE I38V and AmiE I122L showed statistically indistinguishable maximum turnover rates (kcat) and Michaelis constants (KM) compared with WT (**Figure 3.1C, Table A 2**). Furthermore, size exclusion chromatography showed no oligomeric differences between AmiE WT and AmiE I38V or AmiE I122L (**Figure**



**Figure 3.1: Design of deep mutational scanning experiment** A. A graphical overview of this study. Two AmiE enzyme variants with single point mutants (I38V and I122L) with WT catalytic function and lower probabilities of folding *in vivo* were computationally designed and validated experimentally. Constitutive expression of each enzyme from a plasmid was tuned such that the growth rate of our bacterial growth selection strain in selection media was dependent on the expression of functional AmiE. Deep mutational scanning was performed on these variants and compared with WT AmiE. B-G. Design and validation of AmiE variants. B. A graphical representation of the computational enzyme design. C. Enzyme velocity as a function of acetamide and Michaelis-Menten parameters determined relative to WT AmiE. Error bars = 1 s.d., n = 2, p-values obtained using Student's t-test. D. Structural modeling of the cavities introduced into AmiE by designed mutations. E. Enzyme yield following *E. coli* auto-induction expression. Error bars = 1 s.d., n = 3, \* = p-value = 0.0003, # = p-value = 0.002. F. Specific growth rates of strains in M9 (unselective) and in M9 with 10 mM acetamide as sole nitrogen source (selective) (pEDA2 - low expression, pAG - high expression). Error bars = 1 s.d., n ≥ 3. G. Comparison of enzyme reaction velocities at substrate saturation relative to a folded control. Grey dots represent biological replicates, Error bars = 1 s.d., n = 2.

**A 1)** in PBS at 30  $\mu$ M, suggesting that the variants maintain the expected homohexameric quaternary structure. Finally, the secondary structure of the three AmiE proteins was analyzed using far-UV circular dichroism spectroscopy (**Figure A 2**), revealing indistinguishable spectra in the folded and unfolded states.

AmiE I38V removes a methyl group to open a small cavity in the core, while I122L modulates hydrophobic core packing in the monomer subunit (**Figure 3.1D**, **Figure A 3**). Both mutations are located in the hydrophobic core distal from the dimeric and homohexameric contacts necessary for quaternary assembly (**Figure A 3**). Based on Rosetta analysis, we predict that these mutations disrupt the core of AmiE resulting in thermodynamic destabilization of the native monomer. However, mutations in the stability cores of proteins can disrupt the hierarchy of folding<sup>27</sup> through the destabilization of folding intermediates<sup>28</sup>, by limiting the intermediate states accessible during folding<sup>29</sup>, and by decreasing the thermodynamic stability<sup>30</sup> of the monomeric subunits outside of the quaternary structure<sup>31</sup>.

To distinguish among these possibilities, we attempted tryptophan fluorescence unfolding measurements using guanidinium-HCl (Gdn-HCl) as a denaturant to determine the effective thermodynamic stability. However, AmiE WT aggregated in moderate Gdn-HCl concentrations under most conditions (data not shown), and under conditions of no aggregation and complete unfolding no isosbestic point was recovered (**Figure A 4**). This lack of an isosbestic point indicates more complicated reversible folding at 4°C than simple 2-state models. We also performed thermal shift assays with the purified homohexameric enzymes in a series of dilutions (10, 5, 1, 0.5, 0.25  $\mu$ M) (**Figure A 5**) to measure thermal stabilities. Two-state irreversible unfolding curves were obtained at 10 and 5  $\mu$ M. Analysis of the melting curves reveals statistically indistinguishable melting temperatures between WT and variants (**Table A 3**). This

data suggests that all have a single transition from homohexamer into unfolded monomers. It is well established that oligomeric proteins are often more stable than in their natively folded monomeric or dimeric forms<sup>31</sup>. Thus it is likely that the thermal melt is measuring the stability of homohexameric assembly, which would be expected to be identical between WT and variants as neither mutation resides at an oligomeric interface. To identify AmiE concentrations for which the monomeric form is favored, we reasoned that hexameric dissociation would result in inactive enzyme<sup>32</sup>, which could be measured colorimetrically using our established activity assay. Indeed, activity analysis of dilute solutions of AmiE WT at 500 nM showed larger decreases in activity than 900 nM over moderate incubation periods (**Figure A 6**). Unfortunately, usable signal was not detected at protein concentrations of less than 5 μM in the thermal shift assays (**Figure A 5**). To assess thermal denaturation at lower enzyme concentrations, we performed circular dichroism thermal melts using a protein concentration of 1 μM (**Figure A 7**). Under these conditions both AmiE I38V and AmiE I122L have modest but significantly lower melting temperatures than AmiE WT (p-value 0.01 for I38V and 0.03 for I122L; **Table A 4**).

While we were only able to establish moderate decreases in the thermal stabilities of the AmiE variants, complementary *in vitro* and *in vivo* experiments strongly support that both I122L and I38V variants have lower *in vivo* folding probabilities than WT in the general order: I38V<I122L<WT. All synonymous codons encoding the I38V and I122L mutations had fitness metric below zero (**Table A 1**), suggesting the loss of fitness is a result of changes at the protein level rather than effects resulting from the codon used<sup>25</sup>. When driven from the same T7 promoter under identical Studier auto-induction<sup>33</sup> protein expression conditions, both AmiE I38V and I122L have statistically significant lower purification yields of soluble protein than WT (**Figure 3.1E, Table A 2**). Furthermore, *E. coli* harboring the AmiE variants expressed from

the same plasmid – pEDA2<sup>25</sup> maintaining the same constitutive promoter, ribosome-binding site (RBS), and 5' untranslated region (5' UTR) – showed lower specific growth rates than WT when grown with acetamide as the sole nitrogen source (**Figure 3.1F, Table A 2**). These results suggest that the *in vivo* folding probability upon translation for these variants is lower than for WT. Finally, while denatured WT can refold into active enzyme at 14.2% yield, both I38V (0.06%) and I122L (0.11%) have vastly lower refolding yields (**Figure 3.1G, Table A 2**). These results together support a model where the I38V and I122L mutations result in a lower probability of correctly folding into active homohexameric enzyme *in vivo*.

#### *Deep mutational scans for AmiE variants*

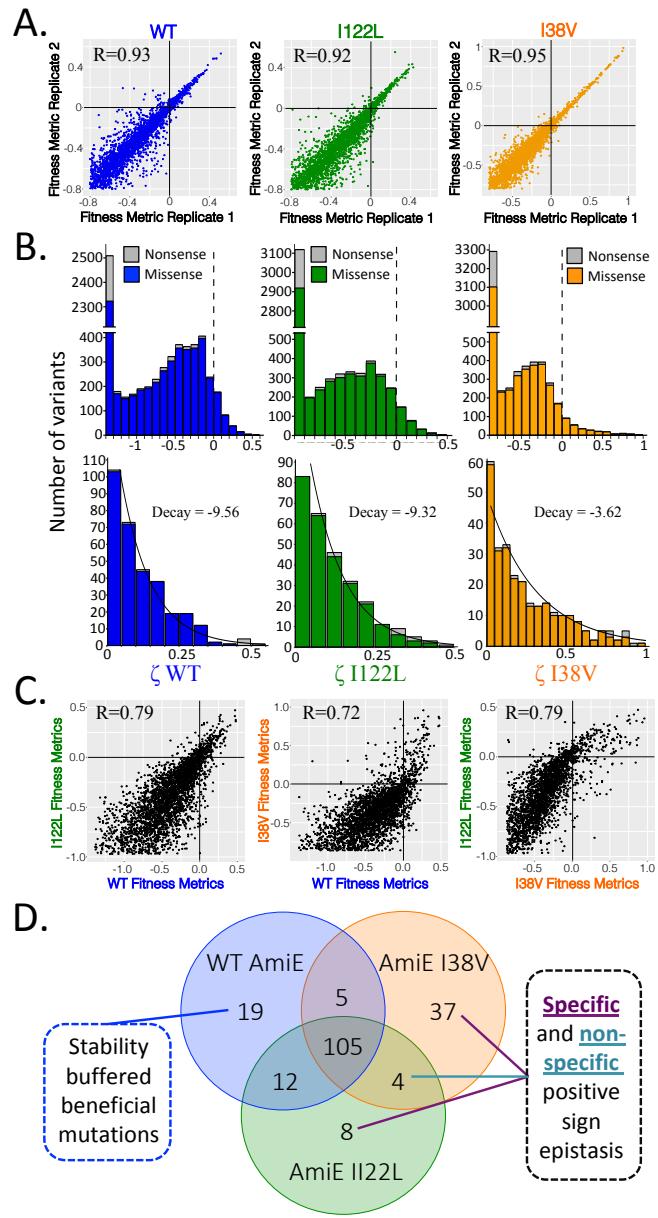
Deep mutational scanning of these variants was performed using a previously developed growth selection<sup>25</sup> in media with 10 mM acetamide as the sole nitrogen source. These growth selections required tuning the constitutive amidase expression such that the specific growth rate of variant i expressed in *E. coli* MG1655 rph+ in the selection media relative to that in defined minimal media ( $\mu_s, i / \mu M9, i$ ) is 0.4-0.6. However, plasmid pEDA2 used for AmiE WT selections did not support high enough growth rates for the I38V and I122L variants (**Figure 3.1F, Table A 2**). Thus, we screened additional promoters for AmiE I38V and AmiE I122L while maintaining the same 5' UTR and RBS for all constructs in order to minimize potential variant-dependent mRNA effects on fitness (**Table A 5**). Plasmid pAG with a stronger constitutive promoter than pEDA2 supported a growth rate ratio of  $0.49 \pm 0.03$  for AmiE I122L and  $0.42 \pm 0.02$  for AmiE I38V (**Figure 3.1F, Table A 2**). By contrast, pAG AmiE WT had a nearly 2-fold higher growth rate ratio of  $0.91 \pm 0.04$  (**Figure 3.1F, Table A 2**).

Next, we generated near comprehensive single-site saturation mutant libraries using nicking mutagenesis<sup>34</sup> (**full library statistics are shown in Table A 6 and Table A 7**). For AmiE I38V mutations at residues 32-44 flanking the site of the disrupting mutation were not constructed, while for AmiE I122L mutations at residues 115-130 and 132 were not made. Plasmids expressing mutant enzyme libraries were electroporated into E. coli MG1655 rph+ under conditions minimizing double transformants. Then, strains harboring AmiE libraries underwent growth selections in replicate with initial population sizes of  $>6 \times 10^6$  cells for approximately 8 generations at 37°C. A biological replicate for AmiE WT covering residues 171-255 was also performed to compare with previous published results<sup>25</sup>. The pre- and post-selection populations were barcoded and deep sequenced. The resulting data was processed using PACT<sup>35</sup> to obtain the relevant fitness metrics for each mutant in the library. The depth of sequencing ranged from 155 to 300-fold coverage for the libraries (**Figure A 8**). In total, we recovered the relative fitness metrics for 93.7% and 91.8% of all possible non-synonymous mutants for AmiE I122L and AmiE I38V, respectively (**Table A 7**).

To estimate reproducibility, we compared the AmiE WT replicate selections performed here with data from an identical selection experiment performed in Wrenbeck et al.<sup>25</sup>. Correlation coefficients between mutation-specific fitness values in selections are  $\geq 0.90$  (**Figure A 9**), which is comparable to correlation between replicates performed for this work (AmiE I122L - 0.921; AmiE I38V - 0.952) (**Figure 3.2A**). Additionally, there was essentially no correlation between relative fitness and pre-selection frequency of a given mutant in the library, (WT AmiE – R = 0.011, AmiE I122L – R = 0.026, AmiE I38V – R = -0.0376) (**Figure A 10**) indicating that pre-selection read counts do not bias the fitness metrics obtained.

*Distribution of beneficial fitness effects are largely insensitive to initial in vivo protein folding probability*

The shape of the distribution of fitness effects (DFE) governs the local protein fitness landscape. Realizing that beneficial mutations are rare, the likelihood of finding beneficial mutations was predicted by Orr<sup>36</sup> to follow the Pareto family of distributions. Using the set of beneficial mutations – variants with relative fitness above wild-type under selective media – we were previously able to describe the shape of the DFE for beneficial mutations as exponential with high statistical power<sup>25</sup>. The new datasets allow us to ask directly whether the shape of DFE changes with respect to enzyme *in vivo* folding probability. Consistent with expectations, all variants have very similar distributions of fitness effects (**Figure 3.2B** and **Figure A 11**) with a tight range of total possible mutations that are beneficial. For all variants the Pareto family of functions also describes their distributions of beneficial fitness effects (**Table A 8**). Thus, given approximately the same relative fitness, the probability of finding rare beneficial mutations is independent of initial likelihood of native folding.



**Figure 3.2: Local fitness landscapes are nearly insensitive to initial protein folding probability *in vivo*.** A. Correlation between AmiE variant technical replicates. B. Distributions of fitness effects (DFE) of nonsense and missense mutations for the AmiE variants. Upper plots show full DFE, while lower plots include only beneficial mutants with best-fit exponential curve. C. Correlation of fitness between AmiE variants following the combination of replicate datasets and reprocessing using PACT. D. Venn diagram for all unique and shared beneficial mutations for the respective variants.

### *Moderate epistasis observed with decreasing enzyme *in vivo* folding probability*

How does the local fitness landscape change in response to a single deleterious point-mutation that alters only the *in vivo* folding probability of an enzyme? If mutations were completely additive with the disrupting mutations, we would expect the comparison of the local fitness landscapes for the enzymes to have 1:1 correlations and approach the correlation coefficients found between replicates ( $R \sim 0.92$ ). On the other hand, complete non-additivity of mutations would lead to minimal correlation. We were able to compare 2,813 mutations above the lower bound of relative fitness (45.4% of possible mutations) shared between the three datasets. Pearson's correlation analysis of the DFE finds that the WT local fitness landscape is reasonably correlated with that of the variants (WT vs. I38V  $R = 0.72$ , WT vs. I122L  $R = 0.79$ ), and this correlation is similar to that between I122L vs. I38V ( $R = 0.79$ ) (**Figure 3.2C**). Notably, these correlation coefficients are lower than for replicates. Furthermore, linear regression best fits show lower slopes between variants than within replicates (**Figure 3.2A** and **3.2C**).

Next, we tested for increased negative epistasis in our distribution of fitness effects. We predicted that the greater *in vivo* folding probability of AmiE WT provides a buffering effect, which could temper many deleterious mutations. To test this hypothesis, we generated empirical cumulative distribution functions (ECDF) for the deleterious mutations for each of the three enzymes (**Figure A 12 and Table A 9**). This analysis is limited to the range of deleterious mutations quantitatively captured in our experimental system<sup>23</sup>. Within this range, the application of the Kolmogorov-Smirnoff test failed to reject the null hypothesis that the ECDF of AmiE WT is below that of AmiE I38V (**Table A 9**). By the same test, we were unable to discriminate the ECDFs of AmiE WT and AmiE I122L (**Table A 9**). Therefore, AmiE I38V has a greater number of more deleterious mutations than AmiE WT. This indicates that the lower *in*

*vivo* folding probability of the AmiE I38V background increases the likelihood of negative epistatic effects when compared to the more resilient AmiE WT.

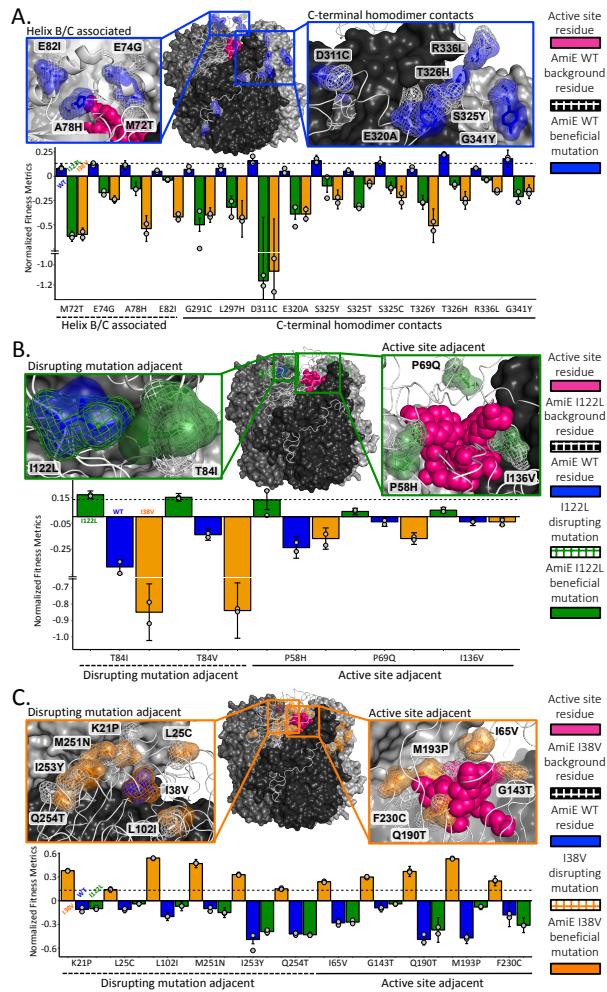
Precise measurements of negative and positive epistasis are complicated by the relatively narrow range of fitness our experimental system captures. However, we can determine the sign of fitness in our datasets with high precision. To evaluate the relative prevalence of sign epistasis, we defined any mutant as beneficial if  $\zeta_i > 0$  for both replicates and if  $\zeta_i > 0$  within a 95% confidence interval (see **Materials and Methods**). Conversely, we define a mutant as deleterious if  $\zeta_i < 0$  for both replicates and if  $\zeta_i < 0$  within a 95% confidence interval. We used these cutoffs to sort beneficial variants into the seven possible fitness bins (**Figure 3.2D**). Using a stricter requirement - that beneficial mutants are defined as those with a  $\geq 10\%$  increase in specific growth rate over the genetic background - leads to similar results (**Figure A 13**).

#### *Most beneficial mutations in the WT background are shared*

Previous studies found that stable proteins can buffer destabilizing mutations that are otherwise beneficial<sup>37,38</sup>. We are able to assess the extent of this phenomenon in our datasets. We find that 122/141 (86.5%) of beneficial mutations in the WT background are also beneficial in the I38V and/or the I122L genetic background (**Figure 3.2D**). Of these, six globally beneficial mutations (S9A, A28R, R89E, I165C, V201M, A234M) have been previously characterized biophysically<sup>25</sup> and are known to improve specific amidase flux under the selection conditions of 10 mM acetamide at 37°C. Conversely, only 19 of 141 beneficial mutations (13.5%) are specific in the WT background (**Figure 3.2D**). Therefore, beneficial mutations that are buffered in the stable background are present but in the minority.

The 19 WT-specific beneficial mutations map to two predominant locations: eleven (G291C, L297H, D311C, E320A, S325Y/T/C, 326Y/H, R336L, G341Y; **Figure 3.3A**) are at the extreme C-terminus that creates extensive homodimer contacts, while four (M72T, E74G, A78H, E82I; **Figure 3.3A**) are located on helix B and helix C distal to any oligomeric contacts in the homohexamer. The majority of the C-terminal mutations are adjacent to the homodimerization interface, which we speculate could lead to subtle structural rearrangements in the AmiE active site. Probable mechanisms behind the helix B/C mutations are more obscure as three of these mutations are at surface exposed positions over 10 Å away from any active site residue. We note that active site-induced effects may still be quite strong even this far away<sup>14,25,39-41</sup>. Regardless of the exact mechanisms, these findings indicate localized regions where small-scale mutational perturbations lead to increased fitness in a more stable genetic background.

To avoid information loss from the simple categorical separation of mutations into bins, normalized linear regression analysis was performed on correlation plots for the beneficial mutations shared in all backgrounds, in which the best-fit linear regression for the respective correlation plots was determined and normalized such that  $Y = 1X + 0$  (**Figure A 14**). In these correlation plots the mutational effects are highly disperse and widely deviate from the predicted fitness metric function for non-epistatic mutational combinations. The regression normalized AmiE I38V (**Figure A 14A**) and AmiE I122L (**Figure A 14B**) datasets compared with WT both show decreased correlations (AmiE I38V Pearson's  $R = 0.62$ ; AmiE I122L Pearson's  $R = 0.66$ ) compared to within replicates (Pearson's  $R > 0.92$ ). These distributions indicate that both classes of epistatic effects are present and imposing large impacts on this population of beneficial mutations.



**Figure 3.3: Positive sign epistatic mutations are spatially segregated and specific.**

A-C. Bar graphs of fitness metrics. Error bars represent 95% confidence intervals calculated from Poisson errors inherent in deep sequencing (Klesmith et al. 2015), while grey dots are fitness metric for each replicate. Horizontal dotted lines represent the cutoff value for mutations that increase the growth rate by  $\geq 10\%$ . Models show: trimer of dimers (wire + surface models), background residues for a respective enzyme (white mesh + sticks), the active site residues (magenta spheres), and where applicable the original mutation (green or orange mesh + sticks). In AmiE I122L and AmiE I38V the unique beneficial mutations tend to cluster around the disrupting mutations or near the active site. A. AmiE WT unique beneficial mutations are located at the C-terminal tail or in the B/C helices. B. Location of AmiE I122L unique beneficial mutations segregate to either positions adjacent to position 122 (T84I/V) or adjacent to the active site. C. Locations of a subset of AmiE I38V unique beneficial mutations. In B and C the blue transparent surface with sticks represents the residue in AmiE WT, and green or orange transparent surfaces with sticks represent respective unique beneficial mutations.

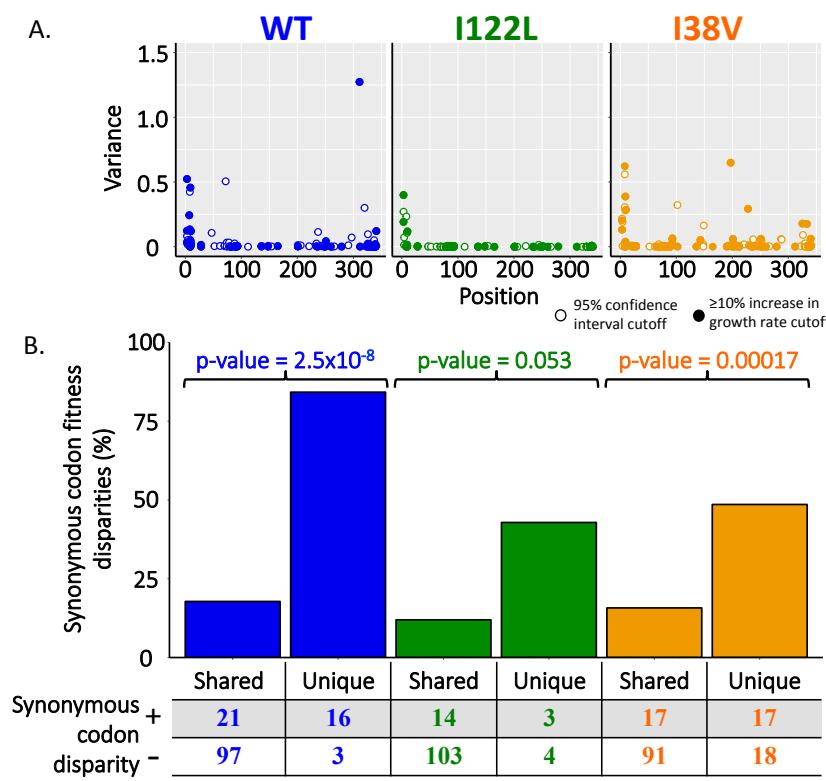
### *Positive sign epistasis is overwhelmingly specific*

Our datasets allow direct comparisons between the prevalence of specific and nonspecific epistasis in a folding impaired background. Contrary to previous literature on different model proteins<sup>2,15,38, 42-45</sup>, we find that specific reciprocal positive sign epistatic mutations dominate nonspecific mutations in the I38V and I122L backgrounds. In particular, we find 8 specific – unique beneficial – mutations for AmiE I122L and 37 specific mutations for AmiE I38V, compared with only 4 nonspecific – shared beneficial – mutations (**Figure 3.2D and 3.3B-C**). Analysis of the locations of the unique mutations in the structures of AmiE I122L and AmiE I38V suggest biophysical interpretations of their epistatic mechanism. For AmiE I122L the strongest specific mutations T84I/V directly contact 122L (**Figure 3.3B**), while similar mutations (K21P, L25C, L102I, M251N, I253Y, Q254T) occur on loops adjacent to I38V (**Figure 3.3C**). Other specific mutations line the enzyme active site. For AmiE I122L (**Figure 3.3B**) there are three (P58H, P69Q, I136V), while in AmiE I38V (**Figure 3.3C**) there are five (I65V, G143T, Q190T, M193P, F230C). Notably, since our datasets do not include immediately adjacent mutations for I38V and I122L, the extent of specific positive sign epistasis is probably underestimated.

### *A plurality of unique beneficial mutations is codon-dependent*

Fitness conferred by a weak-link enzymes depends on intrinsic protein biophysics but also on mRNA sequence-dependent effects. Perhaps the best appreciated of these are synonymous mutations in the first ten codons of a polypeptide because they can substantially alter mRNA stability and access to the ribosome binding sites in bacteria<sup>25,46</sup>. Additionally, synonymous codons can differentially affect cotranslational folding<sup>47</sup> and impact fitness. Here

we mapped the variance between synonymous codons encoding beneficial missense mutations (**Figure 3.4A**). For all three variants, the majority of high variance codons occur in the first 10 codons, as expected (**Figure A 15A**). However, there were localized punctae of high variance at several downstream positions for both WT and I38V datasets corroborated in replicate measurements. While overall variance in beneficial synonymous codons is weakly or not statistically significant among datasets (**Figure A 15B**), contingency table analysis of synonymous codon fitness disparities for the unique and shared beneficial mutations finds correlation with unique beneficial mutations in the WT and I38V backgrounds (WT p-value =  $2.5 \times 10^{-8}$ , I122L p-value = 0.053, I38V p-value = 0.00017; 2-tailed Fisher exact probability test)



**Figure 3.4: Unique beneficial mutations have high percentages of synonymous codon fitness disparities.** A. Variance of fitness metrics for synonymous codons of beneficial mutations as a function of position in the primary sequence. B. Percentage of shared and unique beneficial mutations with synonymous codon fitness metric disparities. p-values reported are from contingency table analysis with 2-tailed Fisher exact probability test.

**(Figure 3.4B).** In fact, the vast majority of WT (84.2%) - and many of the I122L (42.8%) and I38V (48.6%) - unique beneficial mutations have synonymous codon fitness sign disparities **(Figure 3.4B).** This indicates that transcriptional-translational effects impose significant evolutionary constraints.

## Discussion

In this study we used deep mutational scanning to analyze how the probability of attaining the folded, active state impacts local fitness landscapes. AmiE I38V and AmiE I122L were designed and validated to have identical catalytic parameters to WT but have a lower probability of folding under selection conditions. We found that the DFE for both variants was largely similar to the AmiE WT, and that most fitness-enhancing mutations are shared. However, there were two major surprises found when analyzing the set of mutations exhibiting positive sign epistasis.

First, we expected there to be a larger subset of beneficial mutations shared only between the I122L and I38V datasets, as current models of stability-induced epistasis posit that many nonspecific globally-distributed mutations can improve the probability of folding<sup>3</sup>. In contrast, we found that sign epistatic mutations were overwhelmingly specific for the I122L and I38V backgrounds. It is possible that beneficial non-specific sign epistatic mutations are the minority because of the alterations to in vivo expression that are required for the deep mutational scans. In the non-promoter tuned E. coli the impaired enzymes provide very weak cell growth. It is possible that non-specific epistatic effects may have arisen to be equivalent with, or dominant to, specific epistatic phenomenon if experiments could be performed without increasing the in vivo expression of the impaired enzymes. A limitation of our deep mutational scanning experiments is

that cell growth must be proportional to enzyme function and must operate within a window of growth rate ratios<sup>23</sup>. This requires increased *in vivo* expression of the impaired enzymes. By increasing the *in vivo* expression it is possible that the selective advantages that the non-specific mutations might possess could be dampened. Alternatively, the protein folding pathway for homohexameric AmiE in *E. coli* is potentially much more complicated than model systems of monomeric, single domain proteins that have built much of the current intuition about stability and epistasis. Therefore, the sparsity of non-specific beneficial mutations could be an artifact of our experimental system. These considerations also may explain the lack of a neutral peak for fitness expected from previous theoretical<sup>48</sup> and experimental<sup>49</sup> datasets on other proteins. More careful measurements on a wider array of oligomeric proteins should resolve this seeming contradiction.

As a second surprise, we found that unique beneficial mutations strongly depend on codon choice, as approximately 50% of sign epistatic mutations in the I38V background show sign disparities. We speculate that this unexpected result arises from the complicated co-translational folding *in vivo* of the homohexameric AmiE. Local, specific nonsynonymous mutations may recover on-target folding trajectories more efficiently than nonspecific, globally stabilizing mutations. Similarly, on-pathway folding kinetics may differ considerably between variants, which can be selectively modulated by codon choice<sup>50,51</sup>. As an alternative explanation, Kudla and colleagues recently report that synonymous codons can exert fitness effects through RNA toxicity itself<sup>52</sup> through an unknown mechanism.

While we were able to determine that both I122L and I38V mutations decrease the probability of active AmiE expression at 37°C, we were unable to measure the relative stability of the monomeric proteins. For the thermodynamic studies, AmiE aggregated under most

conditions, and where suitable conditions were found the lack of an isosbestic point hampered analysis. The thermal melts showed a single transition, most likely due to the hexameric dissociation. Lowering the AmiE concentration supported monomer formation at 25°C but yielded too weak a signal for analysis of monomer thermal stability. CD melts at low AmiE concentrations did show a modest but significant decrease in Tm in the variants relative to WT. Nevertheless, our results show that simple biophysical models currently used to model protein evolution are incomplete and that biophysical models may need to use kinetic models to account for the folding probability *in vivo*.

## Materials and methods

### *Reagents*

All antibiotics were purchased from GoldBio and all purchased enzymes were from New England Biolabs. All other chemicals were purchased from Sigma-Aldrich. Primers and mutagenic oligos were purchased from Integrated DNA Technologies and were designed using either Benchling ([www.benchling.com](http://www.benchling.com)) or the Agilent QuikChange Primer Design Program ([www.agilent.com](http://www.agilent.com)).

### *Computational design of folding impaired mutants*

The FilterScan Rosetta script<sup>26</sup> was modified to predict mutations that would decrease thermodynamic stability without altering catalytic efficiency. The structural coordinates for AmiE<sup>53</sup> (PDB: 2UXY, 341 residues per monomer) were taken from the Protein Data Bank and prepped for use in Rosetta scripts through the ‘clean\_pdb\_keep\_ligand.py’ script released with Rosetta 3<sup>54</sup>. The crystal structure data was refined through the “refine.xml” Rosetta scripts

(unaltered) from Goldenzweig et al.<sup>26</sup>. To avoid impacting catalytic efficiency residues within 8 Å of the active site were excluded from the FilterScan protocol. Additionally, surface residues were predicted and excluded from computational testing to avoid disturbing the native homohexameric state. The FilterScan script was modified to remove the input from the position specific scoring matrix (PSSM) evolutionary conservation term. Mutants with scores that predicted destabilization of the enzyme and were also shown to decrease relative fitness in a previous study<sup>24</sup> were selected for biophysical analysis.

#### *Plasmid construction*

The variants selected for biophysical analysis were constructed by mutating the AmiE WT sequence using the single mutation protocol of Nicking Mutagenesis<sup>34</sup>. The pEDA3 constitutive expression plasmid from Wrenbeck et al.<sup>34</sup> was used as the vector for the mutagenesis. Variants were subcloned from the pEDA3 background into protein expression plasmid pET-29b(+) (Novagen) or into the constitutive expression plasmid pEDA2<sup>25</sup> at *NdeI* and *XhoI* sites using classic restriction cloning. Plasmid pAG was constructed by mutating the -10 and -35 promoter regions of the pEDA3 plasmid using the multi-site nicking mutagenesis protocol from Wrenbeck et al.<sup>34</sup>. Following mutagenesis, the mutant promoter libraries were transformed into the *E. coli* growth selection strain MG1655 rph+ [F- λ-] (Coli Genetic Stock Center #7925, CGSC strain designation: BW30270). All AmiE variant DNA and protein sequences are listed in **Supplementary Notes A 1 and 2**.

### *Near comprehensive single-site mutant library construction*

Near comprehensive mutant libraries for AmiE I38V and AmiE I122L were constructed using the comprehensive nicking mutagenesis protocol from Wrenbeck et al.<sup>34</sup>. The genes for both impaired enzymes were broken into the following tiles: tile 1 (residues 1-85), tile 2 (residues 86-170), tile 3 (residues 171-255), and tile 4 (residues 256-341) in pAG. For AmiE I38V 13 residues were excluded from mutagenesis (residues 32-44), while for AmiE I122L 17 residues were excluded from mutagenesis (residues 115-130 and 132). Nicking mutagenesis products were amplified and purified as in Klesmith et al.<sup>24</sup>. 10 ng each of the respective DNA libraries was transformed into *E. coli* MG1655 *rph*<sup>+</sup> by electroporation performed with either a 1 mm electroporation cuvette at 1200 V (AmiE I122L; AmiE WT), or a 2 mm electroporation cuvette at 1600 V (AmiE I38V) using an Eppendorf Eporator. All experimental and control libraries were transformed into the selection strain with greater numbers than that required for theoretical complete library coverage (**Table A 6**). The transformation procedure for the selection strain was optimized to minimize double plasmid transformants as described in Kowalsky et al.<sup>23</sup>. -80°C freezer cell stocks of the libraries were prepared as detailed in Klesmith et al.<sup>24</sup>.

### *Protein expression and purification*

pET29(b) constructs harboring genes encoding AmiE variants were transformed into *E. coli* BL21\*(DE3) cells (Invitrogen) and expressed using Studier auto-induction<sup>33</sup> at 22°C for 16-18 hours. Cultures were pelleted and frozen at -80°C. Proteins were purified from cell pellets by Ni-NTA affinity chromatography exactly as described in Klesmith et al.<sup>24</sup>. Purified enzymes were desalting into phosphate buffered saline (PBS; 10 mM Na<sub>2</sub>HPO<sub>4</sub>, 1.8 mM KH<sub>2</sub>PO<sub>4</sub>, 2.7 mM

KCl, 137 mM NaCl, pH 7.4) using disposable PD-10 desalting columns (GE Healthcare), sterilized through a 0.22  $\mu$ m syringe filter, and stored at 4°C until analysis. Purified enzymes showed as a single band by SDS-PAGE (**Figure A 16**). Purified enzyme solutions were quantified using the absorbance at 280 nm in 1x PBS using a published<sup>55</sup> theoretical  $A_{280}$  molar extinction coefficient of 56,980 M<sup>-1</sup> cm<sup>-1</sup>.

For quantitative comparison of purified product yields under the T7 promoter system in BL21\* *E. coli*, slight alterations to the above induction scheme were used. For this analysis induction cultures were always started at an OD<sub>600</sub> of 0.005 from 1 mL LB + 50  $\mu$ g/mL kanamycin cultures grown at 37°C overnight. Next, 500 mL induction cultures were inoculated at an OD<sub>600</sub> of 0.005 with the overnight cultures and grown at 37°C with shaking at 250x rpm for 6 hours, and following this growth step cultures were moved to 22°C and induced for ~17 hours. Induction cultures were then pelleted and the wet cell weights of the pellets recorded. Cultures were then purified, desaltsed, and quantified as described above.

#### *Biophysical analysis of proteins*

Analysis of the growth rates of the AmiE variants in the respective constitutive expression plasmids in MG1655 rph+ *E. coli* were performed exactly as in Wrenbeck et al.<sup>25</sup>. Assessment of the oligomeric state of the purified enzymes was performed using SEC-FPLC. Approximately 3 mL of 30  $\mu$ M of the purified enzymes in PBS were run on an AKTA-FPLC system at 1 mL/min on an HiLoad 16/600 Superdex 200 column equilibrated with PBS. Enzyme kinetics (K<sub>M</sub> and k<sub>cat</sub>) was determined via phenol-alkaline hypochlorite end-point activity assays exactly as in Wrenbeck et al.<sup>25</sup>. Kinetic analysis of purified WT and folding impaired AmiE variants was performed within 6 days of purification.

To test for an isosbestic point the AmiE variants were denatured in guanidinium-HCl (GDN-HCl) and native tryptophan fluorescence was detected. Purified AmiE WT was diluted to 52  $\mu$ M in increasing amounts of ice cold GDN-HCl (0 – 4 M) in PBS with 1 mM DTT and mixed gently. 200  $\mu$ L of the denaturation mixtures were placed into opaque black 96 well plates and covered with optical film. Plates containing the samples were incubated at 4°C for 8 hours. Next, the native tryptophan fluorescence was measured in uncovered plates using an excitation of 290nm and emission of was detected over a range of wavelengths: 310 nm – 370 nm. This lack of an isosbestic point indicates an unfolding model that is more complex than a two-state model, presumably because of monomer folding and homohexamer association.

Thermal shift analysis was performed exactly as described in Wrenbeck et al.<sup>25</sup> with incubations for 2.5 hours at 25°C prior to addition of the dye and initiation of the thermal shift assays; the data was processed as in Huynh et al.<sup>56</sup>. In brief 10, 5, 2.5, 1, 0.5, and 0.25  $\mu$ M purified AmiE in 1x PBS was incubated at 25°C for 2.5 hours. Next, the samples had 5  $\mu$ L of 200x SYPRO-orange dye (Life Technologies) added to 45  $\mu$ L of the diluted enzymes in 0.1 mL MicroAmp<sup>®</sup> 96-Well Reaction Plate (Life Technologies) and covered with MicroAmp<sup>™</sup> optical film (Life Technologies). Thermal melt analysis was performed in a QuantStudio 6 Flex RT-PCR device (ThermoFisher). The melt ranged from 25°C to 98°C with 1°C change per minute and with a 2 minute incubation at the first and last temperatures. Thermal melt analysis with circular dichroism, and the far-UV spectral analysis of the folded and unfolded enzymes, was performed on a Chirascan plus spectrophotometer (Applied Photophysics). Purified AmiE was buffer exchanged into 10 mM phosphate buffer pH 7.5 using PD-10 desalting columns (GE Healthcare) approximately 24 hours prior to analysis and stored at 4°C. 1  $\mu$ M samples were diluted in the 10 mM phosphate buffer and stored in a 25°C water bath for at least 12 hours prior

to analysis. Prior to initiating the thermal melt analysis circular dichroism spectra were obtained from 180 nm to 260 nm at either 15°C (10 µM samples) or at 25°C (1 µM samples) in a 0.5 mm cuvette with 0.5 seconds per-time-point. Thermal melt analysis of the 1 µM samples started at 25°C and ramped to 95°C at a rate of 1°C/min with 0.5°C steps and a tolerance of 0.2°C. Signal at a wavelength of 222 nm was measured every 24 seconds throughout the melt. After the samples had reached the maximum temperature, the far-UV spectra were measured as before but with the temperature set at 95°C for all samples. Following circular dichroism analysis the obtained spectra were adjusted for their respective buffer blanks and smoothed using the Savitsky-Golay filter. The Window Size of the Savitsky-Golay filter was set to 14 for smoothing all buffer blanks prior to adjusting the sample spectra, and to 6 when smoothing the adjusted sample spectra. Finally, the obtained data was converted to Mean Residue Ellipticities (millideg \* cm<sup>2</sup> \* dmol<sup>-1</sup>) prior to curve fitting. Boltzmann curve fitting was performed to determine the apparent T<sub>m</sub> for both thermal melt experiments using GraphPad Prism ([www.graphpad.com](http://www.graphpad.com)).

To assess relative activity of the refolded enzymes, enzymes were first denatured in ice cold 3 M GDN-HCl in PBS supplemented with 1 mM DTT for 16 hours at 4°C at a final concentration of 50 µM. The solution was then diluted 50-fold into PBS with 1 mM DTT and 0.1% (w/v) BSA at 4°C in 96-well PCR plates that had been blocked with 1% (w/v) BSA in PBS for 1 hour at 37°C, resulting in a total protein concentration of 1 µM. Refolding mixtures were then incubated at 4°C for 5 minutes, and warmed to 37°C over 4.5 minutes. Samples were then held at 37°C for 20 minutes and then cooled to 4°C and held there until assaying. Immediately prior to assaying, samples were diluted in fresh ice-cold PBS with 1 mM DTT to ensure linearity in the activity assays. Enzymes were assayed using the phenol-alkaline hypochlorite end-point assay with 20 mM acetamide as the substrate. As a control, enzymes went through the same steps

as above except without initial GDN-HCl denaturation. The percent activity of the refolded enzyme was determined relative to the control sample. Two biological replicates were performed for all reactions. Refolding experiments were performed within 7 days of purification of the enzymes from the pellets.

#### *Growth selections*

Growth selections were performed exactly as in Wrenbeck et al.<sup>25</sup>. Briefly, starter cultures were grown overnight in the non-selective media (M9 minimal media: 47.6 mM Na<sub>2</sub>HPO<sub>4</sub>, 22 mM KH<sub>2</sub>PO<sub>4</sub>, 8.54 mM NaCl, 18.68 mM NH<sub>4</sub>Cl, 50 µg/mL carbenicillin, pH 7.0) and the following day the cells were washed in ice-cold M9 salt solution without ammonium chloride. Next, 3 mL of non-selective or selective media (M9 minimal media without ammonium chloride supplemented with 10 mM acetamide) was inoculated with the washed cells at an initial OD<sub>600</sub> = 0.02 (~6x10<sup>6</sup> cells) in Hungate tubes. Cultures were grown at 37°C with shaking at 250x rpm for approximately 8 generations. Continuous exponential growth was ensured by harvesting cells after the first 4 generations and re-inoculating with 3 mL fresh media + antibiotic at OD<sub>600</sub> = 0.02 prior to growth for the final 4 generations. Following 8 generations of growth the cells were stored and plasmid DNA extracted as in Klesmith et al.<sup>24</sup>. Unique selections started from the same unselected overnight culture were performed as replicates. To evaluate reproducibility between growth selections performed here and previous work performed on AmiE WT<sup>25</sup>, a deep mutational scan of residues 171-255 in AmiE WT was performed in parallel to each growth selection performed in the present work.

### *Deep mutational scanning*

AmiE variant DNA collected from the pre- and post-selection libraries was prepared for 300 BP paired end Illumina MiSeq sequencing as in Kowalsky et al.<sup>23</sup>. Primers used in the PCR reactions in preparation for Illumina sequencing are listed in (**Table A 10**). Sequencing of the variants was performed at the University of Illinois Chicago sequencing core. AmiE WT deep mutational scanning unprocessed sequencing results from Wrenbeck et al.<sup>25</sup> were downloaded from the SRA. Respective technical, and biological, replicates were processed independently. Regression analysis between the three starting points were performed by combining replicates within each variant. All data was processed using PACT<sup>35</sup> with the following changes from the default options entered into the configuration file: fast\_filter\_translate: qaverage = 20, and qlimit = 0; enrichment: ref\_count\_threshold = 5, sel\_count\_threshold = 0, strict\_count\_threshold = True. Normalized fitness metrics ( $\zeta_i$ ) were calculated by PACT as outlined in Kowalsky et al.<sup>23</sup>. To summarize, PACT calculates an enrichment ratio ( $\varepsilon_i$ ) for mutations by assessing the pre- and post-selection counts of each mutant:

$$\varepsilon_i = \log_2\left(\frac{f_{fi}}{f_{oi}}\right) \quad (2)$$

Where  $f_{fi}$  is the frequency of mutant  $i$  in the post-selection population and  $f_{oi}$  is the frequency in the pre-selection population. The normalized fitness metric for each mutant  $i$  ( $\zeta_i$ ) was next calculated using the population-averaged number of doublings during selection ( $g_p$ ) and the enrichment ratios of the mutant ( $\varepsilon_i$ ) and the unmutated starting variant ( $\varepsilon_{ref}$ ):

$$\zeta_i = \log_2\left(\frac{\left(\frac{\varepsilon_i}{g_p}\right) + 1}{\left(\frac{\varepsilon_{ref}}{g_p}\right) + 1}\right) \quad (3)$$

Lower bound fitness metrics - the cut-off below which fitness metrics cannot be discriminated from one another - were calculated as in Wrenbeck et al.<sup>25</sup> by using the median

read count for the pre-selection library and other statistics produced by PACT (**Figure A 8**). Lower fitness metrics were calculated to be: -1.38 for AmiE WT, -0.97 for AmiE I122L, and -0.86 for AmiE I38V.

AmiE I122L had seven outlier mutations removed because the difference in fitness metrics between replicates was greater than the 99.977% confidence intervals determined from sequencing depth of coverage (**Figure A 17**). No other mutants were removed from any of the datasets.

To account for global differences in fitness effects, we also generated normalized scatter plots for globally beneficial mutations. For each pairing of protein variants with AmiE WT (WT/I38V, WT/I122L), we performed a linear regression with the fitness of the first variant as the independent variable and the fitness of the second variant as the dependent variable. We then inverted the linear transformation obtained and applied it to the second fitness. Specifically, if the fitness of the second variant Y was modeled as  $Y = mX + B$ , then the normalized fitness  $\hat{Y}_\text{norm}$  was computed as  $(Y-B)/m$ . Thus, after normalization, the least-squares regression associated with each pairwise plot has a slope of one and an intercept of zero.

#### *Data availability*

Raw sequencing reads have been deposited in the Sequencing Read Archive (SRA [SAMN11258744 – SAMN11258771](#)). The processed data sets are available in **Figures A 18 -20** and in **Tables A 11 – 20**. The AmiE I38V pAG and AmiE I122L pAG plasmids used in mutant library generation have been deposited in Addgene (Addgene ID: 129791 and 129792), while the AmiE WT base construct modified in Wrenbeck et al.<sup>25</sup> was previously deposited in Bienick et al.<sup>55</sup> (Addgene ID: 59837).

## **Acknowledgements**

Thanks to Dr. J. Klesmith for his PACT troubleshooting help, E. Maurer and J. Hosten for their help with assorted tasks, and members of the Whitehead lab for providing feedback on ideas and figures. This work was supported by NSF CBET Career Award #1254238 to T.A.W.

## **REFERENCES**

## REFERENCES

1. Phillips PC: **Epistasis – the essential role of gene interactions in the structure and evolution of genetic systems.** *Nat Rev Genetics* 2008, 9:855-867.
2. Bershtain S, Segal M, Bekerman R, Tokuriki N, Tawfik DS: **Robustness- epistasis link shapes the fitness landscape of a randomly drifting protein.** *Nature* 2006, 444:929–932.
3. Starr TN, Thornton JW: **Epistasis in protein evolution.** *Protein Science* 2016, 25:1204-1218.
4. Breen MS, Kemeny C, Vlasov PK, Notredame C, Kondrashov FA: **Epistasis as the primary factor in molecular evolution.** *Nature* 2012, 490:535-538.
5. Baker D, Agard DA: **Kinetics versus thermodynamics in protein folding.** *Biochemistry* 1994, 33(24):7505-7509.
6. Shakhnovich EI: **Theoretical studies of protein-folding thermodynamics and kinetics.** *Curr Opin in Struc Biol* 1997, 7:29-40.
7. Goldenzweig A, Fleishman SJ: **Principles of protein stability and their application in computational design.** *Annu Rev Biochem* 2018, 87:105-129.
8. Jordan DM, Frangakis SG, Golzio C, Cassa CA, Kurtzberg J, Task Force for Neonatal Genomics, Davis EE, Sunyaev SR, Katsanis N: **Identification of cis-suppression of human disease mutations by comparative genomics.** *Nature* 2015, 524:225-229.
9. Tokuriki N, Stricher F, Serrano L, Tawfik DS: **How Protein Stability and New Functions Trade Off.** *PLoS Comp Biol* 2008, 4:e1000002.
10. Campbell E, Kaltenback M, Correy GJ, Carr PD, Porebski BT, Livingstone EK, Afriat-Jurnou L, Buckle AM, Weik M, Hollfelder F, Tokuriki N, Jackson CJ: **The role of protein dynamics in the evolution of new enzyme function.** *Nat Chem Biol* 2016, 12:944-950.
11. Kumar A, Natarajan C, Moriyama H, Witt CC, Weber RE, Fago A, Storz JF: **Stability-Mediated Epistasis Restricts Accessible Mutational Pathways in the Functional Evolution of Avian Hemoglobin.** *Mol Bio Evo* 2017, 34:1240-1251.
12. Serohijos AWR, Shakhnovich EI: **Contribution of Selection for Protein Folding Stability in Shaping the Patterns of Polymorphisms in Coding Regions.** *Mol Biol Evol* 2014, 31:165–176.
13. Tokuriki N, Jackson CJ, Afriat-Jurnou L, Wyganowski KT, Tang R, Tawfik DS: **Diminishing returns and tradeoffs constrain the laboratory optimization of an enzyme.** *Nat Comm* 2012, 3:1257.
14. Klesmith JR, Bacik JP, Wrenbeck EE, Michalczyk R, Whitehead TA: **Trade-offs between**

**enzyme fitness and solubility illuminated by deep mutational scanning.** *Proc Natl Acad Sci U S A* 2017, 114:2265-2270.

15. Bloom JD, Labthavikul ST, Otey CR, Arnold FH: **Protein stability promotes evolvability.** *Proc Natl Acad Sci U S A* 2006, 103:5869–5874.
16. Yu H, Dalby PA: **Exploiting correlated molecular-dynamics networks to counteract enzyme activity-stability trade-off.** *Proc Natl Acad Sci U S A* 2018, 115:E12192- E12200.
17. Ashenberg O, Gong LI, & Bloom JD: **Mutational effects on stability are largely conserved during protein evolution.** *Proc Natl Acad Sci U S A* 2013, 110:21071-21076.
18. Yu H, Dalby PA: **Coupled molecular dynamics mediate long- and short-range epistasis between mutations that affect stability and aggregation kinetics.** *Proc Natl Acad Sci U S A* 2018, 115:E11043-11052.
19. Shah P, McCandlish DM, Plotkin JB: **Contingency and entrenchment in protein evolution under purifying selection.** *Proc Natl Acad Sci U S A* 2015, 112:E3226-E3235.
20. Dasmeh P, Serohijos AWR: **Estimating the contribution of folding stability to nonspecific epistasis in protein evolution.** *Proteins* 2018, 86:1242-1250.
21. Pollock DD, Thiltgen G, Goldstein RA: **Amino acid coevolution induces an evolutionary Stokes shift.** *Proc Natl Acad Sci U S A* 2012, 109:E1352-E1359.
22. Fowler DM, Fields S: **Deep mutational scanning: a new style of protein science.** *Nat Methods* 2014, 11:801–807.
23. Kowalsky CA, Klesmith JR, Stapleton JA, Kelly V, Reichkitzer N, Whitehead TA: **High-Resolution Sequence-Function Mapping of Full-Length Proteins.** *PLoS One* 2015, 10:e0118193.
24. Klesmith JR, Bacik JP, Michalczyk R, Whitehead TA: **Comprehensive Sequence-Flux Mapping of a Levoglucosan Utilization Pathway in *E. coli*.** *ACS Synth Biol* 2015, 4:1235–1243.
25. Wrenbeck EE, Azouz LR, Whitehead TA: **Single-mutation fitness landscapes for an enzyme on multiple substrates reveal specificity is globally encoded.** *Nat Comm* 2017, 8:15695.
26. Goldenzweig A, Goldsmith M, Hill SE, Gertman O, Laurino P, Ashani Y, Dym O, Unger T, Albeck S, Prilusky J, Lieberman RL, Aharoni A, Silman I, Sussman JL, Tawfik DS, Fleishman SJ: **Automated Structure-and Sequence-Based Design of Proteins for High Bacterial Expression and Stability.** *Mol Cell* 2016, 63:337-346.
27. Raschke TM, Marqusee S: **The kinetic folding intermediate of ribonuclease H resembles the acid molten globule and partially unfolded molecules detected under native conditions.** *Nature Struc Biol* 1997, 4:298-304.

28. Raschke TM, Kho J, Marqusee S: **Confirmation of the hierarchical folding of Rnase H: a protein engineering study.** *Nature Struc Biol* 1999, 6:825-831.
29. Karshikoff A, Nilsson L, Ladenstein R: **Rigidity versus flexibility: the dilemma of understanding protein thermal stability.** *FEBS* 2015, 282:3899-3917.
30. Robic S, Berger JM, Marqusee S: **Contributions of folding cores to the thermostabilities of two ribonucleases H.** *Protein Science* 2002, 11:381-389.
31. Scholl ZN, Yang W, Marszalek PE: **Direct observation of multimer stabilization in the mechanical unfolding pathway of a protein undergoing oligomerization.** *ACS nano* 2015, 9:1189-1197.
32. Cervoni L, Egistelli L, Mocan I, Giartosio A, Lascu I: **Quaternary structure of Dictyostelium discoideum nucleoside diphosphate kinase counteracts the tendency of monomers to form a molten globule.** *Biochemistry* 2003, 42:14599-14605.
33. Studier FW: **Protein production by auto-induction in high-density shaking cultures.** *Protein Exp Purif* 2005, 41:207–234.
34. Wrenbeck EE, Klesmith JR, Stapleton JA, Adeniran A, Tyo KEJ, Whitehead TA: **Plasmid-based one-pot saturation mutagenesis.** *Nat Methods* 2016, 13:928-930.
35. Klesmith JR, Hackel BJ: **Improved mutation function prediction via PACT: Protein Analysis and Classifier Toolkit.** *Bioinformatics* 2018, doi/10.1093/bioinformatics/bty1042.
36. Orr HA: **The population genetics of beneficial mutations.** *Phil Trans R Soc B* 2010, 365:1195-1201.
37. Tokuriki N, Tawfik DS: **Stability effects of mutations and protein evolvability.** *Curr Opin Struct Biol* 2009, 19:596–604.
38. Gong LI, Suchard MA, Bloom JD: **Stability-mediate epistasis constrains the evolution of an influenza protein.** *eLife* 2013, 2:e00631.
39. Abriata LA, Palzkill T, Dal Peraro M: **How structural and physicochemical determinants shape sequence constraints in a functional enzyme.** *PloS One* 2015, 10:e0118684.
40. Jack BR, Meyer AG, Echave J, Wilke CO: **Functional sites induce long-range evolutionary constraints in enzymes.** *PLoS Biol* 2016, 14:e1002452.
41. Mayorov A, Dal Peraro M, Abriata LA: **Active site-induced evolutionary constraints follow fold polarity principles in soluble globular enzymes.** *Mol Biol Evol* 2019, doi:10.1093/molbev/msz096.
42. Huang W, Palzkill T: **A natural polymorphism in beta-lactamase is a global suppressor.** *Proc Natl Acad Sci U S A* 1997, 94:8801-8806.

43. Sideracki V, Huang W, Palzkill T, Gilbert HF: **A secondary drug resistance mutation of TEM-1 beta-lactamase that suppresses misfolding and aggregation.** *Proc Natl Acad Sci USA* 2001, 98:283-288.
44. Bloom JD, Silberg JJ, Wilke CO, Drummond DA, Adami C, Arnold FH: **Thermodynamic prediction of protein neutrality.** *Proc Natl Acad Sci USA* 2005, 102:606-611.
45. Tokuriki N, Tawfik DS: **Chaperonin overexpression promotes genetic variation and enzyme evolution.** *Nature* 2009, 459:668-675.
46. Kristofich J, Morgenthaler AB, Kinney WR, Ebmeier, CC, Snyder DJ, Old WM, Cooper VS, Copley SD: **Synonymous mutations make dramatic contributions to fitness when growth is limited by a weak-link enzyme.** *PLoS genetics* 2018, 14(8):e1007615.
47. Buhr F, Jha S, Thommen M, Mittlestaet J, Kutz F, Schwalbe H, Rodnina MV, Komar AA: **Synonymous codons direct co-translational folding towards different protein conformations.** *Mol Cell* 2016, 61:341-351.
48. Wylie CS, Shakhnovich EI: **A biophysical protein folding model accounts for most mutational fitness effects in viruses.** *Proc Natl Acad Sci USA* 2011, 108:9916-9921.
49. Stiffler MA, Hekstra DR, Ranganathan R: **Evolvability as a function of purifying selection in TEM-1  $\beta$ -lactamase.** *Cell* 2015, 160:882-892.
50. Weatheritt RJ, Babu MM: **Evolution. The hidden codes that shape protein evolution.** *Science* 2013, 342:1325-1326.
51. Firnberg E, Labonte JW, Gray JJ, Ostermeier M: **A comprehensive, high-resolution map of a gene's fitness landscape.** *Mol Biol Evol* 2014, 31:1581-1592.
52. Mittal R, Brindel J, Stephen JB, Kudla, G: **Codon usage influences fitness through RNA toxicity.** *Proc Natl Acad Sci USA* 2018, 115(34):8639-8644.
53. Andrade J, Karmali A, Carrondo MA, Frazão C: **Structure of Amidase from Pseudomonas aeruginosa Showing a Trapped Acyl Transfer Reaction Intermediate State.** *J Biol Chem* 2007, 282:19598-19605.
54. Leaver-Fay A, Tyka M, Lewis SM, Lange OF, Thompson J, Jacak R, Kaufman KW, Renfrew PD, Smith CA, Sheffler W, et al.: **Rosetta3: An Object-Oriented Software Suite for Simulation and Design of Macromolecules.** *Methods in Enzymology* 2011, 487:545-574.
55. Bienick MS, Young KW, Klesmith JR, Detwiler EE, Tomek KJ, Whitehead TA: **The Interrelationship between Promoter Strength, Gene Expression, and Growth Rate.** *PLoS One* 2014, 9:e109105.
56. Huynh K, Partch CL: **Analysis of protein stability and ligand interactions by thermal shift assay.** *Curr Protoc Protein Science* 2015, 79:28.9.1-28.9.14.

## **CHAPTER 4**

### **Saturation mutagenesis genome engineering of infective ΦX174 bacteriophage via unamplified oligo pools and golden gate assembly**

This work is adapted with permission from a publication that is in preparation by Matthew S. Faber, James T. Van Leuven, Martina M. Ederer, Holly H. Wichman Craig R. Miller, and Timothy A. Whitehead.

## **Abstract**

Here we present the first steps in the preparation of a novel protocol for the construction of saturation single-site - and massive multi-site - mutant libraries of a bacteriophage. We segmented the ΦX174 genome into 15 non-toxic and non-replicative fragments compatible with golden gate assembly. We next used nicking mutagenesis with oligonucleotides prepared from unamplified oligo pools with individual segments as templates to prepare near-comprehensive single-site mutagenesis libraries of genes encoding the F capsid protein (421 amino acids scanned) and G spike protein (172 amino acids scanned). Libraries possessed greater than 99% of all 11,860 programmed mutations. Golden Gate cloning was then used to assemble the complete ΦX174 mutant genome libraries, and libraries transformed into *E. coli* C were infective. This protocol will expand reverse genetics experiments possible for studying viral evolution and, with some modifications, can be applied for engineering of therapeutically relevant bacteriophages with larger genomes.

## **Introduction**

Predicting the tempo and trajectory of evolutionary change in the complex environments encountered by viruses and bacteria remains a challenge<sup>1-3</sup>. Understanding such changes are important for fundamental evolutionary studies as well as biotechnology applications like phage therapy for multi-drug resistant bacteria<sup>4</sup>. Two obstacles preventing better predictions are the oversimplification of the environment in experimental evolution studies compared to wild conditions and the vast number of possible mutational combinations that can occur, even in short adaptive walks and small genomes. Advances in DNA sequencing and synthesis have opened new ways to study microbial evolution and overcome some of these obstacles<sup>5,6</sup>. Unfortunately, suitable methods do not exist for generating comprehensive bacteriophage mutant libraries. This technology gap hinders those seeking to engineer phage for biotechnological applications and for those seeking a deeper understanding of how viruses evolve.

Methods for genetically engineering phages have recently been reviewed<sup>7</sup> and commonly involve homologous recombination and recombineering. While such approaches can be improved by clever incorporation of CRISPR-Cas systems<sup>8-10</sup>, the overall modest efficiencies largely limits mutagenesis to the generation and single-site mutation of chimeric genomes<sup>11,12</sup>, gene deletions<sup>13</sup>, and limited multi-site mutagenesis<sup>14</sup>. By contrast, human viruses often have reverse genetics systems in place where comprehensive mutant libraries can be prepared for single genes<sup>15</sup> or regions within a gene<sup>16,17</sup> using replicative plasmids that encode whole viruses or viral components. Similarly, deep mutational scanning can be performed on plasmid-encoded phage proteins like the MS2 capsid protein<sup>18</sup>.

Advances in the technologies for generating mutant libraries, synthesizing DNA, and for assembly of large DNA fragments<sup>19</sup> potentially allow for the facile construction and assembly of

user-defined mutagenesis of long nucleic acids. Nicking mutagenesis (NM) can be used to construct comprehensive single-site or other user-defined mutant libraries<sup>20</sup> using plasmid dsDNA as a template. In NM an oligo encodes the desired mutations by mismatch with the parental template. Recently, on-chip ink-jet printed oligo pools have been integrated into NM<sup>21</sup>, so one can now construct heterogeneous libraries of oligos that contain tens to hundreds of thousands of high fidelity, unique sequences<sup>21,22</sup> with a low per base pair cost.

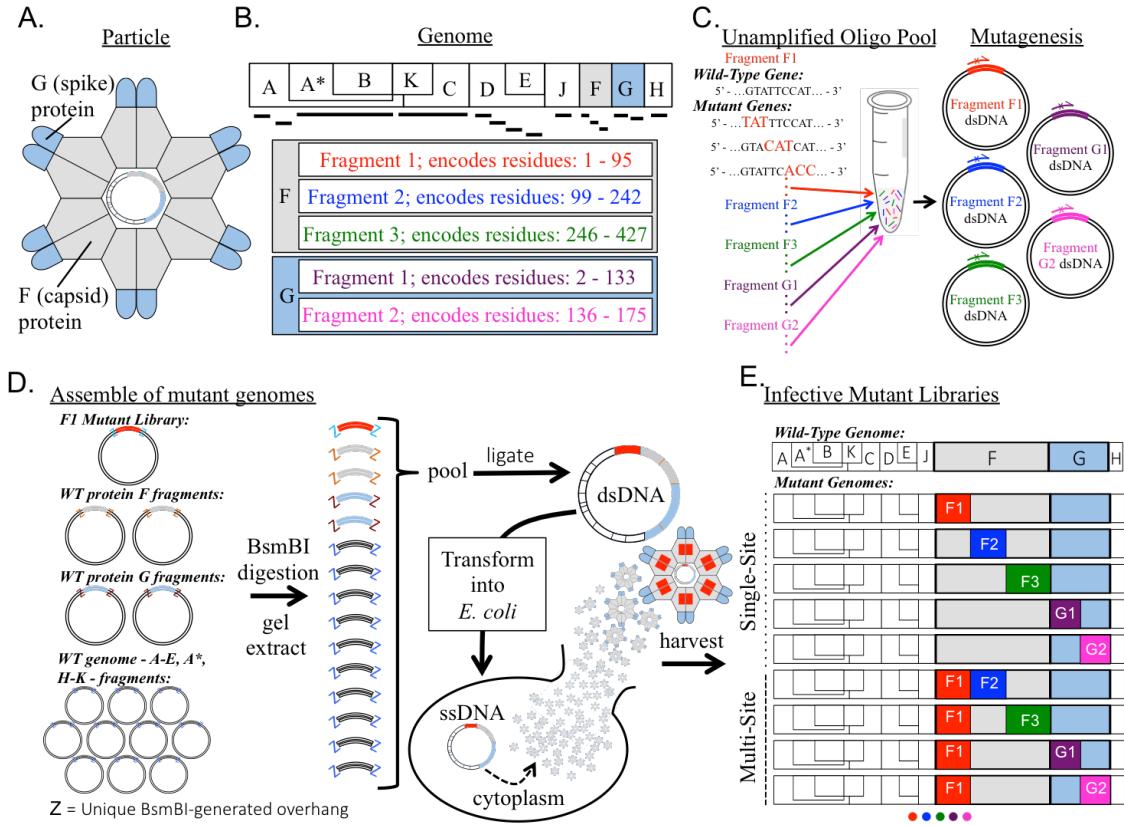
We wondered whether combining the advantages of replicative plasmids with advances in DNA synthesis and mutagenesis would allow us to create large user-defined libraries of bacteriophages. As a model bacteriophage system we selected  $\Phi$ X174 because it has been extensively studied<sup>23-26</sup>, has high resolution X-ray crystallography structures of the capsid and spike proteins<sup>27,28</sup>, and has a very small genome<sup>29</sup> of 5386 nucleotides. These attractive features have led to proof-of-principle demonstrations from other groups for synthetic genome assembly from oligonucleotides<sup>30</sup> and complete refactoring of the phage<sup>31</sup>. This method combines Golden Gate assembly, nicking mutagenesis, and oligo pool technology to construct near comprehensive single-site and expansive multisite mutant libraries for the genes encoding the entire capsid F protein and spike G protein of the bacteriophage  $\Phi$ X174. To our knowledge this is the first time near comprehensive single site mutant libraries of full capsid and spike proteins have been generated for an infective bacteriophage.

## Results

Our strategy for generating large user-defined mutagenesis libraries of bacteriophage  $\Phi$ X174 is shown in **Figure 4.1**. The circular 5386 nucleotide genome compactly encodes 11 genes where the F gene encodes a capsid and G encodes the spike protein (**Figure 4.1A**). The

genome was segmented onto 15 separate plasmids, including 3 plasmids for the F gene and 2 for the G gene (**Figure 4.1B**). User-defined comprehensive mutations on F and G were encoded using nicking mutagenesis with a single unamplified oligo pool containing all 11,860 mutagenic oligonucleotides (**Figure 4.1C**). ΦX174 mutant genomes are to be reconstituted from individual plasmids using Golden Gate cloning and transformed into a host *E. coli* strain (**Figure 4.1D**). Harvested phages will contain single or double non-synonymous mutations in the F and G genes depending on the mutagenized segments used for Golden Gate cloning (**Figure 4.1E**).

We first sought a reverse genetics system for ΦX174 wherein the virus chromosome was segmented and encoded on individual plasmids. Our method of construction closely followed that for the assembly of human coronavirus<sup>32</sup> NL63, where each plasmid contained unique BsmB1 type IIS restriction endonuclease sites flanking unique five nucleotide overlaps of wild-type (WT) ΦX174 sequence. This architecture allows for faithful assembly of the complete genome via Golden Gate cloning. However, as plasmids are replicated in *E. coli*, which ΦX174 naturally infects, the ΦX174 chromosome was encoded into 15 separate nontoxic plasmids where genes were separated



**Figure 4.1:  $\Phi$ X174 mutant library assembly.** A. Anatomy of the viral external surface. B. A linear schematic of the circular  $\Phi$ X174 ssDNA genome. Black lines delineate portions of the virus encoded on each of the 15 replicative plasmids. C. An unamplified oligo pool, containing all user-defined mutagenic oligonucleotides for F and G genes, was created and stored as a single mixture. This primer set was used in nicking mutagenesis for the generation of single-site saturation mutant libraries of genes encoding F and G proteins. D. Golden Gate cloning was used to assemble the mutant genomes. The dsDNA genome was transformed into *E. coli* and allowed a single burst phase. The resulting phage particles were collected and sequenced using SMRT sequencing. E. Linear schematics of the single site mutant libraries and some of the possible multi-site mutant libraries.

from their promoters and larger genes were segmented (full sequences for all plasmid inserts are given in Note B 4.1).  $\Phi$ X174 phage could be reconstituted by digesting all plasmids with BsmB1, ligating inserts overnight, and transforming into electrocompetent *E. coli* C cells. Following incubation, phage plaques were tabulated. Sequencing of the genome showed that the recombinant phage encoded the intended sequence.

The F and G genes were targeted for saturation mutagenesis. We chose to introduce mutations by nicking mutagenesis, which requires the presence of a unique BbvCI nicking site on the dsDNA plasmid. The F gene plasmid F3 encoding residues 246-427 of the F gene product contained a unique BbvCI sequence, while the remaining four plasmids (F1, F2, G1, G2) required introduction of BbvCI nicking sites in the vector backbone. The presence in each plasmid of the unique BbvCI nicking site was verified with the successful generation of circular ssDNA from the modified plasmids by BbvCI.Nt and exonuclease digestion (**Figure B 1**).

In nicking mutagenesis, desired mutations are encoded through libraries of mutagenic oligonucleotides, which can be sourced from unamplified oligo pools<sup>21</sup>. We custom synthesized a single oligo pool containing 11,840 oligos encoding nearly every non-synonymous single mutation in the F and G genes. Application of nicking mutagenesis using this same oligo pool for different plasmids (F1, F2, F3, G1, G2) resulted in at least 14-fold excess transformants required for 99.9% theoretical coverage of the desired library (**Table B 1**). The diversity of the mutant libraries was validated using deep sequencing on an Illumina MiSeq platform and all libraries were found to have >99% coverage of all possible single mutations (full library statistics are given in **Table 4.1**).

**Table 4.1: Mutant library NGS statistics.** Summary table of the of the libraries prior to viral genome assembly.

Fragment	F1		F2		F3		G1		G2
Tile Number	1	1	2	1	2	1	2	1	
Residues	1-95	99-157	158-242	246-336	337-427	2-53	55-133	136-175	
Sequencing reads post quality filter	1163959	681935	1714460	800377	757177	747800	969275	813221	
Fold oversampling of codon combinations	441.1	499.8	851	364.8	345.1	632.8	517.6	821.9	
<b>Percent of reads with:</b>									
No nonsynonymous mutations	23.0%	63.3%	52.3%	54.0%	55.0%	68.3%	48.8%	25.2%	
One nonsynonymous mutation	52.6%	27.5%	36.3%	33.1%	32.1%	24.1%	39.8%	59.7%	
Multiple nonsynonymous mutation	24.4%	9.2%	11.3%	12.9%	12.8%	7.6%	11.4%	15.1%	
Coverage of possible single nonsynonymous mutations	100%	100%	99.8%	100%	99.9%	99.6%	100%	100%	

ϕX174 genomes were assembled by Golden Gate cloning with 1 plasmid encoding a mutant library combined with the remaining 14 WT plasmids. All five ϕX174 libraries prepared

contained active phage. Golden Gate reaction products were transformed in *E. coli* C cells with the number of plaques resulting from each reaction tabulated (**Table 4.2**). 50-60% of all assembled genomes contained at least one mutation as estimated by Sanger sequencing of at least 30 plaques. In contrast, for the control where all plasmids in the Golden Gate reaction were WT, 0% of the phage contained mutations in the encoded regions. The percent theoretical coverage ranged from 9-94% (**Table 4.2**). The pooled G1 and G2 libraries were the most successful but only reach 94% the required number of transformants for near complete coverage. The least efficient, F2, provided only 9% of the needed number of transformants.

**Table 4.2: Mutant genome assembly statistics.** A summary table of the number of possible mutant viruses, percent theoretical coverage, the number of viable viruses recovered from one transformation, and the number of mutant and wild type genomes recovered.

Fragment	F1	F2	F3	G1	G2	WT
Number of possible mutants	1,900	2,880	3,640	3,440		NA
Percent theoretical coverage	45	9	55	94		NA
Number of viable transformants	2,300	460	5,840	19,200		6,400
Number of plaques sequenced	34	25	30	60		30
Number of WT plaques	17	11	15	30		30

## Discussion

Here we present the first steps in the development of a novel method for generating comprehensive single-site saturation – and massive multi-site – virulent mutant libraries of the spike and capsid proteins of the bacteriophage ΦX174. This method uses unamplified oligo pools, nicking scanning mutagenesis, and Golden Gate cloning. We are fortunate that our model virus has a small genome, under 6,000 nts, which we predict will enable us to assemble the 15

non-toxic fragments with sufficient transformants for near-complete coverage of the user-defined mutations. However, most biotech-relevant phages have larger genomes, and extension of this method to other viruses would encounter the following technical challenges. First, increasing genome sizes decrease transformational efficiencies, and larger genomes have increased susceptibility to DNA shearing<sup>33</sup>. Both size-dependent effects can disrupt the integrity and coverage of the phage libraries<sup>34</sup>. Second, our method includes restrictions on the DNA sequences that can be used. Golden Gate cloning requires a genome without unique Type IIS restriction sequences, while nicking mutagenesis requires that there be only one orientation of the BbvCI nicking site within the template DNA fragment. Third, ligating 15 fragments is close to the upper limit of Golden Gate cloning. Additionally, increasing the size of individual fragments is difficult as the presented method depends on replicating plasmids in *E. coli*, and larger fragments are more likely to be toxic *in vivo*.

Based on the above considerations, adapting this method to larger viruses will require modifications for the genome assembly steps. There are a variety of methods available for attempting to improve genome library assembly and amplification. For improving genome assembly, we speculate that a combination of hierarchical assembly<sup>35</sup> with other yeast based assembly methods<sup>36</sup> will allow for large viral genomes to be efficiently assembled. Finally, we expect the Tx.Tl cell free expression system<sup>37</sup> will be able to produce the viral libraries with less bias than *in vivo* amplification in *E. coli*. In summary, we have presented a complete method for efficient deep mutational scanning of the bacteriophage ΦX174. We anticipate this method will find utility in fundamental molecular evolution studies as well as translate to potential medicinal applications.

## Materials and methods

### *Reagents*

All purchased enzymes and DNA purification kits were from New England Biolabs, antibiotics were purchased from GoldBio, other chemicals were purchased from Sigma-Aldrich. Individual primers were purchased from Integrated DNA Technologies.

### *Segmentation of the $\Phi$ X174 genome*

A phage assembly platform for  $\Phi$ X174 was devised following Donaldson et al.<sup>32</sup>. The  $\Phi$ X174 chromosome was divided into 15 genomic fragments designed to avoid host cell toxicity by separating genes from their promoters and breaking large genes into multiple segments (**Note B 4.1**). Each segment is flanked by unique five nucleotide overlaps of WT  $\Phi$ X174 sequence so that they can be amplified from the ancestral  $\Phi$ X174 using PCR primers designed to incorporate terminal BsmB1 restriction sites. Amplicons were cloned into pCR2.1 using the Invitrogen TOPO TA cloning system (Life Technologies, Grand Island, NY).

### *Introduction of nicking site*

The gene fragments – F1, F2, G1, G2 – in the pCR2.1-TOPO plasmid (**Notes B 4.1 and 4.2**) had the BbvCI nicking site introduced via overhang PCR, type I restriction enzyme cutting, and ligation. First, PCR was performed with overhang primers (**Table B 2**) to introduce the BbvCI site. Standard Phusion Polymerase HF reaction conditions were used with 4 ng of template DNA, cycling is follows: 98°C - 1 min, 25x cycles of: 98°C – 10 seconds, 67°C 15 seconds, 72°C 2.5 minutes, followed by 72°C for 10 minutes. PCR products were run on a 1% agarose gel stained with SYBR™ safe stain (Invitrogen) and the DNA bands at ~4600 bp were

extracted using a Monarch® DNA Gel Extraction Kit. Next, 1 µg of the PCR product was digested with KpnI (20 U) in NEB Buffer 1.1 at 37°C for two hours. Digested DNA was then clean and concentrated with a Monarch® PCR and DNA Clean and Concentrate Kit and eluted into 20 µL nuclease free H<sub>2</sub>O. One microliter of the purified and digested DNA was ligated using T4 DNA ligase at ~25°C for 1 hour in standard conditions in a 20 µL reaction. Five microliters of the ligation reaction were transformed into chemically competent XL1-Blue *E. coli* via standard protocols. Cells were plated on LB agar containing 100 µg/mL carbenicillin and 50 µg/mL kanamycin and grown at 37°C for ~16 hours. Cells were picked from transformation plates and grown in 50 mL TB with 100 µg/mL carbenicillin and 50 µg/mL kanamycin for ~12 hours, cells were pelleted, and DNA purified using compact midi-preps.

#### *Comprehensive single site mutant library construction*

Comprehensive mutant libraries were generated using nicking mutagenesis (NM) as in Wrenbeck et al.<sup>20</sup> with modifications for using oligo pool mutagenic primers as noted in Medina et al.<sup>21</sup>. A single oligo pool encoding for all possible single missense and nonsense substitutions in F and G was designed using the custom python scripts from Medina et al.<sup>21</sup> and custom synthesized by Agilent (full sequences of all oligos are given in the associated file “VMA\_supplemental\_data.csv” in the published text). Oligo pools were designed with 20-24 bases of gene overlap flanking the mutated codon. Codons were chosen based on *E. coli* codon usage frequency. This oligo pool was used directly in NM without further amplification and using 2 mg of the relevant golden gate plasmid as a template. Libraries were transformed into high efficiency electrocompetent XL1-Blue *E. coli* (Agilent cat #: 200228) using 1 mm electroporation cuvettes at 1200V. Cells were plated on large bioassay plates (245mm x 245mm

x 25mm, Sigma-Aldrich) containing LB agar + 100 µg/mL carbenicillin and 50 µg/mL kanamycin and grown at 37°C for ~16 hours. Plates were scraped in 10 mL plain LB, broken into ~1.2 mL aliquots, pelleted, and stored -80°C. DNA was purified from 1x aliquot using a Monarch® Plasmid Mini-Prep Kit.

#### *Illumina sequencing prep and analysis*

Purified library DNA was prepared for deep sequencing as in Kowalsky et al.<sup>38</sup> using the primers listed in **Table B 3** and gene tiling as specified in **Table 4.1**. DNA was Illumina sequenced on a Mi-Seq platform with 250 BP paired end reads. The University of Colorado BioFrontiers Sequencing Core performed the Illumina sequencing. Data was processed using PACT<sup>39</sup> to determine the library coverage with the following changes to the default options in the configuration file: fast\_filter\_translate: qaverage = 20, qlimit = 0; enrichment: ref\_count\_threshold = 5, sel\_count\_threshold = 0, strict\_count\_threshold = True. The heatmaps of counts can be found in **Figure B 2 - 9**.

#### *Assembly of mutant genomes*

We pooled plasmid DNA containing all 15 of the phage DNA fragments in equimolar amounts and digested them with BsmB1 (Fermentas Fast Digest, Life Technologies, Grand Island, NY) for 30 minutes to 1 hour at 37°C. The digested plasmids were subjected to agarose gel electrophoresis for 10 to 15 minutes using a 1.2% agarose gel to separate the vector from the inserts. The inserts were excised from the gel, purified using the GeneJET gel extraction kit (Fermentas), ligated overnight at 14°C with T4 DNA ligase (Promega Corporation, Madison, WI), and transformed by electroporation into 100 µl competent *E. coli* C cells. The

transformation mix was resuspended with 1 ml of ΦLB and either plated immediately or incubated for about 20 minutes at 37°C to allow for one viral burst. The ΦLB was added to 3 ml of ΦLB top agar and plated onto a ΦLB agar plate. After four to five hours of incubation at 37°C, recombinant phage plaques were visible and plates were removed from the incubator. To verify that the recombinant phage encoded the intended sequence, we picked about 30 plaques for each intended mutational target and Sanger sequenced the entire targeted gene. We also sequenced the F and G genes for about 30 wild type plaques to assure that no mutations were naturally accumulating. Briefly, individual plaques were picked with sterile toothpicks and placed in 200 uL ΦLB and gently swirled. 1 uL of this mix was used to PCR amplify approximately ½ of the ΦX174 genome using ΦX-0F (5'-GAGTTTATCGCTTCATG-3') and ΦX-2953R (5'-CCGCCAGCAATAGCACC-3') primers. Internal sequencing primers ΦX-979F (5'-CGGCCCTTACTTGAGG-3') and ΦX-1500R (5'-TTGAGATGGCAGCAACGG-3') were used to sequence gene F. ΦX-2953R was used to sequence gene G. PCR cleanups and sequencing was done at Eurofins Genomics. PCR reaction conditions were; 5 uL10X Taq buffer, 2.5 uL 10 uM ΦX-0F primer, 10 uM ΦX-2953F primer, 0.8 uL 12.5 uM dNTPs, 0.5 uL Taq polymerase (NEB #M0273), 1 uL template, 37.7 uL H<sub>2</sub>O. Thermocycling conditions were 1 cycle at 95°C for 2 min, 30 cycles at 95°C for 15 sec, 52°C for 30 sec, 68°C for 2 min, 1 cycle at 68°C for 5min.

## **REFERENCES**

## REFERENCES

1. Holmes EC: **What can we predict about viral evolution and emergence?**. *Curr Opin in Vir* 2013, 3:180-184.
2. Neher RA: **Genetic draft, selective interference, and population genetics of rapid adaptation.** *Annu Rev Ecol Evol Syst* 2013, 44:195-215.
3. Lieberman TD, Flett KB, Yelin I, Martin TR, McAdam AJ, Priebe GP, Kishony R: **Genetic variation of a bacterial pathogen within individuals with cystic fibrosis provides a record of selective pressures.** *Nature Gen* 2014, 46:82-87.
4. Dedrick RM, Guerrero-Bustamante CA, Garlena RA, Russell DA, Ford K, Harris K, Gilmour KC, Soothill J, Jacobs-Sera D, Schooley RT, Hatfull GF, Spencer H: **Engineered bacteriophages for treatment of a patient with a disseminated drug-resistant *Mycobacterium abscessus*.** *Nature Medicine*. 2019, 25:730-733.
5. Barrick JE, Lenski RE: **Genome dynamics during experimental evolution.** *Nat Rev Genetics* 2013, 14:827-839.
6. Fowler DM, Fields S: **Deep mutational scanning: a new type of protein science.** *Nat Methods* 2014, 11:801-807.
7. Pires DP, Cleto S, Sillankorva S, Azeredo J, Lu TK: **Genetically engineered phages: a review of advances over the last decade.** *Microbiol and Mol Biol Rev* 2016, 80:523-543.
8. Kiro R, Shitrit D, Qimron U: **Efficient engineering of a bacteriophage genome using the type I-E CRISPR-Cas system.** *RNA Biol* 2014, 11:42-44.
9. Lemay ML, Tremblay DM, Moineau S: **Genome engineering of virulent lactococcal phages using CRIPSR-Cas9.** *ACS Synth Biol* 2017, 6:1351-1358.
10. Schilling T, Dietrich S, Hoppert M, Hertel R: **A CRISPR-Cas9-based toolkit for fast and precise in vivo genetic engineering of *Bacillus subtilis* phages.** *Viruses* 2018, 10:241.
11. Doore SM, Fane BA: **The kinetic and thermodynamic aftermath of horizontal gene transfer governs evolutionary recovery.** *Mol. Biol. Evo.* 2015, 32:2571-2584.
12. Doore SM, Schweers NJ, Fane BA: **Elevating fitness after a horizontal gene exchange in bacteriophage ΦX174.** *Virology* 2017, 501:25-34.
13. Uchiyama A, Heiman P, Fane BA: **N-terminal deletions of the ΦX174 external scaffolding protein affect the timing and fidelity of assembly.** *Virology* 2009, 386:303-309.
14. Stockdale SR, Collins B, Silvia S, Douillard FP, Mahony J, Cambillau C, van Sinderen D: **Structure and assembly of TP901-1 Virion unveiled by mutagenesis.** *PLOS ONE* 2015,

10:e0131676.

15. Lee JM, Huddleston J, Doud MB, Hooper KA, Wu NC, Bedford T, Bloom JD: **Deep mutational scanning of hemagglutinin helps predict evolutionary fates of human H3N2 influenza variants.** *Proc. Natl. Aca. Sci. USA* 2018, 115:E8276-E8285.
16. Doub MB, Lee JM, Bloom JD: **How single mutations affect viral escape from broad and narrow antibodies to H1 influenza hemagglutinin.** *Nature Comm.* 2018, 9:1386.
17. Wu NC, Qi H: **Application of deep mutational scanning in hepatitis C virus.** *Met. in Mol. Biol.* 2019, 1911:183-190.
18. Hartman EC, Jakobson CM, Favor AH, Lobba MJ, Álvarez-Benedicto E, Francis MB, Tullman-Ercek D: **Quantitative characterization of all single amino acid variants of a viral capsid based drug delivery vehicle.** *Nature Comm.* 2018, 9:1385.
19. van Dolleweerd CJ, Kessans SA, Van de Bittner KC, Bustamante LY, Bundela R, Scott B, Nicholson MJ, Parker EJ: **MIDAS: a modular DNA assembly system for synthetic biology.** *ACS Synth Biol* 2018, 7:1018-1029.
20. Wrenbeck EE, Klesmith JR, Stapleton JA, Adeniran A, Tyo KEJ, Whitehead TA: **Plasmid-based one-pot saturation mutagenesis.** *Nat. Methods* 2016, 13:928-930.
21. Medina-Cucurella A, Steiner PJ, Faber MS, Beltrán J, Borelli A, Kirby MB, Cutler SR, Whitehead TA: **User-defined single pot mutagenesis using unpurified oligo pools.** *Prot. Eng. Des. Sel.* 2019, TBD:TBD.
22. Kosuri S, Church GM: **Large-scale *de novo* DNA synthesis: technologies and applications.** *Nat. Methods* 2014, 11:499-507.
23. Aoyama A, Hamatake RK, Hayashi M: **Morphogenesis of ΦX174: *in vitro* synthesis of infectious phage from purified viral components.** *Proc. Natl. Aca. Sci. USA* 1981, 12:7285-7289.
24. Hafenstein S, Fane BA: **ΦX174 genome-capsid interactions influence the biophysical properties of the virion: evidence for a scaffolding-like function for the genome during the final stages of morphogenesis.** *J. of Viro.* 2002, 76:5350-5356.
25. Bernal RA, Hafenstein S, Esmeralda R, Fane BA, Rossmann MG: **The ΦX174 protein J mediates DNA packaging and viral attachment to host cells.** *J. Mol. Biol.* 2004, 337:1109-1122.
26. Rokyta DR, Burch CL, Caudle SB, Wichman HA: **Horizontal gene transfer and the evolution of microvirid coliphage genomes.** *J. of Bacter.* 2006, 188:1134-1142.
27. McKenna R, Xia D, Willingmann R, Ilag LL, Krishnaswamy S, Rossmann MG, Olson NH, Baker TS, Incardona NL: **Atomic structure of single-stranded DNA bacteriophage ΦX174 and its functional implications.** *Nature* 1992, 355:137-143.

28. Dokland T, McKenna R, Ilag LL, Bowman BR, Incardona NL, Fane BA, Rossmann MG: **Structure of a viral procapsid with molecular scaffolding.** *Nature* 1997, 389:308-313.
29. Sanger F, Coulson AR, Friedmann T, Air GM, Barrell BG, Brown NL, Fiddes JC, Hutchinson CA 3<sup>rd</sup>, Slocombe PM, Smith M: **The nucleotide sequence of bacteriophage phiX174.** *J Mol Biol* 1978, 125:225-246.
30. Smith HO, Hutchinson III CA, Pfannkoch C, Venter JC: **Generating a synthetic genome by whole genome assembly: ΦX174 bacteriophage from synthetic oligonucleotides.** *Proc. Natl. Aca. Sci. USA* 2003, 100:15440-15445.
31. Jaschke PR, Lieberman EK, Rodriguez J, Sierra A, Endy D: **A fully decompressed synthetic bacteriophage ΦX174 genome assembled and archived in yeast.** *Virology* 2012, 434:278-284.
32. Donaldson EF, Yount B, Sims AC, Burkett S, Pickles RJ, Baric RS: **Systematic Assembly of a Full-Length Infectious Clone of Human Coronavirus NL63.** *J. Virol.* 2008, 82:11948–11957.
33. Rogers SO, Bendich AJ: **Extraction of DNA from plant tissues.** *Plant Mole Biol Manual* 1988, A6:1-10.
34. Edwards RA, Rohwer F: **Viral metagenomics.** *Nature Reviews Micro* 2005, 3:504-509.
35. Lee, ME, DeLoache WC, Cervantes B, Dueber JE: **A highly characterized yeast toolkit for modular, multipart assembly.** *ACS Synth Biol* 2015, 4:975-986.
36. Ando H, Lemire S, Pires DP, Lu TK: **Engineering modular viral scaffolds for targeted bacterial population editing.** *Cell Systems* 2015, 1:187-196.
37. Garamella J, Marshall R, Rustad M, Noireaux V: **The all *E. coli* Toolbox 2.0: a platform for cell-free synthetic biology.** *ACS Synth Biol* 2016, 5:344-355.
38. Kowalsky CA, Klesmith, JR, Stapleton JA, Kelly V, Reichkitzer N, Whitehead TA: **High-Resolution Sequence-Function Mapping of Full-Length Proteins.** *PLoS One* 2015, 10:e0118193.
39. Klesmith JR, Hackel BJ: **Improved mutation function prediction via PACT: Protein Analysis and Classifier Toolkit.** *Bioinformatics* 2018, doi/10.1093/bioinformatics/bty1042.

## **CHAPTER 5**

**Conclusions and perspectives**

This thesis focuses on the generation and analysis of large mutational datasets for understanding molecular evolution. Deep mutational scanning (DMS) is a technique that combines saturation mutagenesis libraries with high throughput selection and deep sequencing to assess the impact of nearly all single-point mutations on protein function<sup>1</sup>. Here, DMS is used to study how folding probability constrains the molecular evolution of enzymes. This thesis also expands the research questions DMS experiments can address through a novel method for generating mutant genome libraries of virulent phages. This chapter provides a summary of the work performed and a discussion of the broader impacts of this PhD thesis.

DMS experiments produce large datasets that quantitatively describe the impacts of single mutations<sup>2-5</sup>. In **Chapter 2**, we reviewed the use of DMS datasets - as well as datasets from several other experimental and evolutionary sources – in the forward engineering of protein therapeutics. We describe the use of these datasets in the identification of mutations that can tune the specificity, affinity, and solubility of a target protein. As more comprehensive studies of mutational impacts are performed, we expect the integration of the resultant data to provide better predictions of which mutations will, and will-not, work toward a specific goal. Eventually, we expect guidelines like Lipinski’s rule of five for small molecule drugs<sup>6</sup> to emerge for protein therapeutics. Indeed, work aiming to establish these guidelines are already underway<sup>7</sup>.

In **Chapter 3**, I demonstrated how DMS can be used to understand how a single biophysical parameter of an enzyme constrains its molecular evolution. In this study I performed near-comprehensive single-mutation analysis on two unique single-point mutants of the aliphatic amidase AmiE from *P. aeruginosa*. The point-mutants have decreased *in vivo* folding probabilities and statistically indistinguishable catalytic activities compared to the starting enzyme. In a previous study<sup>5</sup> the starting enzyme had undergone DMS with the same selection

conditions applied in this study, allowing us to compare the obtained datasets. Because the only biophysical parameter modulated is the *in vivo* folding probability, we are able to assess how only this parameter alters the types of mutations that can arise. Additionally, this study provides a near comprehensive analysis of how single-point mutations combine, providing a detailed view of epistasis, the non-additive combination of mutational effects<sup>8</sup>. The novelty of this study is found in the aforementioned details and in the fact that this study is the most comprehensive analysis of how a single initial biophysical parameter impacts mutational outcomes to date.

Comparing the datasets obtained for our two disrupted enzymes with that of the initial starting enzyme provides intriguing insights into molecular evolution. First, in all genetic backgrounds the distributions of beneficial fitness effects are described by the General Pareto distribution. Thus, the likelihood of uncovering a beneficial mutation is insensitive to initial *in vivo* folding probability in our experiments. Interestingly, most beneficial mutations are shared, and only 19 mutations are likely to require an increased *in vivo* folding probability to buffer stability penalties.

Epistasis was found in both the beneficial and deleterious populations of mutations following decreases in the *in vivo* folding probabilities. Interestingly, in the population of double mutants that are deleterious individually and beneficial when combined, phenotype-specific interactions dominated those that are phenotype-independent. This suggests that mutations rescuing *in vivo* folding probability are more likely to be in direct contact with the disrupting mutation, and less likely to alter the global properties of the enzyme. Additionally, beneficial mutations with synonymous codon fitness disparities - opposite mutational outcomes for synonymous codons - were correlated with background specific beneficial mutations. This

correlation reveals the importance of the mRNA sequence selected, and how it significantly restricts mutational outcomes for identical peptides.

Taken together, these results indicate that models of protein stability and evolution must include all aspects of the protein life cycle, from transcription to degradation, to fully understand the impact of a single mutation. It is uncertain if these findings are unique to this experimental system, or if they are descriptive of protein evolution in general. Additional studies similar to this one, using other proteins and with different biophysical parameters altered, will improve our ability to predict mutational outcomes.

The final goal of this thesis was methodological. In **Chapter 4** we expanded the experimental range of DMS by generating comprehensive single-site saturation, and massive multi-site, mutant genome libraries of the bacteriophage  $\phi$ X174. The mutant genomes were generated by performing nicking scanning mutagenesis<sup>9</sup> on non-toxic truncations of the genome and assembling complete viral genomes with Golden Gate cloning. These mutant genomes will allow researchers to probe how host specificity is encoded into the primary sequence, how intra- and inter-molecular epistasis is involved in viral evolution and micro-compartment formation, and expand our understanding of the molecular evolution of viruses.

In conclusion, this thesis provides significant insights into molecular evolution and how to study it. All aspects of this thesis involve the generation and/or analysis of large mutational datasets with the goal of understanding how proteins change over evolutionary time. In the coming decades the mutational data obtained from more DMS experiments, and from other experimental and archival sources, will provide valuable insights into molecular evolution. I expect this information will be combined and integrated into a new generation of predictive algorithms, which will significantly improve the protein engineering process.

## **REFERENCES**

## REFERENCES

1. Fowler DM, Fields S: **Deep mutational scanning: a new style of protein science.** *Nat Methods* 2014, 11:801–807.
2. Kowalsky CA, Klesmith JR, Stapleton JA, Kelly V, Reichkitter N, Whitehead TA: **High-Resolution Sequence-Function Mapping of Full-Length Proteins.** *PLoS One* 2015, 10:e0118193.
3. Klesmith JR, Bacik JP, Michalczyk R, Whitehead TA: **Comprehensive Sequence-Flux Mapping of a Levoglucosan Utilization Pathway in *E. coli*.** *ACS Synth Biol* 2015, 4:1235–1243.
4. Klesmith JR, Bacik JP, Wrenbeck EE, Michalczyk R, Whitehead TA: **Trade-offs between enzyme fitness and solubility illuminated by deep mutational scanning.** *Proc Natl Acad Sci U S A* 2017, 114:2265-2270.
5. Wrenbeck EE, Azouz LR, Whitehead TA: **Single-mutation fitness landscapes for an enzyme on multiple substrates reveal specificity is globally encoded.** *Nat Comm* 2017, 8:15695.
6. Lipinski CA: **Lead-and drug-like compounds: the rule-of-five revolution.** *Drug Discov Today Technol* 2004, 1:337-341.
7. Rabia LA, Desai AA, Jhajj HS, Tessier PM: **Understanding and overcoming trade-offs between antibody affinity, specificity, stability and solubility.** *Biochem Eng J* 2018, 137:365-374.
8. Starr TN, Thornton JW: **Epistasis in protein evolution.** *Protein Science* 2016, 25:1204-1218.
9. Wrenbeck EE, Klesmith JR, Stapleton JA, Adeniran A, Tyo KEJ, Whitehead TA: **Plasmid-based one-pot saturation mutagenesis.** *Nat Methods* 2016, 13:928-930.

## **APPENDICES**

## **APPENDIX A**

### **Chapter 3 supporting information**

**Note A 1: Amino acid sequences of the AmiE variants.** Mutations are highlighted in red.

>AmiE\_WT

MRHGDISSNDTVGVAVVNYKMPRLHTAAEVLDNARKIAEMIVGMKQGLPGMDLVVF  
PEYSLQGIMYDPAEMMETAVAIPGEETEIFSRACRKANWGVFSLTGERHEEHPRKAPY  
NTLVLIIDNNGEIVQKYRKIIPWCPIEGWYPGGQTYVSEGPKGKISLIICDDGNYPEIWRD  
CAMKGAEELIVRCQGYMYPAKDQQVMMAKAMAWANNCYVAVANAAGFDGVYSYFG  
HSAIIGFDGRTLGECEEEEMGIQYAQLSLSQIRDARANDQSQNHLFKILHRGYSGLQASG  
DGDRGLAECPFEFYRTWVTDAEKARENVERLRTSTTGVHQCPVGRLPYEGLEHHHHH

>AmiE\_I122L

MRHGDISSNDTVGVAVVNYKMPRLHTAAEVLDNARKIAEMIVGMKQGLPGMDLVVF  
PEYSLQGIMYDPAEMMETAVAIPGEETEIFSRACRKANWGVFSLTGERHEEHPRKAPY  
NTLVLLDNNGEIVQKYRKIIPWCPIEGWYPGGQTYVSEGPKGKISLIICDDGNYPEIWR  
DCAMKGAEELIVRCQGYMYPAKDQQVMMAKAMAWANNCYVAVANAAGFDGVYSYFG  
GHSAIIGFDGRTLGECEEEEMGIQYAQLSLSQIRDARANDQSQNHLFKILHRGYSGLQASG  
GDGDRGLAECPFEFYRTWVTDAEKARENVERLRTSTTGVHQCPVGRLPYEGLEHHHHH  
H

>AmiE\_I38V

MRHGDISSNDTVGVAVVNYKMPRLHTAAEVLDNARKVAEMIVGMKQGLPGMDLVVF  
FPEYSLQGIMYDPAEMMETAVAIPGEETEIFSRACRKANWGVFSLTGERHEEHPRKAP  
YNTLVLIIDNNGEIVQKYRKIIPWCPIEGWYPGGQTYVSEGPKGKISLIICDDGNYPEIWR  
DCAMKGAEELIVRCQGYMYPAKDQQVMMAKAMAWANNCYVAVANAAGFDGVYSYFG  
GHSAIIGFDGRTLGECEEEEMGIQYAQLSLSQIRDARANDQSQNHLFKILHRGYSGLQASG  
GDGDRGLAECPFEFYRTWVTDAEKARENVERLRTSTTGVHQCPVGRLPYEGLEHHHHH  
H

**Note A 2: DNA sequences of the AmiE variants.** Codons mutated are underlined and bold, while mutations are highlighted in red.

>AmiE\_WT

ATGAGACATGGCGATATTAGCTCGTCAAATGATACCGTAGGCGTAGCCGTGGTGAA  
TTACAAGATGCCCGTTACATACTGCTGAAGTCCTGGATAATGCCGAAAAT  
TGC~~GG~~AAATGATCGTTGGTATGAAGCAAGGTCTGCCGGCATGGATCTGGTTGTGTT  
TCCTGAATATTCTTACAGGGTATTATGTACGACCCTGCTGAAATGATGGAAACAGC  
CGTGGCGATTCCAGGC~~A~~AGAACGGAAATCTTAGCCGTGCTGTAGAAAAGCAA  
ATGTTGGGGTGTGTTCTCCCTGACCGCGA~~C~~ACGT~~C~~ATGAAGAACACCC~~T~~AGAAAGG  
CACCATACAACACTCTGGTCTTGATCGATAACAAACGGTGAAATCGTACAAAAGTAC  
AGAAAGATCATCCCATGGTGTCCGATTGAAGGCTGGTATCCAGGTGGCCAGACATA  
CGTCTGAAGGTCCGAAAGGCATGAAGATCTCATTAAATTATCG~~G~~ATGACGGTAA  
TTATCCGAAATTGGAGAGATTGTGCCATGAAGGGT~~G~~CGGAATTGATCGTTGCTG  
CCAAGGCTATATGTACCCTGCTAAAGACCAACAAAGTTATGATGGCTAAGGCAATGG  
CCTGGCGAATAACTGTTATGTCGCTGTAGCAAACGCTGCAGGTTTGATGGCGTT  
ATAGCTACTTCGGTCATAGTGCATTATCGGTTTGACGGCCGTACTCTGGTGAAT  
GCGCGAAGAACGGAAATGGCATTCAATACGCGCAGTTGTCTGT~~C~~ACAAATCCGC  
GATGCCGTGCGAATGACCAAAGTCAGAACCATTGTTAAAATCTGCACAGAGGT  
TACTCCGGTTG~~C~~AGGCTTCGGCGATGGC~~G~~ACCGTGGTCTGGCAGAATGCCATT  
GAATTCTACCGTACCTGGGTTACTGATGCTGAAAAGGCAAGAGAAAACGTGGAACG  
CCTGACTCGCTCCACAAACAGGTGTCGCCAATGCCAGTAGGTCG~~T~~GTGCC~~T~~ATGA  
AGGCCTCGAGCACCACCACCAC

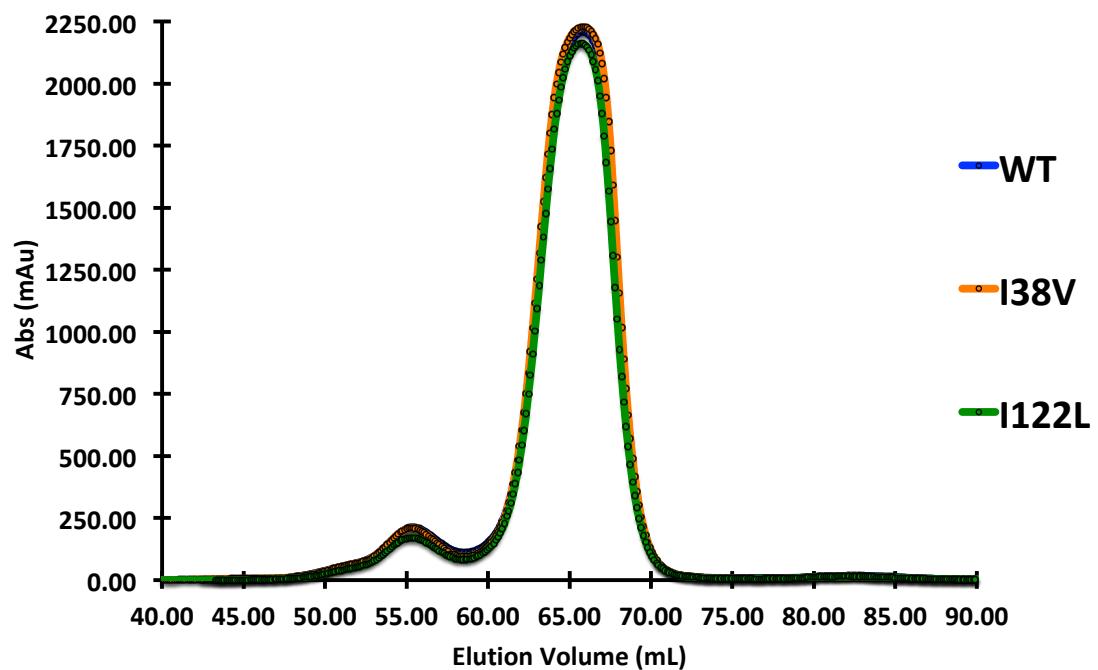
>AmiE\_I122L

ATGAGACATGGCGATATTAGCTCGTCAAATGATACCGTAGGCGTAGCCGTGGTGAA  
TTACAAGATGCCCGTTACATACTGCTGAAGTCCTGGATAATGCCGAAAAT  
TGC~~GG~~AAATGATCGTTGGTATGAAGCAAGGTCTGCCGGCATGGATCTGGTTGTGTT  
TCCTGAATATTCTTACAGGGTATTATGTACGACCCTGCTGAAATGATGGAAACAGC  
CGTGGCGATTCCAGGC~~A~~AGAACGGAAATCTTAGCCGTGCTGTAGAAAAGCAA  
ATGTTGGGGTGTGTTCTCCCTGACCGCGA~~C~~ACGT~~C~~ATGAAGAACACCC~~T~~AGAAAGG  
CACCATACAACACTCTGGTCTTG~~C~~TCGATAACAAACGGT~~G~~AAATCGTACAAAAGTAC  
AGAAAGATCATCCCATGGTGTCCGATTGAAGGCTGGTATCCAGGTGGCCAGACATA  
CGTCTGAAGGTCCGAAAGGCATGAAGATCTCATTAAATTATCG~~G~~ATGACGGTAA  
TTATCCGAAATTGGAGAGATTGTGCCATGAAGGGT~~G~~CGGAATTGATCGTTGCTG  
CCAAGGCTATATGTACCCTGCTAAAGACCAACAAAGTTATGATGGCTAAGGCAATGG  
CCTGGCGAATAACTGTTATGTCGCTGTAGCAAACGCTGCAGGTTTGATGGCGTT  
ATAGCTACTTCGGTCATAGTGCATTATCGGTTTGACGGCCGTACTCTGGTGAAT  
GCGCGAAGAACGGAAATGGCATTCAATACGCGCAGTTGTCTGT~~C~~ACAAATCCGC  
GATGCCGTGCGAATGACCAAAGTCAGAACCATTGTTAAAATCTGCACAGAGGT  
TACTCCGGTTG~~C~~AGGCTTCGGCGATGGC~~G~~ACCGTGGTCTGGCAGAATGCCATT  
GAATTCTACCGTACCTGGGTTACTGATGCTGAAAAGGCAAGAGAAAACGTGGAACG

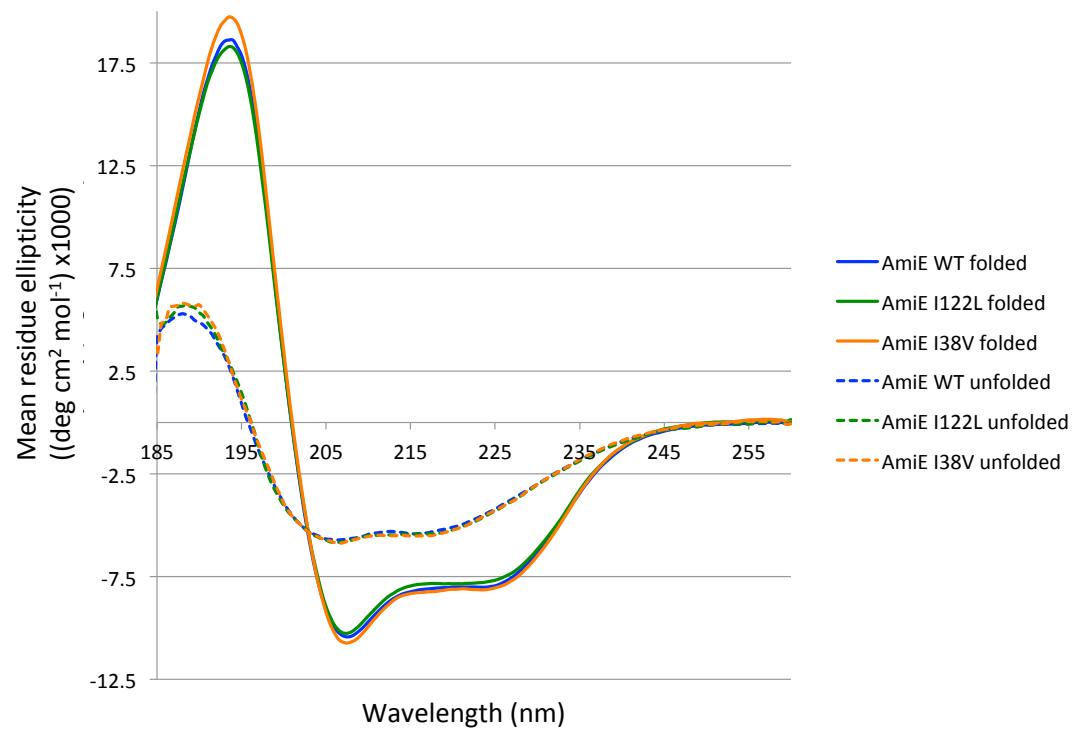
CCTGACTCGCTCCACAACAGGTGTCGCCAATGCCAGTAGGTCGTCTGCCGTATGA  
AGGCCTCGAGCACCACCACCAACCAC

>AmiE\_I38V

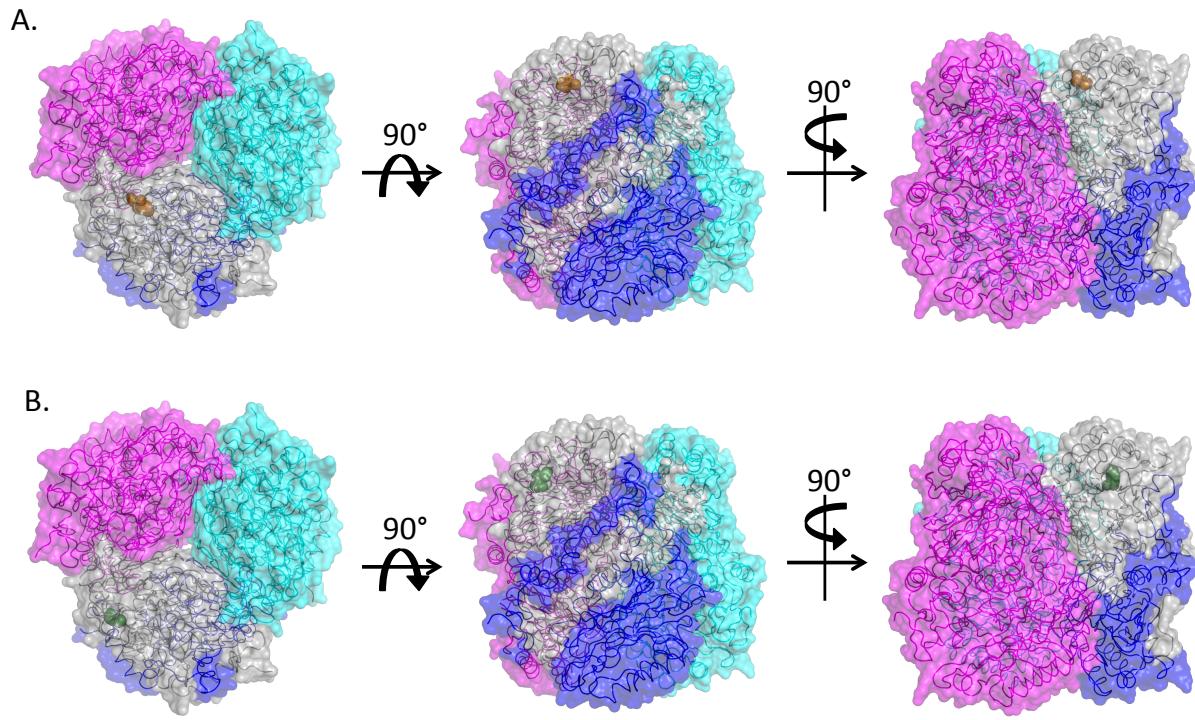
ATGAGACATGGCGATATTAGCTCGTCAAATGATACCGTAGGCGTAGCCGTGGTGAA  
TTACAAGATGCCCGTTACATACTGCTGCTGAAGTCCTGGATAATGCCGCAAAGT  
TCGGGAAATGATCGTGGTATGAAGCAAGGTCTGCCGGCATGGATCTGGTGTGTT  
TCCTGAATATTCTTACAGGGTATTATGTACGACCCCTGCTGAAATGATGGAAACAGC  
CGTGGCGATTCCAGGCAGAACGGAAATCTTAGCCGTGTTGTAGAAAAGCAA  
ATGTTGGGTGTGTTCTCCCTGACCGCGAACGTATGAAGAACACCCCTAGAAAGG  
CACCATACAACACTCTGGTCTGATCGATAACAACGGTGAATCGTACAAAAGTAC  
AGAAAGATCATCCCATTGGTGTCCGATTGAAGGCTGGTATCCAGGTGCCAGACATA  
CGTCTCTGAAGGTCCGAAAGGCATGAAGATCTCATTAATTATCTGCATGACGGTAA  
TTATCCGAAATTGGAGAGATTGTGCCATGAAGGGTGCAGGAAATTGATCGTTCGCTG  
CCAAGGCTATATGTACCCCTGCTAAAGACCAACAAGTTATGATGGCTAAGGCAATGG  
CCTGGCGAATAACTGTTATGTCGCTGTAGCAAACGCTGCAGGTTTGATGGCGTT  
ATAGCTACTCGGTCAAGTGCCATTATCGGTTTGACGGCGTACTCTGGTGAAT  
GCGCGAAGAAGAAATGGGCATTCAATACGCGCAGTTGTCTGTCACAAATCCGC  
GATGCCCGTGCAGAACCAAAGTCAGAACCAAGTTAAAATCTGCACAGAGGT  
TACTCCGGTTGCAGGCTCGGGCGATGGCGACCGTGGTCTGGCAGAATGCCATT  
GAATTCTACCGTACCTGGGTTACTGATGCTGAAAGGCAAGAGAAAACGTGGAACG  
CCTGACTCGCTCCACAACAGGTGTCGCCAATGCCAGTAGGTCGTCTGCCGTATGA  
AGGCCTCGAGCACCACCACCAACCAC



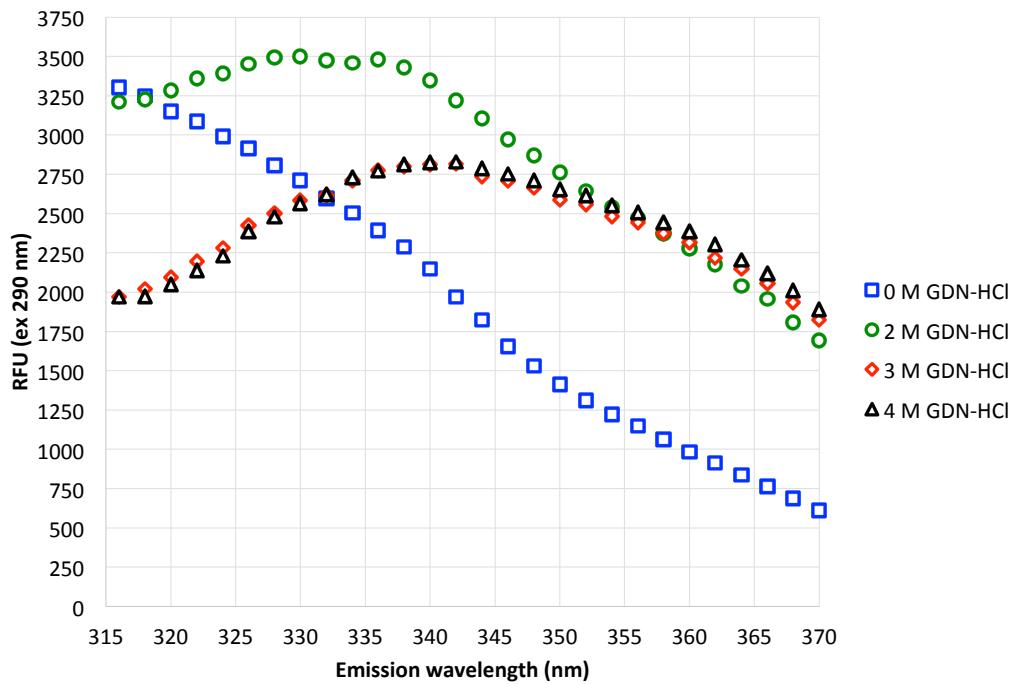
**Figure A 1: Chromatogram of purified AmiE variants.** SEC-FPLC chromatograms of AmiE proteins run on a HiLoad 16/600 Superdex 200 column at 1 mL/min in PBS as the mobile phase. No gross differences in oligomeric state were determined for the proteins.



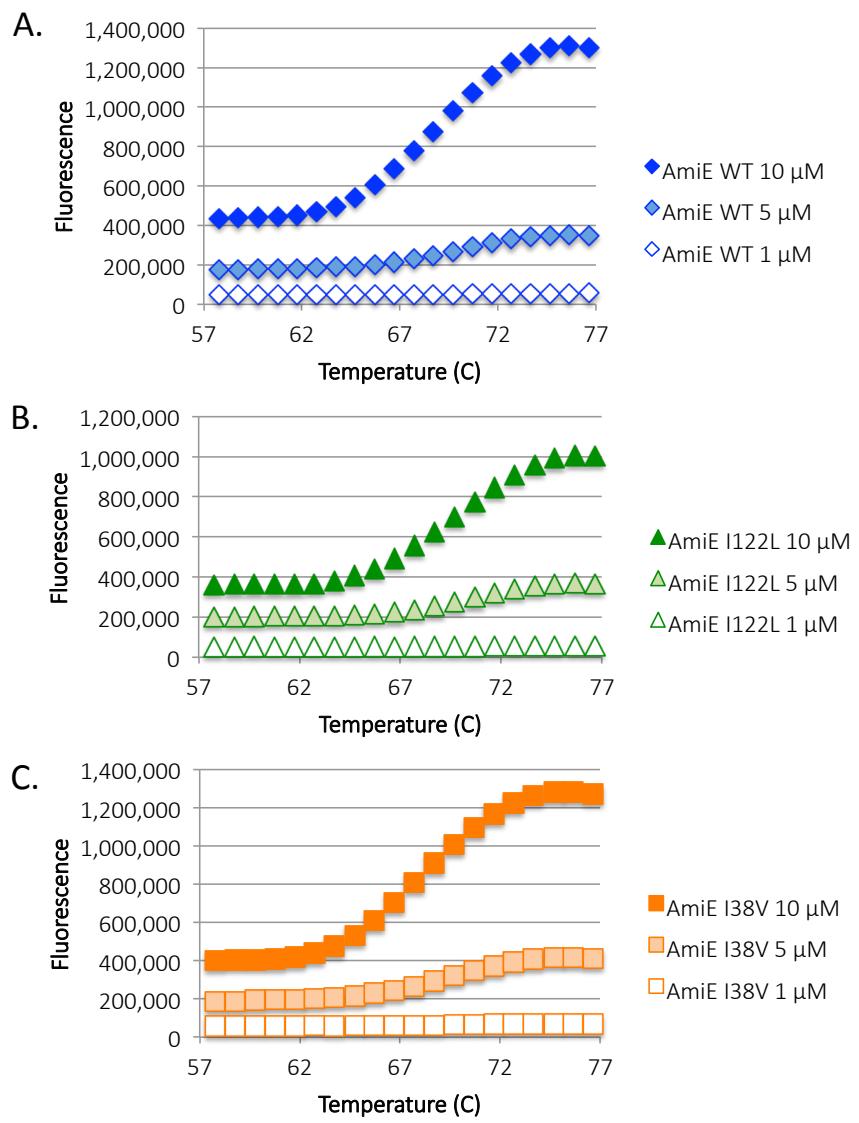
**Figure A 2: Representative far-UV spectra of the folded and unfolded AmiE variants.** Far-UV spectra of folded and unfolded 10  $\mu\text{M}$  AmiE WT, 10  $\mu\text{M}$  AmiE I122L, 10  $\mu\text{M}$  AmiE I38V at 15°C (folded) and 90°C (unfolded) in 10 mM phosphate buffer at pH 7.5.



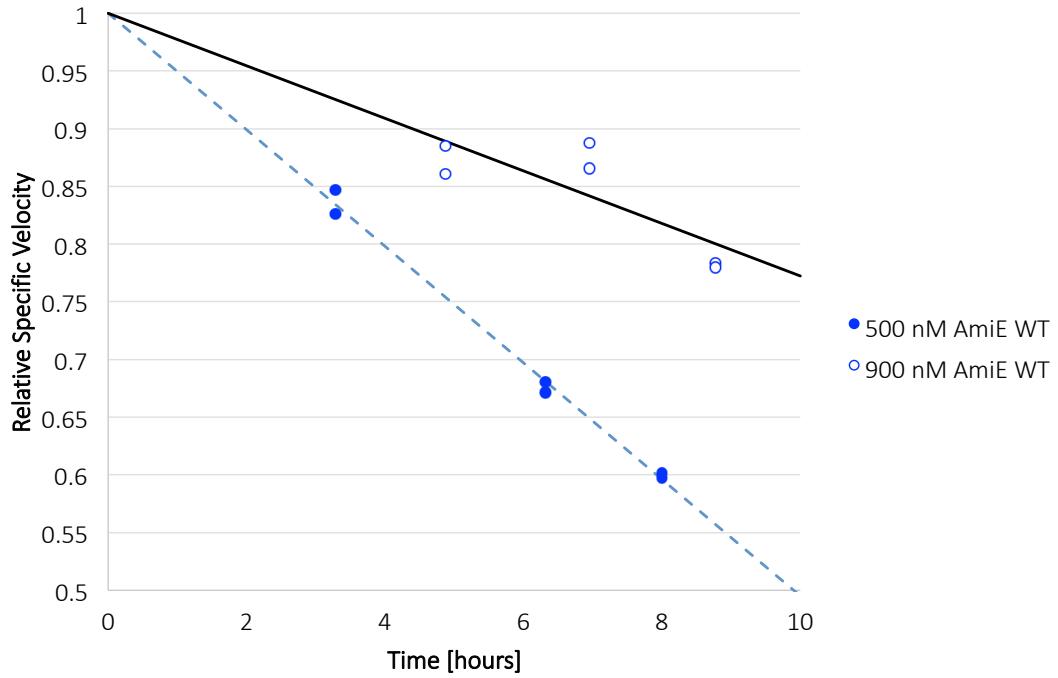
**Figure A 3: Locations of I38V and I122L mutations in AmiE quaternary structure.** A. Three perspectives on the location of the disrupting I38V mutation, the residue I38 is shown as orange spheres. B. Three perspectives on the location of the disrupting I122L mutation, the residue I122 is shown as green spheres. Both mutations are in the monomer core and are not predicted to impact quaternary structure or assembly.



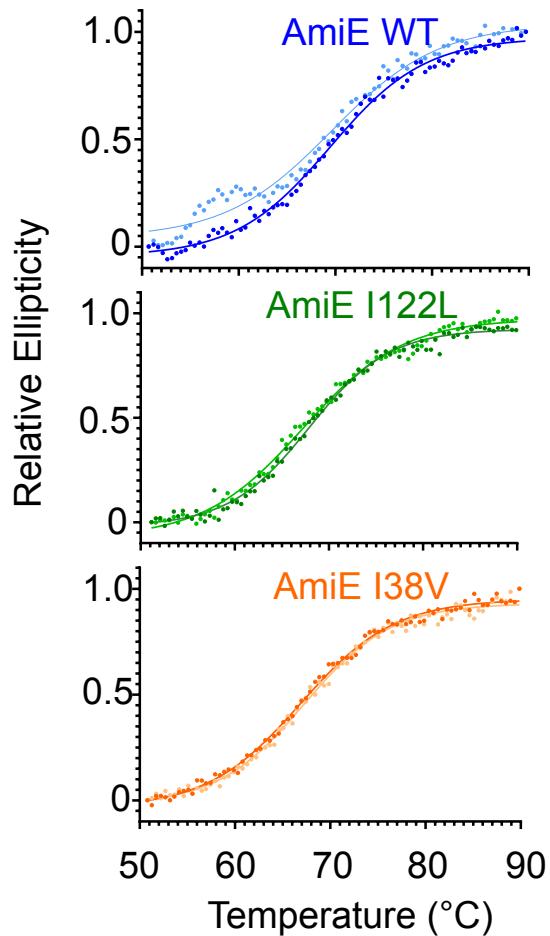
**Figure A 4: AmiE WT tryptophan emission spectra as a function of GDN-HCl concentration.** 52  $\mu$ M Protein was incubated in PBS with the indicated concentration of GDN-HCl at 4oC for 8 hr before fluorescence measurements. The excitation wavelength was 290 nm, while emission was detected at 315-370 nm. The lack of an isosbestic point between spectra indicates more complicated unfolding than a two-state model.



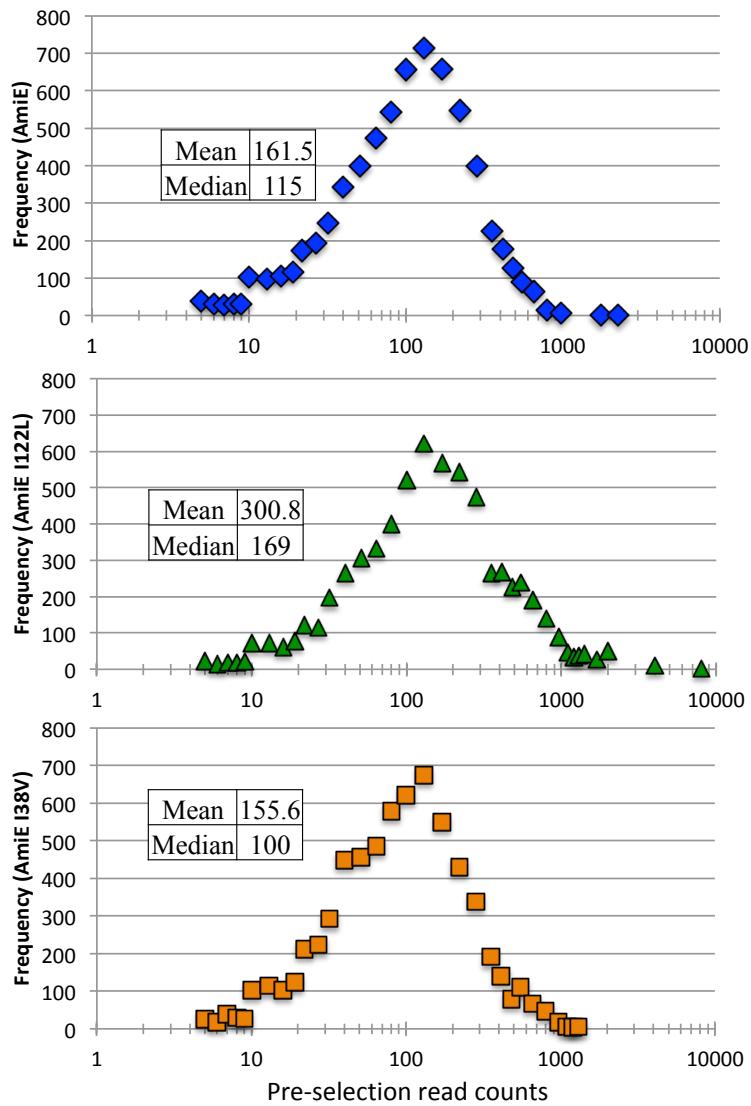
**Figure A 5: Thermal shift analysis of the purified AmiE variants.** A. Melt curves for 10, 5, and 1  $\mu$ M AmiE WT. B. Melt curves for 10, 5, and 1  $\mu$ M AmiE I122L. C. Melt curves for 10, 5, and 1  $\mu$ M AmiE I38V. All experiments were performed with biological (n=2) and technical (n=3) replicates, and representative curves are shown.

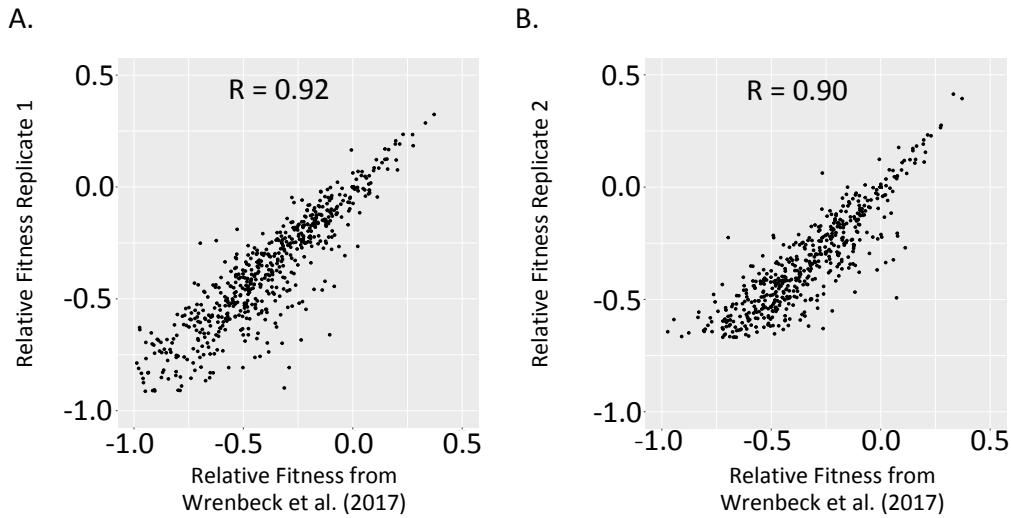


**Figure A 6: Activity loss of AmiE WT following dilution.** Comparison of the relative specific AmiE reaction velocity in a saturating amount of 20 mM acetamide following incubations at dilute concentration in 1x PBS over a time course. Dotted and solid lines represent the best fit linear regressions for the respective data sets (n=6).

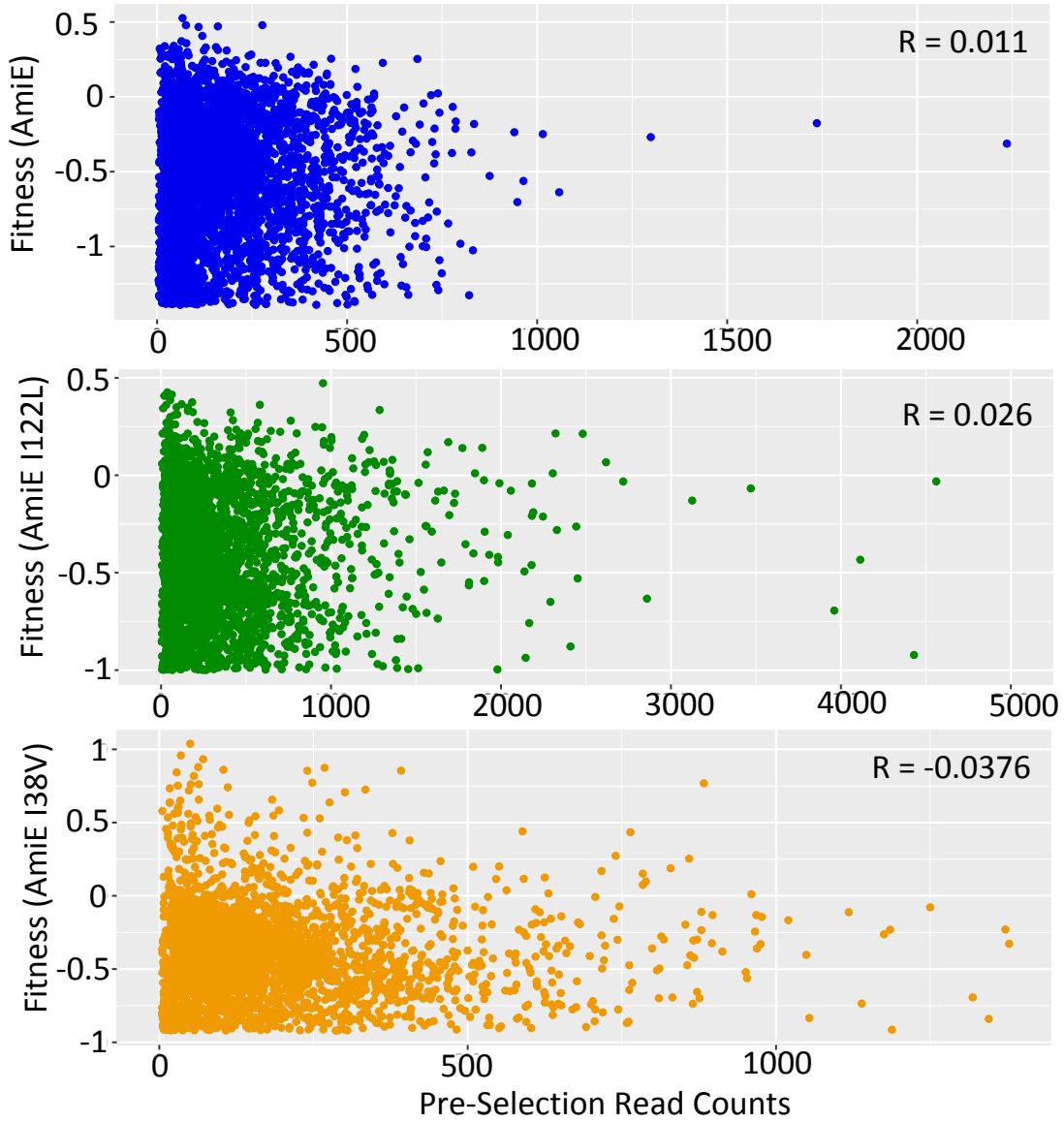


**Figure A 7: Thermal denaturation monitored by far-UV circular dichroism.** Melt curves for AmiE WT (blue), AmiE I122L (green), and AmiE I38V (orange). All experiments were performed by scanning  $\lambda_{222\text{nm}}$  at 1  $\mu\text{M}$  protein concentration with 2 biological replicates (different colors on each panel represent a replicate experiment). Relative ellipticity ranges from completely folded state (0) to unfolded (1).

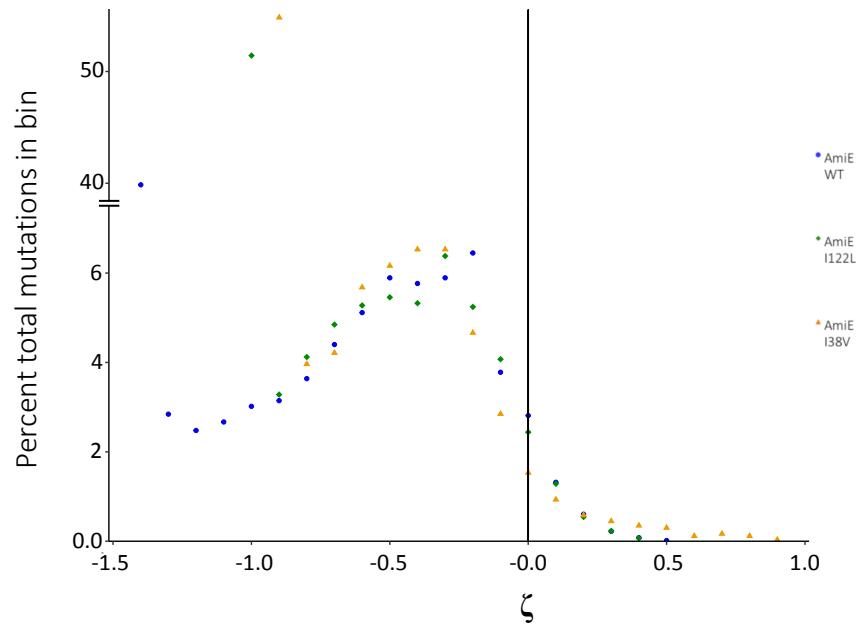




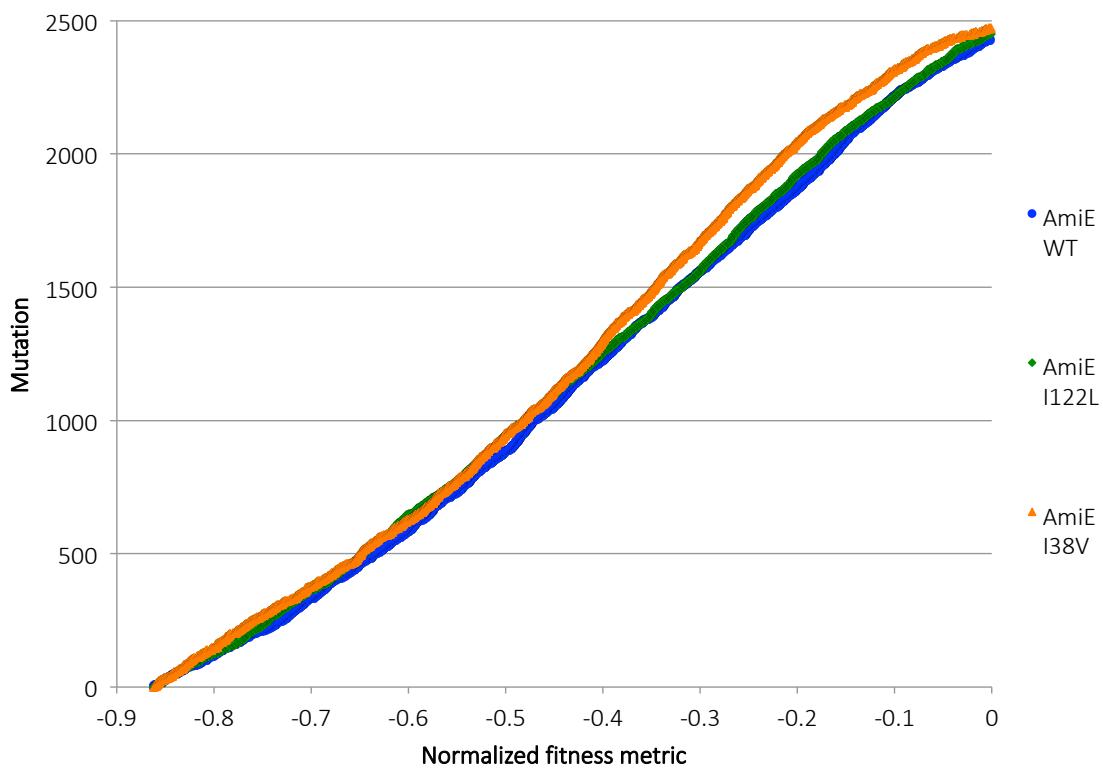
**Figure A 9: Deep mutational scanning replicates for AmiE WT compared with previous literature.** Comparison of relative fitness metrics for a AmiE WT mutational library covering residues 171-255 to those from same mutants reported originally in Wrenbeck et al (2017). These replicate deep mutational scans were performed in parallel to AmiE variant growth selections as internal controls. Pearson's correlation coefficients reported on plots. A. control for AmiE I122L growth selection. B. control for AmiE I38V growth selection.



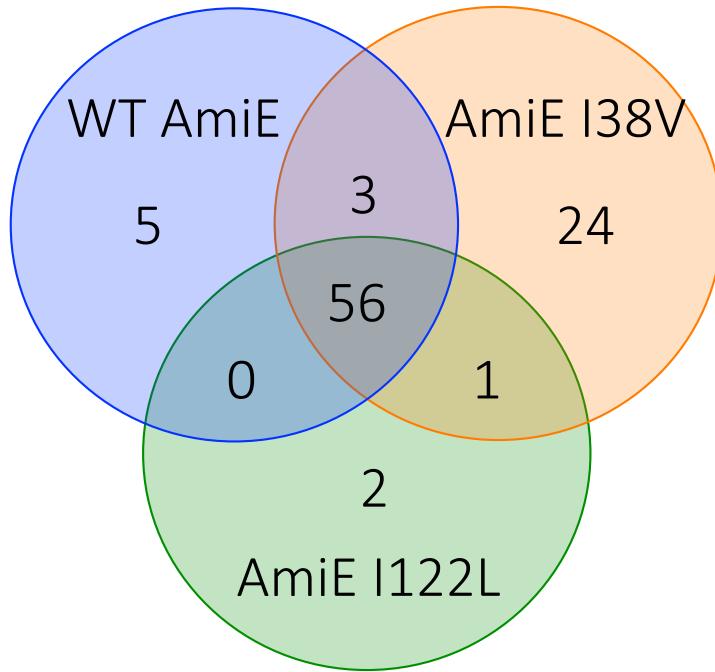
**Figure A 10: Relative fitness metrics for AmiE proteins as a function of pre-selection read counts.** Absolute Pearson's correlation coefficients reported on plots are below 0.04 in all cases, showing that less than 0.2% of the variance can be explained by initial frequency of a given mutant in the library.



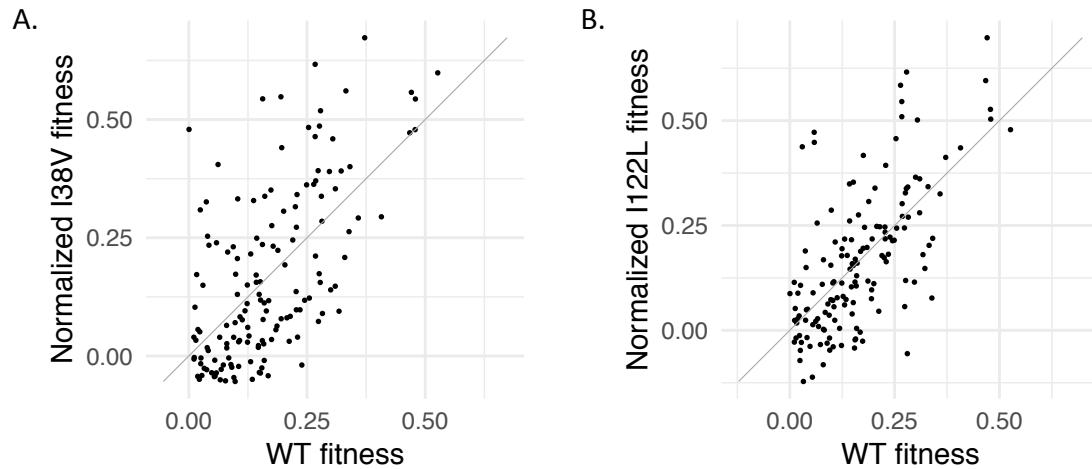
**Figure A 11: Normalized distribution of fitness effects for WT, I122L, and I38V backgrounds.**



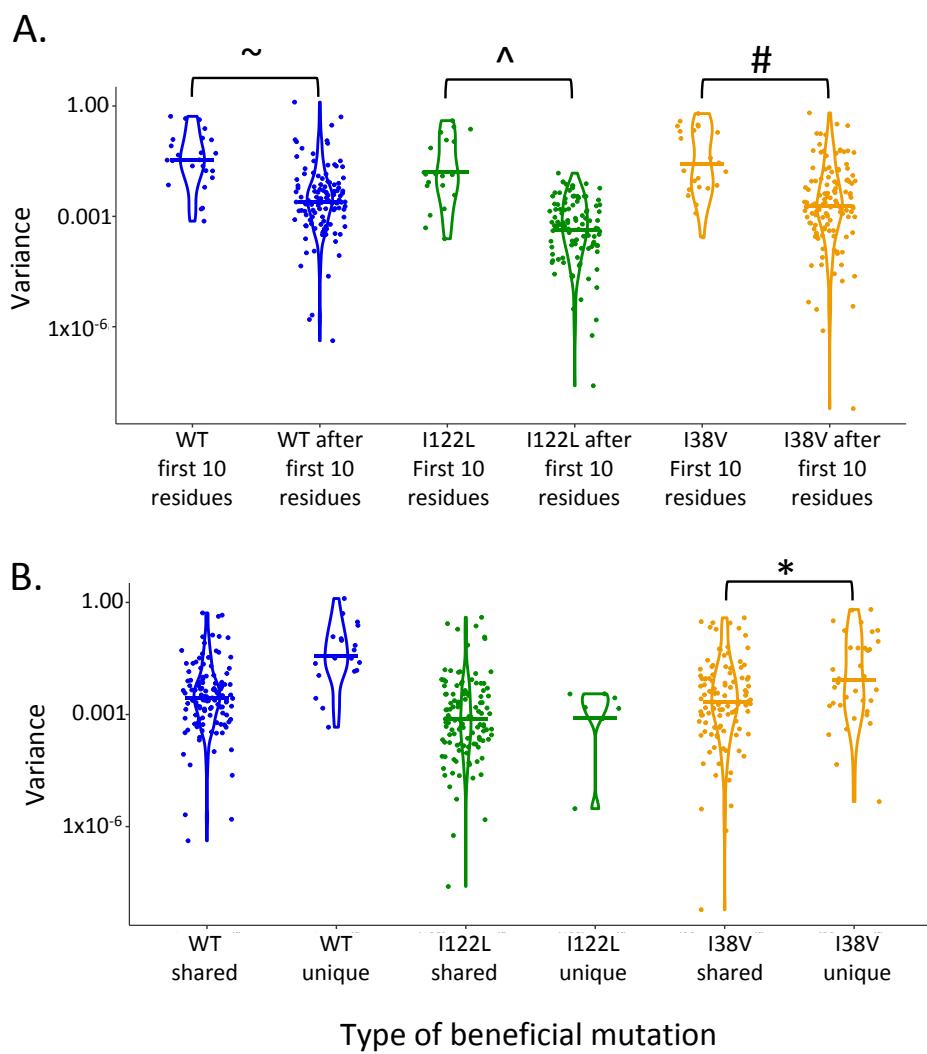
**Figure A 12: Analysis of proportions of deleterious mutations.** Cumulative distribution functions for the deleterious mutations for each enzyme variant are plotted above the lower bounds of experimental measurements.



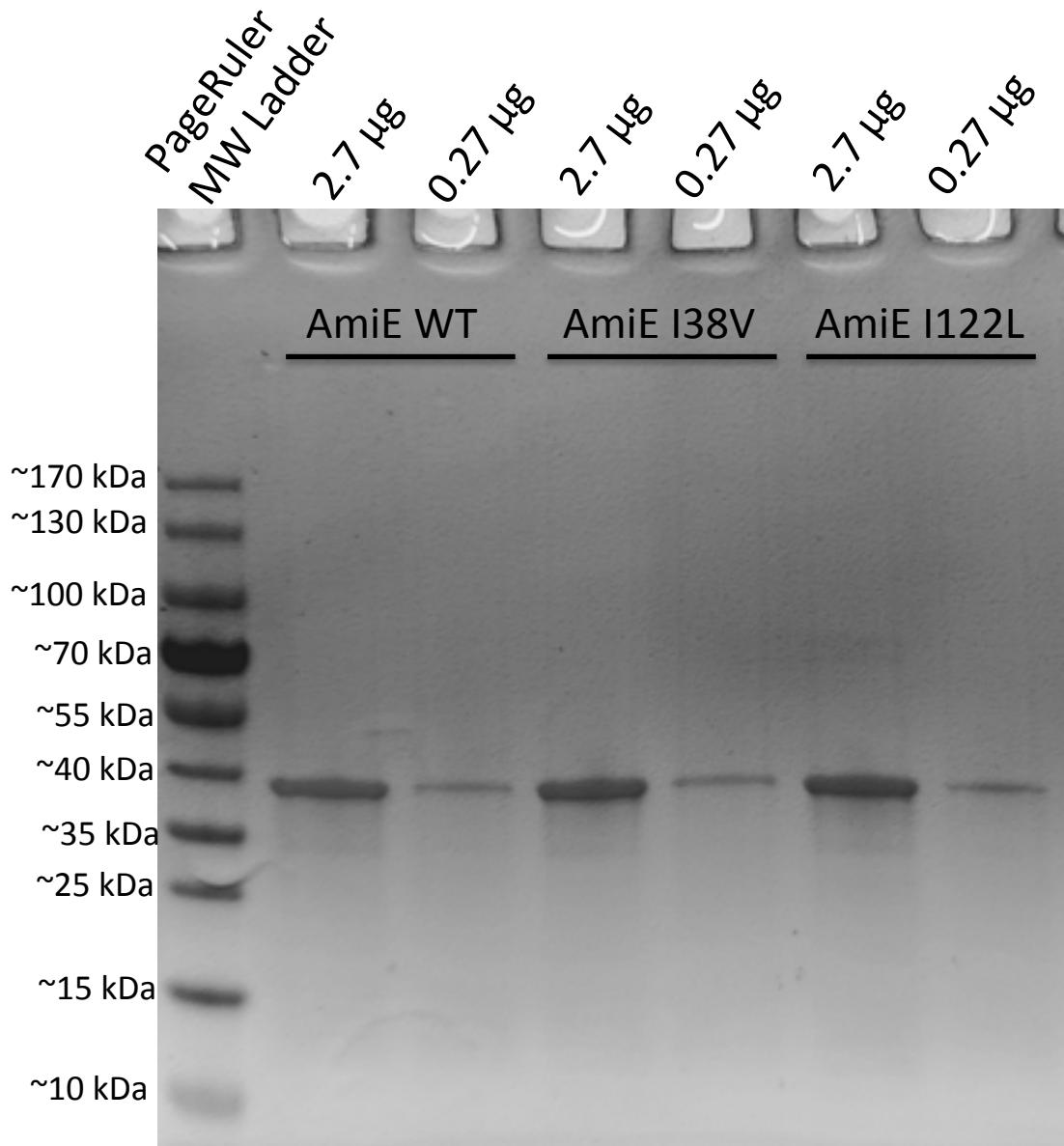
**Figure A 13: Venn diagram of the shared and unique beneficial mutations using the strict cutoff.** Mutations had to improve the growth rate by  $\geq 10\%$  to be classified as a beneficial mutation ( $\zeta_i \geq 0.138$ ).



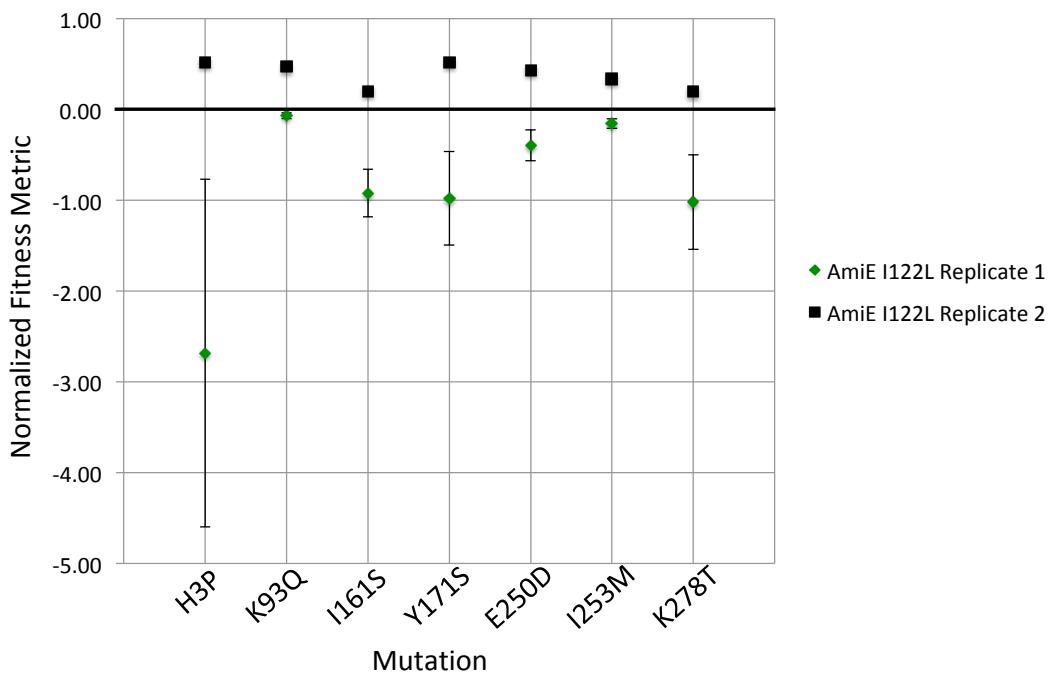
**Figure A 14: Correlation analysis of linear regression normalized shared beneficial mutations.** Comparison of the linear regression normalized fitness metrics for the beneficial mutations shared by all enzymes; solid lines represent the linear function  $Y = 1X + 0$ . A. Comparison of linear regression normalized AmiE I38V fitness metrics with AmiE WT. B. Comparison of linear regression normalized AmiE I122L fitness metrics with AmiE WT.



**Figure A 15: Dot plots of fitness effect synonymous codon variances for beneficial mutations.** A. Dot plots of the fitness effect synonymous codon variances for either the beneficial mutations located in the first 10 residues, or all other beneficial mutations after the first 10 residues (95% confidence interval cutoffs). The distribution of the variances is represented by the violin plot overlay and the colored marker lines represent the mean variance for the population, in all three enzymes the first 10 residues have significantly high variance than those outside of the window:  $\sim$  = WT p-value = 0.02,  $\wedge$  = I122L p-value = 0.01,  $\#$  = I38V p-value = 0.011 from students t-test. B. Dot plots of the fitness effect synonymous codon variances for the shared and unique beneficial mutations for the respective enzymes (95% confidence interval cutoffs) The distribution of the variances is represented by the violin plot overlay and the colored marker lines represent the mean variance for the population. I38V has a significantly high variance in the unique beneficial mutations: \* = I38V p-value = 0.049 from Students t-test.



**Figure A 16: SDS-PAGE analysis of the purity of the purified AmiE variants.** Samples were denatured in SDS-PAGE loading buffer (Laemmli's buffer supplemented to 1.5%  $\beta$ -mercaptoethanol at 1x) at 98° for 10 minutes. Samples run on 4-20% Mini-PROTEAN<sup>®</sup> TGX™ precast gels (Bio-Rad) at 120 V for 1 hour and were washed and stained with SimplyBlue SafeStain (Thermo-Fisher) as described by the manufacturer. Molecular weight ruler used is the PageRuler Prestained Protein Ladder (10-180 kDa, Thermo-Fisher).



**Figure A 17: Comparison of AmiE I122L outlier mutation technical replicates.** Mutations removed from analysis have their fitness metrics for the respective technical replicates shown, error bars represent the 99.977% confidence intervals ( $3.5\sigma$ ) for the fitness metrics. This indicates that errors such as these should occur less than once per dataset for the size of our experiment.

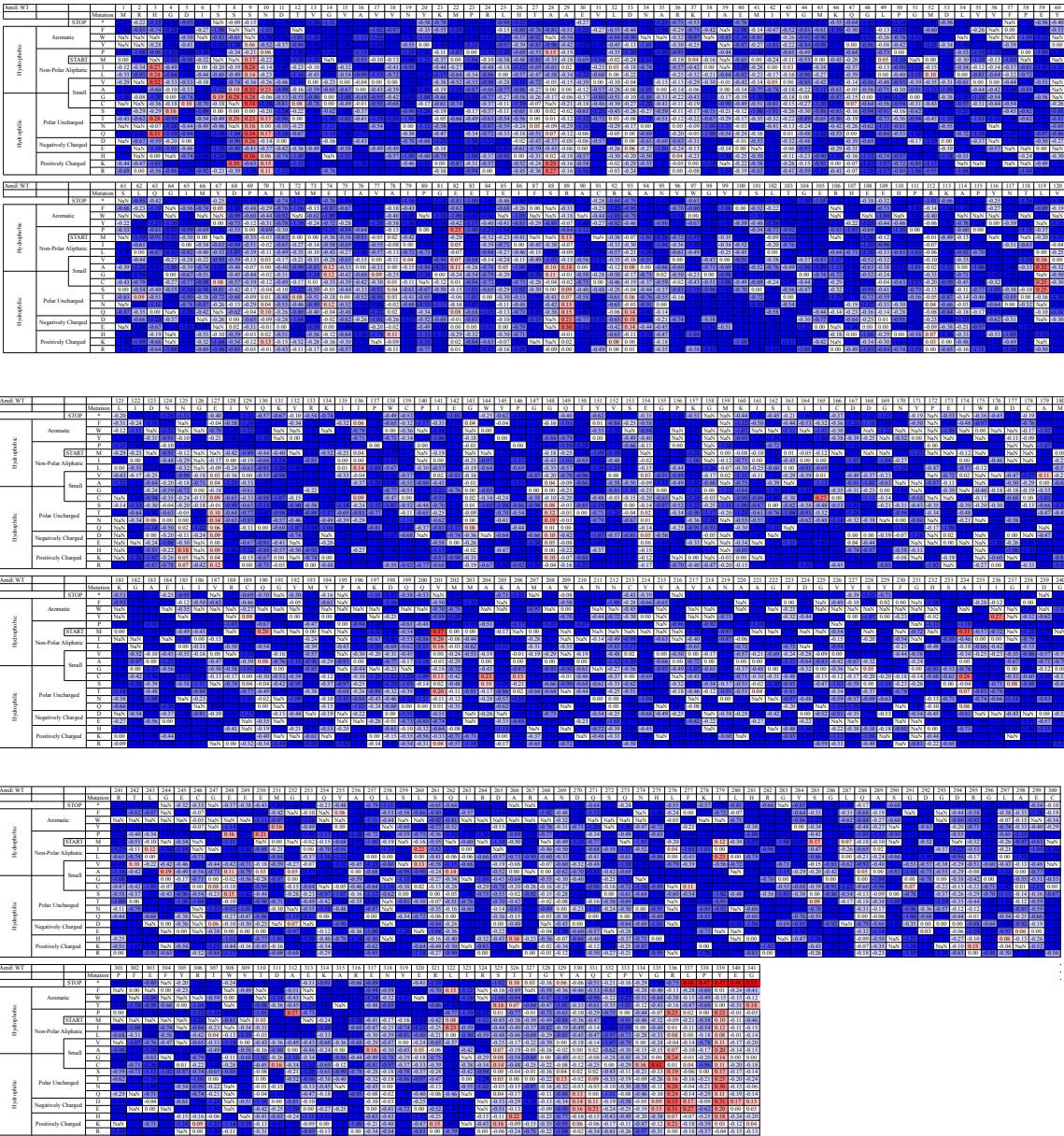


Figure A 18: AmiE WT heatmap.

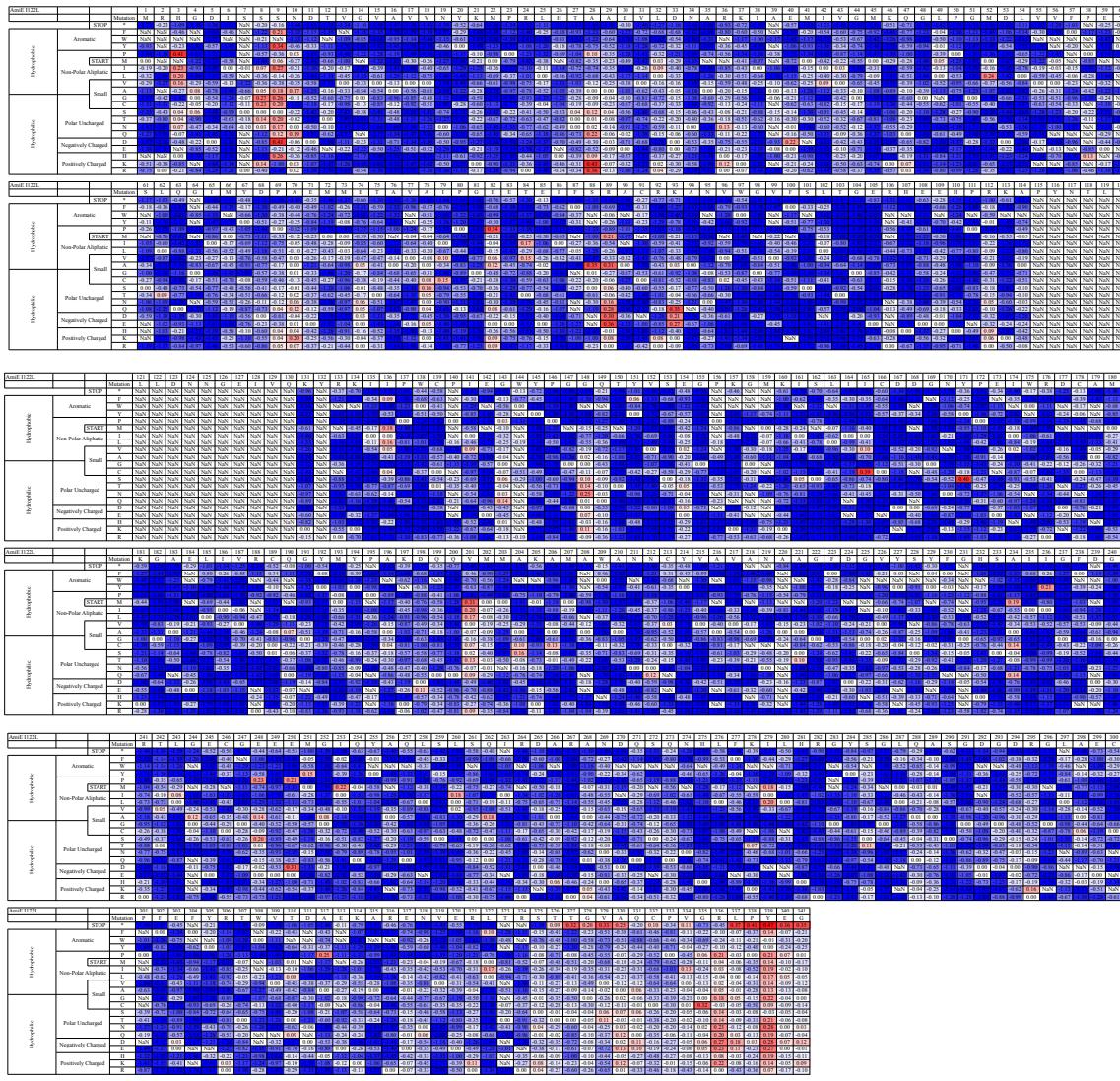


Figure A 19: AmiE I122L heatmap.

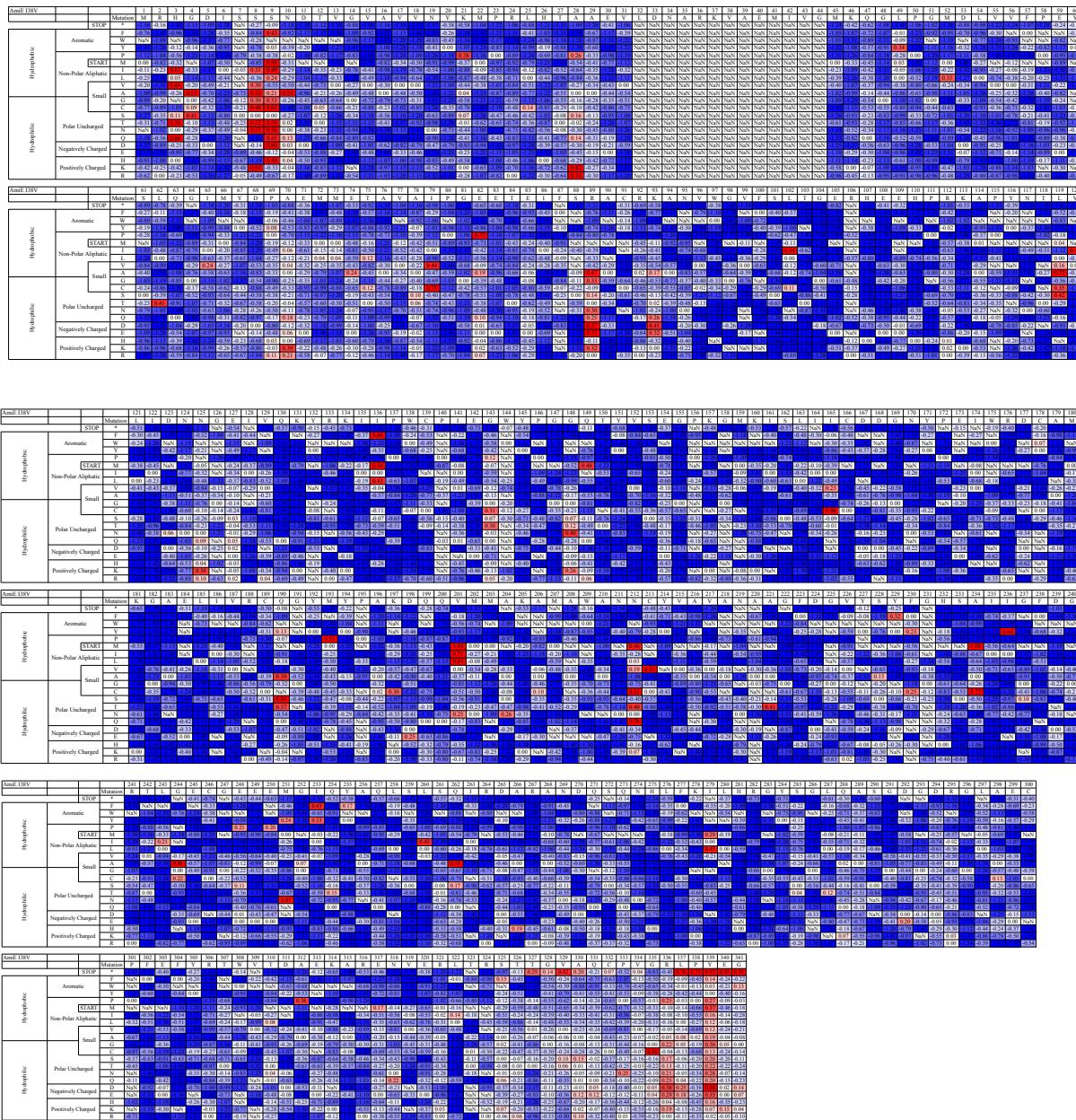


Figure A 20: AmiE I38V heatmap.

**Table A 1: Relative fitness for I38V and I122L synonymous codons in the AmiE WT background.** Green highlight indicates the codons used for the deep mutational scans presented in this work. Data is from Wrenbeck et al.<sup>24</sup>.

	Relative Fitness		Relative Fitness
I38V codon (GTT)	-0.17	I122L codon (TTA)	-0.89
I38V codon (GTG)	-0.43	I122L codon (TTG)	-0.69
I38V codon (GTC)	-1.17	I122L codon (CTT)	-0.46
I38V codon (GTA)	-1.26	I122L codon (CTC)	-0.31
I38V average	-0.35	I122L codon (CTA)	-0.54
		I122L codon (CTG)	-0.69
		I122L average	-0.53

**Table A 2: Biophysical analysis of the AmiE variants.** Error values reported are 1 s.d. except in the kinetic analysis where 95% confidence interval is given.

	AmiE Variant		
	WT	I122L	I38V
Rosetta FilterScan score <sup>^</sup>	0.00	0.40	2.02
Relative $k_{cat}^*$	$1.00 \pm 0.06$	$1.08 \pm 0.06$	$1.02 \pm 0.08$
$k_{cat}$ p-value <sup>#</sup>	1.00	0.57	0.79
Relative $K_M^*$	$1.00 \pm 0.12$	$0.89 \pm 0.15$	$0.82 \pm 0.20$
$K_M$ p-value <sup>#</sup>	1.00	0.64	0.32
Relative protein yield <sup>T</sup>	$1 \pm 0.06$	$0.26 \pm 0.04$	$0.59 \pm 0.08$
Refolded enzyme velocity normalized to folded enzyme (%)	$14.2 \pm 2.47$	$0.114 \pm 0.0494$	$0.064 \pm 0.0250$
Growth Rate in M9 ( $hr^{-1}$ ) - pEDA2 plasmid	$0.78 \pm 0.01$	$0.79 \pm 0.01$	$0.78 \pm 0.02$
Growth Rate in Selection Media ( $hr^{-1}$ ) - pEDA2 plasmid	$0.47 \pm 0.02$	$0.18 \pm 0.01$	$0.11 \pm 0.01$
Growth Rate in M9 ( $hr^{-1}$ ) - pAG plasmid	$0.79 \pm 0.04$	$0.81 \pm 0.03$	$0.79 \pm 0.01$
Growth Rate in Selection Media ( $hr^{-1}$ ) - pAG plasmid	$0.72 \pm 0.02$	$0.40 \pm 0.03$	$0.33 \pm 0.02$

<sup>^</sup> >0 represents a less favorable score to the starting wild-type structure (i.e. decreased likelihood folding)

\* Relative to AmiE WT

<sup>#</sup> Calculated for variant compared to AmiE WT

<sup>T</sup> Expressed via Studiers autoinduction

**Table A 3: Thermal shift analysis data and statistics.** Summary of thermal melting temperatures ( $T_m$ ) determined using Boltzmann curve fit of a 2-state transition. Students t-test p-values are reported for tests comparing AmiE WT with the I38V and I122L variants (N = 6 from two biological replicates of each enzyme;  $\pm$  indicates 1 s.d.).

Variant	Assaying Concentration ( $\mu\text{M}$ )	$T_m$ ( $^{\circ}\text{C}$ )	p-value	Assaying Concentration ( $\mu\text{M}$ )	$T_m$ ( $^{\circ}\text{C}$ )	p-value	Assaying Concentration ( $\mu\text{M}$ )	$T_m$ ( $^{\circ}\text{C}$ )	p-value
AmiE WT	1	N/A	-	5	$69.6 \pm 0.3$	-	10	$68.8 \pm 0.2$	-
AmiE I122L	1	N/A	N/A	5	$69.2 \pm 1.3$	0.42	10	$68.7 \pm 1$	0.74
AmiE I38V	1	N/A	N/A	5	$69.3 \pm 0.4$	0.2	10	$68.5 \pm 0.4$	0.16

**Table A 4: Circular dichroism analysis of thermal denaturation statistics.** Summary of apparent thermal melting temperatures ( $T_m$ ) determined using Boltzmann curve fit of a 2-state folded-unfolded transition. Student's t-test p-values are reported for tests comparing AmiE WT with the I38V and I122L variants (N = 2 from respective biological replicates of each enzyme).

Variant	Average $T_m$ (°C)	Replicate 1	Replicate 2	p-value relative to AmiE WT
AmiE WT	69.5	69.5	69.6	-
AmiE I122L	67.5	67.9	67.2	0.03
AmiE I38V	67.1	66.9	67.4	0.01

**Table A 5: DNA sequences of transcriptional elements in plasmids used for deep mutational scans.** Differences between sequences are shown in red, the RBS in the 5' UTR is underlined.

	-35 Promoter	-10 Promoter	5' UTR
pEDA2 – AmiE	TT <u>CCC</u> CG	TA <u>A</u> TAT	TCAGGGAGACCACAA <u>AC</u> GG
pAG – AmiE I38V	TT <u>GG</u> CG	TAC <u>CA</u> AT	TTTCCCTCTACAAATAATT TTGTTAAC <u>TT</u> C <u>TA</u> GAAA
pAG – AmiE I122L	TT <u>GG</u> CG	TAC <u>CA</u> AT	TAATTT <u>GTT</u> AA <u>CT</u> T TAAGA <u>AG</u> <u>TTT</u> TATACAT

**Table A 6: Summary of MG1655 rph+ transformants obtained during selection strain mutant library preparation.**  $5 \times 10^4$  transformants corresponds to >99% theoretical library coverage.

Selection Strain SSM Library Preparation	Transformants Obtained
AmiE I122L pAG Tile 1	$15 \times 10^5$
AmiE I122L pAG Tile 2	$19 \times 10^5$
AmiE I122L pAG Tile 3	$7.5 \times 10^5$
AmiE I122L pAG Tile 4	$23 \times 10^5$
AmiE I38V pAG Tile 1	$6.9 \times 10^5$
AmiE I38V pAG Tile 2	$5.8 \times 10^5$
AmiE I38V pAG Tile 3	$7.5 \times 10^5$
AmiE I38V pAG Tile 4	$7.8 \times 10^5$
WT amiE pEDA2 Tile 3	$13 \times 10^5$

**Table A 7: Summary of mutational library statistics.**

Screen	Growth				
Enzyme	AmiE WT				
Tile Number	1	2	3	4	Cumulative
Residues	1-85	86-170	171-255	256-341	1-341
Population	Growth Selected	Growth Selected	Growth Selected	Growth Selected	Growth Selected
Number of mutated codons	85	85	85	86	341
Reference sequencing reads post quality filter	448689	393052	561875	511337	1914953
Selected sequencing reads post quality filter	1454211	3016045	2101231	1691121	8262608
<b>Fold oversampling of codon combinations</b>					
Reference total	81	71.6	101.8	91.6	86.5
Selected total	256.6	544	378.1	292.1	367.7
Reference non-synonymous	50.8	38.3	54.5	55.8	49.9
Selected non-synonymous	79.8	95.7	57.4	95	82.0
<b>Percent of reads with:</b>					
No nonsynonymous mutations	37.6	46.9	46.7	39.4	42.7
One nonsynonymous mutation	60.6	52.2	51.9	59.1	56.0
Multiple nonsynonymous mutation	1.8	0.9	1.4	1.4	1.4
Coverage of possible single nonsynonymous mutations	95.2	93.9	89	94.7	93.2
Screen	Growth				
Enzyme	AmiE I122L				
Tile Number	1	2	3	4	Cumulative
Residues	1-85	86-114, 131, 133-170	171-255	256-341	1-114,131,133-341
Population	Growth Selected	Growth Selected	Growth Selected	Growth Selected	Growth Selected
Number of mutated codons	85	68	85	86	324.0
Reference sequencing reads post quality filter	674718	593530	897746	1214190	3380184
Selected sequencing reads post quality filter	1286349	1002404	1115963	931995	4336711
<b>Fold oversampling of codon combinations</b>					
Reference total	118.7	99	146.8	210	143.6
Selected total	221.8	176.5	198.1	156	188.1
Reference non-synonymous	76.1	61.5	90.1	137	91.2
Selected non-synonymous	71.6	35.8	37.8	63.6	52.2
<b>Reference population percent of reads with:</b>					
No nonsynonymous mutations	35.4	35.3	35.2	34	35.0
One nonsynonymous mutation	60.4	55.4	53.7	61.2	57.7
Multiple nonsynonymous mutation	4.3	9.3	11.1	4.8	7.4
Coverage of possible single nonsynonymous mutations	94.5	92.5	93	94.8	93.7
Screen	Growth				
Enzyme	AmiE I38V				
Tile Number	1	2	3	4	Cumulative
Residues	1-31, 45-85	86-170	171-255	256-341	1-31, 45-341
Population	Growth Selected	Growth Selected	Growth Selected	Growth Selected	Growth Selected
Number of mutated codons	72	85	85	86	328
Reference sequencing reads post quality filter	550823	500238	525635	490442	2067138
Selected sequencing reads post quality filter	1373021	1060307	991922	1129295	4554545
<b>Fold oversampling of codon combinations</b>					
Reference total	100.7	91.5	96.1	88.6	94.2
Selected total	236.5	189	173.1	185	195.9
Reference non-synonymous	50.3	45.7	48.3	46	47.6
Selected non-synonymous	96.4	42.6	69.8	81	72.5
<b>Percent of reads with:</b>					
No nonsynonymous mutations	50.5	50.5	50.3	48.6	50.0
One nonsynonymous mutation	48.9	48.9	49.2	50.8	49.5
Multiple nonsynonymous mutation	0.6	0.5	0.6	0.6	0.6
Coverage of possible single nonsynonymous mutations	96.7	88.6	87.9	93.9	91.8

**Table A 8: Goodness-of-fit test statistics for distribution fittings.** For all enzymes the determined p-values indicate a failure to reject the null hypotheses that they can be fit by a general Pareto distribution with a negative shape parameter. All calculations were performed as in Wrenbeck et al.<sup>24</sup>.

<u>General Pareto distribution</u>	AmiE WT	AmiE I122L	AmiE I38V
bootstrap test p-value	0.16	0.29	0.10
shape parameter	-0.28	-0.33	-0.32

**Table A 9: Deleterious empirical cumulative distribution function (ECDF) analysis statistics.** A table of the two-sample Kolmogorov-Smirnov test results for comparing the deleterious mutation ECDFs for the three enzymes.

Pairing (x-y)	K-S Test	D	p-value	Null hypothesis	Alternative hypothesis	AmiE WT (N)	AmiE variant (N)
AmiE WT – AmiE I38V	two sided	0.058	0.0005	is equal to	not equal	2428	2472
AmiE WT – AmiE I38V	greater	0.002	0.99	x not greater than y	the CDF of x lies above that of y	2428	2472
AmiE WT – AmiE I38V	less	0.058	0.0002	x not less than y	the CDF of x lies below that of y	2428	2472
AmiE WT – AmiE I122L	two sided	0.024	0.48	is equal to	not equal	2428	2450
AmiE WT – AmiE I122L	greater	0.010	0.77	x not greater than y	the CDF of x lies above that of y	2428	2450
AmiE WT – AmiE I122L	less	0.024	0.24	x not less than y	the CDF of x lies below that of y	2428	2450

**Table A 10: Inner and outer primers for PCR reactions for Illumina sequencing.** Red indicates overhang regions for attaching Illumina adapter primers (inner PCR primers) or overhangs for attaching to inner PCR product (outer PCR primers), black is the overlap region in the gene or the barcode, blue is the Illumina adapter.

Inner PCR primers	Sequence (5' to 3')
Fwd_Tile_1	gttcagagttctacagtccgacgatcttaactttaagaagttttatacat
Fwd_Tile_2	gttcagagttctacagtccgacgatccgcagaaaacggaa
Fwd_Tile_3	gttcagagttctacagtccgacgatctcgatgacggtaat
Fwd_Tile_4	gttcagagttctacagtccgacgatcaagaatggcattcaatac
Rev_Tile_1	ccttggcacccgagaattcca aagcacggctaaagat
Rev_Tile_2	ccttggcacccgagaattcca ctctccaaattccggata
Rev_Tile_3	ccttggcacccgagaattcca cagagacaactgcgc
Rev_Tile_4	ccttggcacccgagaattcca tggtggtctcgag
<b>Illumina outer primer adapter</b>	aatgatacgccgaccaccgagatctacac
<b>Illumina outer PCR adapters and barcodes</b>	gttcagagttctacagtccga
RPI37 (AmiE WT unselected, Tile 1)	caagcagaagacggcatacgagat ATTCCG gtgactggagttccctggcaccc gagaattcca
RPI22 (AmiE WT unselected, Tile 2)	caagcagaagacggcatacgagat CGTACG gtgactggagttccctggcaccc gagaattcca
RPI39 (AmiE WT unselected, Tile 3)	caagcagaagacggcatacgagat GTATAG gtgactggagttccctggcaccc gagaattcca
RPI40 (AmiE WT unselected, Tile 4)	caagcagaagacggcatacgagat TCTGAG gtgactggagttccctggcaccc gagaattcca
RPI41 (AmiE WT selected, Tile 1 replicate 1)	caagcagaagacggcatacgagat GTCGTC gtgactggagttccctggcaccc gagaattcca
RPI38 (AmiE WT selected, Tile 1 replicate 2)	caagcagaagacggcatacgagat AGCTAG gtgactggagttccctggcaccc gagaattcca
RPI33 (AmiE WT selected, Tile 2 replicate 1)	caagcagaagacggcatacgagat CGCCTG gtgactggagttccctggcaccc gagaattcca
RPI34 (AmiE WT selected, Tile 2 replicate 2)	caagcagaagacggcatacgagat GCCATG gtgactggagttccctggcaccc gagaattcca
RPI43 (AmiE WT selected, Tile 3 replicate 1)	caagcagaagacggcatacgagat GCTGTA gtgactggagttccctggcaccc gagaattcca
RPI40 (AmiE WT selected, Tile 3 replicate 2)	caagcagaagacggcatacgagat TCTGAG gtgactggagttccctggcaccc gagaattcca
RPI44 (AmiE WT selected, Tile 4 replicate 1)	caagcagaagacggcatacgagat ATTATA gtgactggagttccctggcaccc gagaattcca
RPI41 (AmiE WT selected, Tile 4 replicate 2)	caagcagaagacggcatacgagat GTCGTC gtgactggagttccctggcaccc gagaattcca
RPI20 (AmiE I38V unselected, Tile 1)	caagcagaagacggcatacgagat GGCCAC gtgactggagttccctggcaccc gagaattcca
RPI21 (AmiE I38V unselected, Tile 2)	caagcagaagacggcatacgagat CGAACAC gtgactggagttccctggcaccc gagaattcca
RPI27 (AmiE I38V unselected, Tile 3)	caagcagaagacggcatacgagat AGGAAT gtgactggagttccctggcaccc gagaattcca

**Table A 10 (cont'd)**

RPI34 (AmiE I38V unselected, Tile 4)	caagcagaagacggcatacgagatGCCATGtgactggagttccctggcaccc gagaattcca
RPI37 (AmiE I38V selected, Tile 1 replicate 1)	caagcagaagacggcatacgagatATTCCGtgactggagttccctggcaccc gagaattcca
RPI38 (AmiE I38V selected, Tile 1 replicate 2)	caagcagaagacggcatacgagatAGCTAGtgactggagttccctggcaccc gagaattcca
RPI16 (AmiE I38V selected, Tile 2 replicate 1)	caagcagaagacggcatacgagatGGACGGtgactggagttccctggcaccc gagaattcca
RPI14 (AmiE I38V selected, Tile 2 replicate 2)	caagcagaagacggcatacgagatGGAACCTtgactggagttccctggcaccc gagaattcca
RPI48 (AmiE I38V selected, Tile 3 replicate 1)	caagcagaagacggcatacgagatTGCCGAtgactggagttccctggcaccc gagaattcca
RPI23 (AmiE I38V selected, Tile 3 replicate 2 )	caagcagaagacggcatacgagatCCACTCtgactggagttccctggcaccc gagaattcca
RPI32 (AmiE I38V selected, Tile 4 replicate 1)	caagcagaagacggcatacgagatTGAGTGtgactggagttccctggcaccc gagaattcca
RPI44 (AmiE I38V selected, Tile 4 replicate 2)	caagcagaagacggcatacgagatATTATAtgactggagttccctggcaccc gagaattcca
RPI19 (AmiE WT control, AmiE I38V selection, unselected, Tile 3)	caagcagaagacggcatacgagatTTTCACtgactggagttccctggcaccc gagaattcca
RPI26 (AmiE WT control, AmiE I38V selection, selected, Tile 3)	caagcagaagacggcatacgagatGCTCATtgactggagttccctggcaccc gagaattcca
RPI33 (AmiE I122L unselected, Tile 1, sequencing runs 1 and 2)	caagcagaagacggcatacgagatCGCCTGtgactggagttccctggcaccc gagaattcca
RPI46 (AmiE I122L unselected, Tile 2, sequencing runs 1 and 2)	caagcagaagacggcatacgagatTCGGGAtgactggagttccctggcaccc gagaattcca
RPI40 (AmiE I122L unselected, Tile 3, sequencing runs 1and 2)	caagcagaagacggcatacgagatTCTGAGtgactggagttccctggcaccc gagaattcca
RPI29 (AmiE I122L unselected, Tile 4, sequencing runs 1 and 2)	caagcagaagacggcatacgagatTAGTTGtgactggagttccctggcaccc gagaattcca
RPI17 (AmiE I122L selected, Tile 1 replicate 1)	caagcagaagacggcatacgagatCTCTACtgactggagttccctggcaccc gagaattcca
RPI32 (AmiE I122L selected, Tile 1 replicate 2)	caagcagaagacggcatacgagatTGAGTGtgactggagttccctggcaccc gagaattcca
RPI12 (AmiE I122L selected, Tile 2 replicate 1)	caagcagaagacggcatacgagatTACAAGtgactggagttccctggcaccc gagaattcca
RPI42 (AmiE I122L selected, Tile 2 replicate 2)	caagcagaagacggcatacgagatCGATTAtgactggagttccctggcaccc gagaattcca
RPI25 (AmiE I122L selected, Tile 3 replicate 1)	caagcagaagacggcatacgagatATCAGTtgactggagttccctggcaccc gagaattcca
RPI46 (AmiE I122L selected, Tile 3 replicate 2 )	caagcagaagacggcatacgagatTCGGGAtgactggagttccctggcaccc gagaattcca
RPI31 (AmiE I122L selected, Tile 4 replicate 1)	caagcagaagacggcatacgagatATCGTGtgactggagttccctggcaccc gagaattcca
RPI15 (AmiE I122L selected, Tile 4 replicate 2)	caagcagaagacggcatacgagatTGACATtgactggagttccctggcaccc gagaattcca
RPI7 (AmiE WT control, AmiE I122L selection, unselected, Tile 3)	caagcagaagacggcatacgagatGATCTGtgactggagttccctggcaccc gagaattcca
RPI10 (AmiE WT control, AmiE I122L selection, selected, Tile 3)	caagcagaagacggcatacgagatAAGCTAtgactggagttccctggcaccc gagaattcca

**Table A 11: Beneficial mutations shared by all enzymes.**

Location-mutation	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE I122L cumulative normalized fitness metric	AmiE I122L technical replicate 1 normalized fitness metric	AmiE I122L technical replicate 2 normalized fitness metric	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric
9T	0.25	0.24	0.26	0.20	0.19	0.20	0.59	0.60	0.59
9Q	0.34	0.34	0.34	0.12	0.11	0.12	0.45	0.46	0.44
9N	0.16	0.21	0.11	0.17	0.14	0.18	0.56	0.56	0.56
9I	0.28	0.26	0.29	0.27	0.25	0.28	0.49	0.49	0.48
9H	0.36	0.36	0.36	0.26	0.24	0.28	0.49	0.49	0.50
9G	0.28	0.25	0.29	0.26	0.25	0.27	0.33	0.31	0.34
9D	0.26	0.28	0.25	0.41	0.38	0.43	0.60	0.59	0.60
9C	0.34	0.34	0.34	0.20	0.19	0.21	0.65	0.65	0.65
9A	0.32	0.31	0.32	0.18	0.19	0.16	0.21	0.22	0.21
93E	0.14	0.14	0.15	0.27	0.26	0.29	0.32	0.32	0.32
93D	0.18	0.17	0.19	0.21	0.20	0.23	0.41	0.40	0.42
93A	0.08	0.07	0.09	0.03	0.03	0.02	0.17	0.16	0.18
8T	0.20	0.20	0.20	0.14	0.15	0.13	0.52	0.52	0.51
8K	0.30	0.30	0.30	0.14	0.14	0.14	0.64	0.63	0.64
8G	0.28	0.28	0.28	0.27	0.26	0.27	0.30	0.31	0.29
8A	0.10	0.12	0.09	0.05	0.04	0.06	0.55	0.55	0.55
89S	0.09	0.07	0.11	0.06	0.08	0.03	0.14	0.12	0.15
89Q	0.15	0.14	0.16	0.28	0.28	0.27	0.25	0.25	0.25
89N	0.15	0.14	0.16	0.16	0.16	0.17	0.30	0.29	0.31
89E	0.30	0.29	0.31	0.36	0.36	0.36	0.73	0.73	0.74
89D	0.23	0.21	0.24	0.30	0.30	0.30	0.57	0.57	0.55
89A	0.18	0.17	0.18	0.31	0.32	0.30	0.47	0.46	0.49
82Q	0.08	0.07	0.09	0.08	0.08	0.07	0.10	0.09	0.11
82P	0.25	0.25	0.26	0.34	0.33	0.34	0.77	0.77	0.77
82A	0.11	0.10	0.12	0.12	0.13	0.11	0.19	0.18	0.19
79V	0.04	0.05	0.03	0.10	0.11	0.09	0.44	0.44	0.44
78S	0.04	0.02	0.05	0.16	0.16	0.16	0.10	0.11	0.10
74A	0.12	0.13	0.12	0.05	0.05	0.05	0.24	0.22	0.25
70Q	0.10	0.08	0.11	0.12	0.11	0.12	0.18	0.17	0.19
70K	0.13	0.13	0.13	0.20	0.20	0.19	0.39	0.38	0.39
62T	0.09	0.10	0.09	0.09	0.11	0.08	0.41	0.41	0.41
52L	0.10	0.09	0.11	0.24	0.23	0.24	0.33	0.33	0.33
4S	0.16	0.14	0.17	0.06	0.08	0.04	0.41	0.42	0.41
3V	0.32	0.32	0.33	0.16	0.16	0.15	0.64	0.63	0.64
3L	0.24	0.24	0.24	0.20	0.20	0.20	0.05	0.03	0.07
3I	0.27	0.25	0.28	0.23	0.21	0.24	0.61	0.61	0.60
341D	0.13	0.11	0.14	0.12	0.10	0.13	0.14	0.13	0.14
339V	0.11	0.10	0.11	0.14	0.13	0.15	0.12	0.11	0.13
339T	0.23	0.22	0.23	0.21	0.19	0.22	0.22	0.22	0.21

**Table A 11 (cont'd)**

Location-mutation	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE I122L cumulative normalized fitness metric	AmiE I122L technical replicate 1 normalized fitness metric	AmiE I122L technical replicate 2 normalized fitness metric	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric
339Q	0.11	0.10	0.11	0.19	0.19	0.20	0.20	0.21	0.18
339P	0.23	0.22	0.23	0.21	0.21	0.22	0.27	0.28	0.27
339N	0.30	0.28	0.31	0.28	0.26	0.30	0.28	0.29	0.27
339M	0.10	0.09	0.11	0.14	0.13	0.14	0.37	0.36	0.38
339L	0.08	0.08	0.08	0.17	0.17	0.17	0.12	0.11	0.13
339I	0.12	0.13	0.12	0.19	0.18	0.19	0.16	0.15	0.18
339H	0.18	0.16	0.20	0.19	0.17	0.20	0.16	0.15	0.16
339G	0.14	0.12	0.15	0.22	0.21	0.23	0.30	0.30	0.30
339E	0.20	0.18	0.21	0.27	0.27	0.27	0.35	0.36	0.35
339D	0.31	0.29	0.32	0.28	0.28	0.28	0.58	0.58	0.59
339C	0.11	0.08	0.12	0.09	0.08	0.09	0.13	0.13	0.12
339A	0.20	0.19	0.20	0.13	0.12	0.13	0.19	0.19	0.19
337E	0.27	0.26	0.28	0.11	0.09	0.12	0.18	0.20	0.16
337D	0.17	0.16	0.17	0.18	0.18	0.18	0.25	0.25	0.24
336T	0.16	0.15	0.16	0.14	0.13	0.15	0.13	0.12	0.13
336S	0.19	0.17	0.19	0.14	0.14	0.14	0.17	0.17	0.16
336Q	0.24	0.25	0.24	0.20	0.19	0.20	0.25	0.24	0.26
336P	0.25	0.26	0.25	0.21	0.20	0.22	0.25	0.26	0.25
336N	0.28	0.27	0.29	0.23	0.23	0.23	0.21	0.20	0.21
336K	0.21	0.18	0.22	0.22	0.21	0.22	0.19	0.17	0.22
336H	0.07	0.06	0.07	0.08	0.07	0.09	0.04	0.03	0.04
336G	0.24	0.24	0.23	0.18	0.16	0.19	0.22	0.21	0.23
336E	0.31	0.30	0.31	0.23	0.23	0.24	0.29	0.27	0.31
336D	0.33	0.33	0.33	0.27	0.26	0.28	0.38	0.37	0.38
336A	0.07	0.07	0.07	0.05	0.05	0.06	0.05	0.06	0.05
335D	0.09	0.09	0.09	0.06	0.06	0.05	0.05	0.05	0.05
335C	0.41	0.38	0.42	0.32	0.35	0.30	0.50	0.50	0.50
331E	0.21	0.18	0.23	0.10	0.10	0.10	0.12	0.13	0.12
330S	0.02	0.03	0.02	0.07	0.07	0.06	0.15	0.14	0.16
330E	0.16	0.18	0.14	0.13	0.11	0.16	0.12	0.09	0.14
329T	0.13	0.13	0.13	0.11	0.13	0.09	0.06	0.06	0.06
329S	0.04	0.05	0.04	0.06	0.06	0.07	0.10	0.09	0.10
322I	0.23	0.21	0.24	0.17	0.16	0.18	0.14	0.14	0.13
312P	0.27	0.25	0.28	0.25	0.27	0.23	0.38	0.38	0.38
28R	0.27	0.27	0.27	0.36	0.36	0.37	0.74	0.74	0.74
28P	0.15	0.16	0.14	0.10	0.10	0.09	0.26	0.25	0.27
28K	0.28	0.26	0.29	0.43	0.42	0.43	0.82	0.81	0.82
279M	0.12	0.10	0.14	0.18	0.15	0.19	0.29	0.30	0.28
279L	0.23	0.21	0.24	0.20	0.20	0.20	0.47	0.46	0.47

**Table A 11 (cont'd)**

Location-mutation	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE I122L cumulative normalized fitness metric	AmiE I122L technical replicate 1 normalized fitness metric	AmiE I122L technical replicate 2 normalized fitness metric	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric
266H	0.16	0.14	0.17	0.06	0.07	0.05	0.19	0.20	0.18
262A	0.14	0.11	0.15	0.18	0.17	0.18	0.55	0.54	0.55
260I	0.22	0.19	0.23	0.18	0.19	0.16	0.43	0.43	0.43
252A	0.05	0.02	0.07	0.08	0.08	0.09	0.07	0.07	0.06
251Y	0.16	0.20	0.14	0.15	0.15	0.15	0.24	0.20	0.27
250P	0.21	0.21	0.22	0.21	0.21	0.22	0.20	0.20	0.20
248S	0.15	0.14	0.15	0.20	0.20	0.19	0.11	0.11	0.11
248P	0.16	0.17	0.16	0.23	0.23	0.23	0.21	0.22	0.21
244A	0.19	0.21	0.19	0.12	0.12	0.12	0.86	0.86	0.86
243I	0.12	0.09	0.13	0.08	0.07	0.09	0.21	0.21	0.22
236Y	0.27	0.27	0.28	0.21	0.21	0.22	0.64	0.64	0.64
234M	0.33	0.34	0.33	0.19	0.19	0.19	0.88	0.88	0.88
234C	0.28	0.27	0.28	0.14	0.14	0.14	0.77	0.77	0.78
221T	0.04	0.04	0.04	0.10	0.11	0.10	0.41	0.42	0.41
206C	0.15	0.15	0.14	0.13	0.11	0.15	0.10	0.11	0.10
201M	0.37	0.37	0.37	0.31	0.31	0.31	1.04	1.04	1.04
201L	0.16	0.16	0.15	0.17	0.17	0.18	0.85	0.85	0.86
201I	0.20	0.20	0.19	0.20	0.19	0.21	0.71	0.70	0.71
165C	0.27	0.26	0.27	0.39	0.38	0.39	0.96	0.93	0.98
148T	0.12	0.12	0.13	0.14	0.14	0.14	0.12	0.12	0.13
148S	0.08	0.07	0.09	0.10	0.10	0.10	0.07	0.07	0.08
148N	0.19	0.17	0.20	0.25	0.24	0.26	0.40	0.39	0.41
148K	0.10	0.08	0.12	0.11	0.12	0.10	0.26	0.23	0.30
136L	0.14	0.14	0.15	0.16	0.16	0.16	0.43	0.44	0.43
136F	0.06	0.06	0.06	0.09	0.08	0.09	0.66	0.65	0.67
10Q	0.17	0.18	0.17	0.19	0.17	0.20	0.13	0.12	0.13
10A	0.23	0.22	0.23	0.17	0.17	0.17	0.53	0.53	0.53

**Table A 12: Beneficial mutations shared by only AmiE WT and AmiE I122L.**

Location-mutation	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE I122L cumulative normalized fitness metric	AmiE I122L technical replicate 1 normalized fitness metric	AmiE I122L technical replicate 2 normalized fitness metric	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric
74N	0.12	0.11	0.12	0.06	0.06	0.06	-0.07	-0.05	-0.08
82V	0.07	0.07	0.07	0.06	0.06	0.06	-0.09	-0.10	-0.08
88A	0.10	0.09	0.10	0.35	0.35	0.35	-0.09	-0.08	-0.10
92K	0.08	0.07	0.09	0.08	0.07	0.08	-0.13	-0.13	-0.13
112N	0.04	0.03	0.05	0.05	0.05	0.05	-0.05	-0.05	-0.06
136C	0.09	0.09	0.09	0.04	0.04	0.05	-0.11	-0.16	-0.06
154D	0.05	0.04	0.05	0.05	0.06	0.03	-0.11	-0.11	-0.10
201C	0.13	0.14	0.13	0.07	0.06	0.09	-0.51	-0.58	-0.45
201R	0.08	0.09	0.07	0.09	0.08	0.10	-0.11	-0.10	-0.13
201T	0.20	0.20	0.20	0.13	0.12	0.14	-0.19	-0.19	-0.20
234Q	0.06	0.07	0.05	0.14	0.15	0.12	-0.28	-0.34	-0.23
248A	0.11	0.12	0.10	0.14	0.15	0.14	-0.12	-0.15	-0.08

**Table A 13: Beneficial mutations shared by only AmiE WT and AmiE I38V.**

Location-mutation	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE II22L cumulative normalized fitness metric	AmiE II22L technical replicate 1 normalized fitness metric	AmiE II22L technical replicate 2 normalized fitness metric	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric
10H	0.06	0.07	0.05	-0.26	-0.26	-0.26	0.04	0.06	0.03
9L	0.14	0.13	0.15	-0.14	-0.13	-0.16	0.24	0.25	0.23
3Q	0.31	0.31	0.32	-0.07	-0.07	-0.06	0.66	0.65	0.66
339S	0.17	0.16	0.18	-0.03	-0.04	-0.03	0.20	0.20	0.20
237S	0.08	0.04	0.09	-0.19	-0.19	-0.20	0.10	0.13	0.07

**Table A 14: Beneficial mutations shared by only AmiE I38V and AmiE I122L.**

Location-mutation	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE I122L cumulative normalized fitness metric	AmiE I122L technical replicate 1 normalized fitness metric	AmiE I122L technical replicate 2 normalized fitness metric	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric
4A	-0.10	-0.13	-0.08	0.08	0.08	0.08	0.55	0.54	0.55
79C	-0.11	-0.11	-0.12	0.15	0.14	0.15	0.73	0.72	0.73
69H	-0.03	-0.04	-0.02	0.04	0.03	0.05	0.03	0.04	0.03
69R	-0.05	-0.07	-0.03	0.05	0.05	0.05	0.11	0.11	0.11

**Table A 15: AmiE WT unique beneficial mutations.**

Location-mutation	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE I122L cumulative normalized fitness metric	AmiE I122L technical replicate 1 normalized fitness metric	AmiE I122L technical replicate 2 normalized fitness metric	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric
5C	0.10	0.09	0.11	-0.20	-0.18	-0.21	-0.32	-0.32	-0.32
7G	0.19	0.15	0.21	-0.07	-0.02	-0.12	-0.10	-0.14	
47C	0.07	0.08	0.07	-0.44	-0.41	-0.46	-0.53	-0.51	-0.55
67C	0.08	-0.05	0.14	-0.08	-0.04	-0.14	-0.13	-0.12	-0.14
72T	0.08	0.09	0.08	-0.62	-0.60	-0.63	-0.60	-0.57	-0.63
74G	0.12	0.11	0.14	-0.17	-0.15	-0.19	-0.24	-0.25	-0.24
78H	0.11	0.10	0.12	-0.13	-0.13	-0.13	-0.54	-0.59	-0.49
82I	0.05	0.04	0.06	-0.04	-0.04	-0.03	-0.42	-0.45	-0.39
179V	0.11	-0.17	0.17	-0.05	-0.01	-0.10	-0.26	-0.26	-0.27
189Y	0.08	-0.09	0.12	-0.10	-0.05	-0.15	-0.31	-0.31	-0.30
228G	0.05	-0.10	0.09	-0.16	-0.11	-0.23	-0.20	-0.20	-0.21
234T	0.07	0.06	0.08	-0.08	-0.08	-0.09	-0.36	-0.35	-0.37
252D	0.07	-0.25	0.13	-0.21	-0.14	-0.32	-0.54	-0.47	-0.62
277C	0.11	-0.47	0.19	-0.49	-0.35	-0.72	-0.30	-0.29	-0.31
291C	0.07	0.10	0.05	-0.50	-0.36	-0.75	-0.40	-0.42	-0.39
297H	0.08	0.11	0.06	-0.32	-0.42	-0.26	-0.44	-0.46	-0.42
311C	0.16	0.20	0.14	-1.17	-1.11	-1.22	-1.07	-1.28	-0.94
320A	0.05	0.08	0.03	-0.39	-0.52	-0.32	-0.39	-0.44	-0.34
325C	0.14	0.13	0.14	-0.12	-0.13	-0.12	-0.22	-0.18	-0.27
325T	0.05	0.04	0.06	-0.32	-0.30	-0.33	-0.08	-0.06	-0.09
325Y	0.16	0.18	0.14	-0.10	-0.01	-0.23	-0.24	-0.22	-0.28
326H	0.22	0.23	0.22	-0.09	-0.10	-0.08	-0.21	-0.27	-0.16
326Y	0.07	0.09	0.06	-0.27	-0.28	-0.27	-0.51	-0.56	-0.47
336L	0.08	0.07	0.08	-0.04	-0.05	-0.04	-0.16	-0.17	-0.14
341Y	0.18	0.19	0.17	-0.25	-0.23	-0.27	-0.16	-0.14	-0.18

red highlighting indicates samples were dropped from analysis of beneficial mutations due to technical replicate disparities

**Table A 16: AmiE I122L unique beneficial mutations.**

Location-mutation	AmiE I122L cumulative normalized fitness metric	AmiE I122L technical replicate 1 normalized fitness metric	AmiE I122L technical replicate 2 normalized fitness metric	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric
47K	0.07	0.07	0.06	-0.07	-0.04	-0.08	-0.07	-0.08	-0.07
47R	0.03	0.03	0.02	-0.07	-0.07	-0.07	-0.13	-0.12	-0.14
58H	0.13	0.01	0.20	-0.24	-0.20	-0.27	-0.17	-0.14	-0.22
69Q	0.04	0.04	0.05	-0.04	-0.06	-0.03	-0.17	-0.16	-0.18
84I	0.17	0.17	0.16	-0.39	-0.44	-0.35	-0.85	-0.92	-0.79
84V	0.15	0.15	0.14	-0.14	-0.16	-0.13	-0.84	-0.83	-0.85
136V	0.05	0.06	0.04	-0.04	-0.04	-0.04	-0.04	-0.06	-0.03
336V	0.02	0.02	0.03	-0.04	-0.05	-0.04	-0.05	-0.03	-0.08

**Table A 17: AmiE I38V unique beneficial mutations.**

Location-mutation	AmiE I38V cumulative normalized fitness metric	AmiE I38V technical replicate 1 normalized fitness metric	AmiE I38V technical replicate 2 normalized fitness metric	AmiE WT cumulative normalized fitness metric	AmiE WT technical replicate 1 normalized fitness metric	AmiE WT technical replicate 2 normalized fitness metric	AmiE I122L cumulative normalized fitness metric	AmiE I122L technical replicate 1 normalized fitness metric	AmiE I122L technical replicate 2 normalized fitness metric
8V	0.30	0.30	0.31	-0.74	-0.81	-0.61	-0.36	-0.39	-0.33
21A	0.04	0.04	0.05	-0.19	-0.21	-0.18	-0.28	-0.29	-0.27
21P	0.38	0.38	0.38	-0.11	-0.14	-0.09	-0.10	-0.09	-0.11
21S	0.07	0.06	0.08	-0.19	-0.20	-0.19	-0.26	-0.26	-0.26
25C	0.14	0.15	0.14	-0.11	-0.12	-0.10	-0.04	-0.03	-0.04
65V	0.24	0.25	0.23	-0.28	-0.31	-0.27	-0.27	-0.25	-0.29
69Y	0.08	0.08	0.08	-0.12	-0.15	-0.10	-0.27	-0.30	-0.24
72L	0.04	0.03	0.04	-0.35	-0.32	-0.37	-0.43	-0.41	-0.45
75C	0.12	0.13	0.10	-0.39	-0.45	-0.35	-0.19	-0.19	-0.18
75L	0.12	0.11	0.13	-0.27	-0.27	-0.26	-0.23	-0.23	-0.22
102C	0.11	0.09	0.13	-0.24	-0.22	-0.25	-0.41	-0.42	-0.40
102I	0.54	0.53	0.54	-0.20	-0.20	-0.20	-0.07	-0.07	-0.07
143R	0.05	0.01	0.08	-0.19	-0.19	-0.20	-0.13	-0.13	-0.13
143T	0.30	0.29	0.30	-0.09	-0.11	-0.09	-0.04	-0.04	-0.04
149M	0.40	0.39	0.41	-0.25	-0.27	-0.24	-0.25	-0.23	-0.29
149R	0.06	0.04	0.08	-0.16	-0.16	-0.16	-0.36	-0.32	-0.42
190T	0.37	0.38	0.36	-0.49	-0.43	-0.53	-0.37	-0.42	-0.33
193P	0.53	0.54	0.53	-0.47	-0.48	-0.46	-0.08	-0.08	-0.07
197C	0.30	0.34	0.24	-1.20	-1.10	-1.26	-0.81	-0.95	-0.73
201N	0.25	0.27	0.24	-0.11	-0.11	-0.11	-0.07	-0.06	-0.08
204N	0.26	0.26	0.26	-0.28	-0.35	-0.25	-0.16	-0.13	-0.19
212C	0.51	0.50	0.53	-0.37	-0.35	-0.38	-0.33	-0.34	-0.32
212K	0.07	0.09	0.04	-0.33	-0.56	-0.27	-0.60	-0.46	-0.84
212T	0.40	0.40	0.40	-0.25	-0.10	-0.37	-0.24	-0.23	-0.25
212V	0.19	0.20	0.19	-0.48	-0.41	-0.52	-0.37	-0.37	-0.38
228A	0.15	0.13	0.16	-0.32	-0.31	-0.32	-0.25	-0.23	-0.28
230C	0.25	0.25	0.26	-0.18	-0.21	-0.16	-0.31	-0.31	-0.31
251N	0.47	0.47	0.47	-0.10	-0.13	-0.09	-0.15	-0.16	-0.15
253Y	0.33	0.33	0.32	-0.49	-0.63	-0.44	-0.39	-0.40	-0.37
254T	0.15	0.16	0.14	-0.42	-0.43	-0.41	-0.43	-0.43	-0.44
287T	0.12	0.13	0.10	-0.17	-0.17	-0.17	-0.21	-0.21	-0.20
288K	0.07	0.08	0.06	-0.07	-0.07	-0.07	-0.04	-0.04	-0.03
317M	0.17	0.18	0.16	-0.41	-0.34	-0.46	-0.23	-0.17	-0.29
325F	0.15	0.19	0.10	-0.16	-0.16	-0.16	-0.15	-0.09	-0.22
326R	0.06	0.04	0.08	-0.24	-0.22	-0.26	-0.23	-0.23	-0.24
333N	0.25	0.22	0.28	-0.30	-0.31	-0.29	-0.20	-0.21	-0.20
340K	0.15	0.16	0.13	-0.12	-0.11	-0.13	-0.05	-0.04	-0.05

**Table A 18: AmiE WT unique beneficial mutations with synonymous codon fitness disparities.**

Position	Mutation	Codon	Codon Frequency	Reference population count	Cumulative selected population count	Selected population count - technical replicate 1	Selected population count - technical replicate 2	Cumulative normalized fitness metric	Technical replicate 1 normalized fitness metric	Technical replicate 2 normalized fitness metric	Variance for synonymous codons of normalized fitness metrics - cumulative	Variance for synonymous codons of normalized fitness metrics - technical replicate 1	Variance for synonymous codons of normalized fitness metrics - technical replicate 2	Shared-unique beneficial mutation bin
3	L	CTA	0.04	109	2381	1050	1331	0.29	0.29	0.28	0.08	0.09	0.05	Beneficial for all - 95% CI cutoff
3	L	CTG	0.47	1	2	0	2	-0.27	None	-0.12	0.08	0.09	0.05	
3	L	CTT	0.12	29	109	43	66	-0.10	-0.12	-0.09	0.08	0.09	0.05	
3	L	TTC	0.10	1	1	5	5	None	None	None	0.08	0.09	0.05	
3	L	TTG	0.13	6	0	0	None	None	None	0.08	0.09	0.05		
3	Q	CAA	0.34	144	4028	1674	2354	0.13	0.31	0.34	0.52	0.32	1.06	WT and 138V shared beneficial - >10% increase in growth rate cutoff
3	Q	CAG	0.66	16	9	2	2	-0.60	0.47	0.41	0.52	0.32	1.06	Beneficial for all - 10% increase in growth rate cutoff
3	V	GTC	0.10	29	908	368	120	-0.27	0.21	0.21	0.12	0.06	0.34	
3	V	GTC	0.2	9	144	65	79	0.22	0.23	0.22	0.12	0.06	0.34	
3	V	GTC	0.35	3	6	5	1	-0.27	-0.09	-0.09	0.12	0.06	0.34	
3	V	GTC	0.29	2	1	1	None	None	None	0.12	0.06	0.34		
5	C	TGC	0.52	77	799	332	467	0.13	0.13	0.14	0.05	0.06	0.05	WT unique beneficial - 95% CI cutoff
5	C	TGA	0.46	18	48	19	29	-0.19	-0.21	-0.18	0.05	0.06	0.05	
5	C	TGA	0.40	1	1	1	1	-0.19	-0.19	-0.19	0.05	0.06	0.05	
7	G	GGC	0.37	121	1799	639	1160	0.21	0.17	0.23	0.24	0.37	0.21	Unique WT beneficial - >10% increase in growth rate cutoff
7	G	GGC	0.15	9	5	1	4	None	None	None	0.05	0.06	0.05	
8	A	GCA	0.25	1	11	7	4	0.15	0.23	0.05	0.42	0.39	0.46	Beneficial for all - 95% CI cutoff
8	A	GCA	0.26	10	2	1	1	None	None	None	0.42	0.42	0.46	
8	A	GCA	0.26	41	414	203	211	0.11	0.16	0.12	0.39	0.39	0.46	
8	A	GCA	0.18	1	1	1	1	-0.07	-0.17	-0.17	0.42	0.39	0.46	
8	K	AAA	0.75	2	5	0	5	-0.21	None	-0.07	0.13	None	0.07	Beneficial for all - 10% increase in growth rate cutoff
8	K	AAK	0.26	36	877	376	501	0.31	0.31	0.13	0.13	0.06	0.07	
8	K	AAK	0.26	1	1	1	1	-0.07	-0.07	-0.07	0.28	0.28	0.54	WT and 138V shared beneficial - >10% increase in growth rate cutoff
9	L	CTG	0.1	4	2	1	1	-0.74	-0.67	-0.80	0.76	0.28	0.54	
9	L	CTG	0.47	47	5	0	None	None	None	0.76	0.28	0.54		
9	L	CTG	0.10	68	1317	553	764	0.21	0.21	0.21	0.26	0.26	0.54	
9	L	CTG	0.10	69	20	49	62	0.12	0.04	0.17	0.76	0.28	0.54	
10	A	GCA	0.25	97	3524	1496	2028	0.18	0.38	0.38	0.12	0.12	0.12	Beneficial for all - 10% increase in growth rate cutoff
10	A	GCA	0.26	160	1003	463	620	0.04	0.04	0.04	0.12	0.12	0.12	
10	A	GCG	0.31	36	385	151	194	0.14	0.14	0.14	0.12	0.12	0.12	
10	A	GCT	0.18	46	73	30	43	-0.34	-0.35	-0.34	0.12	0.12	0.12	
10	A	GCT	0.26	124	1246	562	640	-0.04	-0.02	-0.02	0.03	0.03	0.03	WT and 138V shared beneficial - 95% CI cutoff
10	H	CAT	0.57	250	1238	664	724	0.21	0.21	0.21	0.11	0.14	0.09	WT unique beneficial - 95% CI cutoff
10	C	TGC	0.54	234	1225	1029	1296	0.12	0.13	0.12	0.11	0.14	0.09	WT unique beneficial - 95% CI cutoff
10	C	TGC	0.46	80	128	44	56	-0.34	-0.40	-0.30	0.11	0.14	0.09	
10	C	TGC	0.00	99	252	542	542	0.21	0.21	0.21	0.00	0.00	0.00	Beneficial for all - 95% CI cutoff
10	L	CTG	0.1	85	785	333	452	0.11	0.11	0.11	0.00	0.00	0.00	
10	L	CTG	0.47	73	621	258	363	0.09	0.08	0.09	0.00	0.00	0.00	
10	L	CTG	0.10	29	19	72	72	0.06	0.07	0.06	0.00	0.00	0.00	
10	L	TIA	0.14	9	67	30	37	0.06	0.07	0.05	0.00	0.00	0.00	
10	L	TIG	0.15	87	491	199	292	0.00	-0.02	0.00	0.00	0.00	0.00	
10	T	ACG	0.25	33	401	171	210	0.17	0.17	0.17	0.01	0.01	0.01	Beneficial for all - 95% CI cutoff
10	T	ACG	0.25	24	187	84	103	0.07	0.08	0.06	0.01	0.01	0.01	
10	T	ACG	0.25	10	66	46	56	0.04	0.04	0.04	0.00	0.00	0.00	
10	S	TGC	0.54	46	427	103	324	0.11	0.02	0.17	1.05	0.04	0.03	WT unique beneficial - 95% CI cutoff
10	C	TGT	0.46	7	1	0	1	None	None	None	0.04	0.04	0.03	
10	T	ACA	0.17	10	62	290	393	0.17	0.13	0.13	0.62	0.31	0.42	WT unique beneficial - 95% CI cutoff
10	T	ACA	0.17	59	602	259	386	0.17	0.13	0.13	0.62	0.31	0.42	
10	T	ACA	0.25	27	122	55	67	-0.06	-0.04	-0.07	0.62	0.31	0.42	
10	T	ACA	0.19	5	1	1	1	None	None	None	0.04	0.04	0.03	
10	G	GCA	0.11	132	655	254	399	-0.03	-0.03	-0.02	0.62	0.31	0.42	WT unique beneficial - 95% CI cutoff
10	G	GCA	0.37	23	86	35	51	-0.10	-0.12	-0.10	0.03	0.03	0.03	
10	G	GCA	0.15	58	1045	448	597	0.23	0.25	0.25	0.03	0.03	0.03	
10	G	GCA	0.26	100	100	49	56	0.06	0.07	0.06	0.00	0.00	0.00	
10	G	GCA	0.43	31	401	101	157	-0.03	-0.05	-0.02	0.03	0.03	0.03	WT unique beneficial - 95% CI cutoff
10	A	GCA	0.26	55	651	266	385	-0.03	-0.05	-0.01	0.01	0.01	0.00	
10	A	GCA	0.33	43	991	420	571	0.12	0.13	0.13	0.00	0.01	0.00	
10	A	GCA	0.19	151	1519	655	864	0.16	0.17	0.15	0.00	0.00	0.00	
10	V	GIA	0.17	14	1	1	1	None	None	None	0.60	0.28	0.65	WT unique beneficial - 95% CI cutoff
10	V	GIC	0.2	226	3082	289	2793	0.16	-0.11	0.21	0.60	0.28	0.65	
10	V	GIA	0.35	2	11	4	4	None	None	None	0.60	0.28	0.65	
10	V	GIA	0.28	16	7	0	2	None	None	None	0.60	0.28	0.65	
10	Y	TAC	0.17	51	1441	60	1381	0.34	-0.13	0.41	3.49	None	2.60	WT unique beneficial - 95% CI cutoff
10	C	TGC	0.41	27	493	186	268	0.07	-0.04	-0.03	0.01	0.02	0.01	WT and 1122I shared beneficial - 95% CI cutoff
10	T	ACT	0.19	58	557	177	380	0.08	0.09	0.08	0.00	0.01	0.01	WT and 1122I shared beneficial - 95% CI cutoff
10	Q	CAA	0.34	49	325	97	228	0.00	-0.01	0.00	0.01	0.01	0.01	WT unique beneficial - 95% CI cutoff
10	Q	CAA	0.66	33	372	121	124	0.14	0.14	0.14	0.01	0.01	0.01	
10	Q	CAA	0.26	62	140	40	44	0.06	0.07	0.06	0.02	0.02	0.02	WT unique beneficial - 95% CI cutoff
10	T	ACG	0.4	295	3407	1004	2403	0.12	0.13	0.02	0.04	0.04	0.02	WT unique beneficial - 95% CI cutoff
10	T	ACG	0.25	34	91	20	29	0.17	-0.23	-0.23	0.02	0.02	0.02	
10	S	TGC	0.25	20	19	7	12	-0.55	-0.49	-0.59	0.11	0.09	0.12	WT and 138V shared beneficial - 95% CI cutoff
10	S	TGC	0.16	401	401	120	281	-0.24	-0.24	-0.24	0.11	0.09	0.12	
10	S	TGC	0.26	33	33	23	23	-0.24	-0.24	-0.24	0.11	0.09	0.12	
10	S	TGC	0.41	282	493	186	268	-0.24	-0.24	-0.24	0.06	0.06	0.04	Beneficial for all - 10% increase in growth rate cutoff
10	S	TGC	0.39	25	147	44	103	-0.03	-0.03	-0.03	0.04	0.06	0.04	WT unique beneficial - 95% CI cutoff
10	A	GCA	0.26	31	191	47	144	0.16	0.19	0.14	0.07	0.09	0.06	Beneficial for all - 95% CI cutoff
10	A	GCA	0.26	45	380	91	258	0.08	0.07	0.10	0.01	0.01	0.01	
10	A	GCG	0.33	11	64	21	43	-0.03	-0.01	-0.04	0.01	0.01	0.01	
10	A	GCT	0.18	18	219	62	157	0.17	0.20	0.16	0.01	0.01	0.01	
10	A	GCT	0.15	53	161	47	94	-0.30	-0.30	-0.30	0.30	0.13	0.22	WT unique beneficial - 95% CI cutoff
10	T	ACA	0.17	18	302	109	193	0.25	0.25	0.25	0.01	0.02	0.01	WT unique beneficial - 95% CI cutoff
10	T	ACA	0.17	18	302	109	193	0.25	0.25	0.25	0.01	0.02	0.01	
10	T	ACA	0.19	41	356	120	236	0.11	0.10	0.13	0.01	0.02	0.01	
10	D	TAI	0.17	147	1353	450	705	0.09	0.11	0.07	0.08	0.11	0.11	WT unique beneficial - 95% CI cutoff
10	D	TAI	0.59	11	16	3	8	-0.35	-0.29	-0.40	0.10	0.08	0.11	
10	S	AGC	0.25	82	429	145	284	-0.01	-0.02	0.00	0.09	0.00	0.09	Beneficial for all - 95% CI cutoff
10	S	AGC	0.16	79	219	161	358	0.05	0.01	0.06	0.00	0.00	0.00	
10	S	AGC	0.17	101	201	49	106	0.05	0.05	0.05	0.00	0.00	0.00	
10	L	TCT	0.12	148	2191									

**Table A 19: AmiE I122L unique beneficial mutations with synonymous codon fitness disparities.**

Position	Mutation	Codon	Codon Frequency	Reference population count	Cumulative selected population count	Selected population count - technical	Cumulative normalized fitness metric	Technical replicate 1 normalized fitness metric	Variance for synonymous codons of normalized fitness metrics - cumulative	Variance for synonymous codons of normalized fitness metrics - technical	Variance for synonymous codons of normalized fitness metrics - cumulative	Shared-unique beneficial mutation bin
3	I	ATA	0.11	14	117	57	60	0.19	0.19	0.18	0.40	0.37
3	I	ATC	0.39	5	160	64	96	0.43	0.40	0.46	0.40	0.37
3	I	ATT	0.49	8	2	1	1	0.50	-0.74	0.39	0.40	0.37
3	L	CTA	0.04	29	438	214	224	0.30	0.30	0.30	0.27	0.14
3	L	CTC	0.1	7	5	3	2	-0.41	-0.34	0.09	0.27	0.14
3	L	CTG	0.47	5	1	0	1	0.85	None	0.01	0.27	0.14
3	L	CTT	0.12	24	93	37	56	0.03	-0.01	0.06	0.27	0.14
3	L	TTA	0.14	10	168	86	82	0.32	0.33	0.31	0.27	0.14
3	L	TTG	0.13	4	1	1	0	0.50	-0.50	None	0.27	0.14
3	V	GTA	0.17	16	105	53	52	0.14	0.15	0.13	0.19	0.12
3	V	GTC	0.2	1	55	25	30	0.52	0.51	0.53	0.19	0.12
3	V	GTG	0.35	2	4	2	2	-0.13	-0.12	-0.14	0.19	0.12
3	V	GTT	0.28	4	2	2	0	0.82	-0.30	None	0.19	0.11
4	A	GCA	0.23	28	141	57	84	0.08	0.05	0.11	0.07	0.07
4	A	GCC	0.26	45	395	197	198	0.20	0.20	0.19	0.07	0.07
4	A	GCG	0.33	19	58	34	24	-0.03	0.02	-0.08	0.07	0.07
4	A	GCT	0.18	32	22	11	11	-0.42	-0.40	-0.44	0.07	0.07
4	S	AGC	0.25	61	335	171	164	0.10	0.12	0.09	0.01	0.01
4	S	AGT	0.16	17	40	26	14	-0.09	-0.02	-0.19	0.01	0.03
4	S	TCA	0.14	14	25	17	8	-0.16	-0.07	-0.29	0.01	0.03
4	S	TCC	0.15	15	69	33	36	0.06	0.06	0.06	0.01	0.03
4	S	TCG	0.14	7	48	29	19	0.15	0.19	0.09	0.01	0.03
4	S	TCT	0.17	5	23	11	12	0.06	0.06	0.06	0.01	0.03
8	A	GCA	0.23	17	82	33	49	0.07	0.04	0.10	0.23	0.17
8	A	GCC	0.26	38	8	5	3	-0.83	-0.72	-0.98	0.23	0.17
8	A	GCG	0.33	154	691	318	373	0.06	0.05	0.07	0.23	0.17
8	A	GCT	0.18	20	214	90	124	0.24	0.21	0.26	0.23	0.17
9	Q	CAA	0.34	16	216	96	120	0.28	0.27	0.29	0.11	0.08
9	Q	CAG	0.66	28	45	26	19	-0.18	-0.13	-0.24	0.11	0.08
10	A	GCA	0.23	122	1343	641	702	0.24	0.24	0.24	0.12	0.11
10	A	GCC	0.26	119	584	269	315	0.08	0.07	0.08	0.12	0.11
10	A	GCG	0.33	99	1582	757	825	0.31	0.31	0.31	0.12	0.11
10	A	GCT	0.18	122	76	34	42	-0.45	-0.47	-0.44	0.12	0.11
47	R	AGA	0.07	154	423	215	208	-0.05	-0.04	-0.07	0.00	0.00
47	R	AGG	0.04	108	462	224	238	0.05	0.05	0.05	0.00	0.00
47	R	CGA	0.07	314	1132	581	551	0.01	0.03	-0.01	0.00	0.00
47	R	CGC	0.36	390	1676	835	841	0.05	0.06	0.04	0.00	0.00
47	R	CGG	0.11	149	607	284	323	0.04	0.03	0.04	0.00	0.00
47	R	CGT	0.36	230	883	442	441	0.02	0.03	0.01	0.00	0.00
69	R	AGA	0.07	102	327	145	182	-0.02	-0.03	0.00	0.00	0.00
69	R	AGG	0.04	51	300	135	165	0.12	0.10	0.13	0.00	0.00
69	R	CGA	0.07	156	589	271	318	0.02	0.01	0.03	0.00	0.00
69	R	CGC	0.36	203	899	445	454	0.06	0.06	0.05	0.00	0.00
69	R	CGG	0.11	124	537	243	294	0.05	0.04	0.06	0.00	0.00
69	R	CGT	0.36	159	780	371	409	0.08	0.08	0.08	0.00	0.00
89	S	AGC	0.25	58	424	240	184	0.15	0.15	0.14	0.01	0.01
89	S	AGT	0.16	52	178	109	69	-0.02	0.01	-0.05	0.01	0.01
89	S	TCA	0.14	27	90	54	36	-0.02	0.00	-0.05	0.01	0.01
89	S	TCC	0.15	23	95	54	41	0.03	0.04	0.01	0.01	0.01
89	S	TCG	0.14	12	97	70	27	0.17	0.22	0.06	0.01	0.01
89	S	TCT	0.17	26	88	56	32	-0.02	0.02	-0.07	0.01	0.01
93	A	GCA	0.23	128	672	371	301	0.08	0.08	0.07	0.00	0.00
93	A	GCC	0.26	128	393	215	178	-0.04	-0.04	-0.04	0.00	0.00
93	A	GCG	0.33	86	410	225	185	0.06	0.06	0.05	0.00	0.00
93	A	GCT	0.18	76	259	157	102	-0.02	0.01	-0.05	0.00	0.00
136	V	GTA	0.17	122	750	395	355	0.11	0.10	0.12	0.00	0.00
136	V	GTC	0.2	169	600	339	261	-0.01	0.00	-0.02	0.00	0.00
136	V	GTG	0.35	95	378	222	156	0.02	0.04	-0.01	0.00	0.00
136	V	GTT	0.28	106	535	321	214	0.07	0.09	0.04	0.00	0.00
148	S	AGC	0.25	69	397	214	183	0.10	0.10	0.10	0.00	0.00
148	S	AGT	0.16	23	133	74	59	0.10	0.10	0.09	0.00	0.00
148	S	TCA	0.14	25	169	82	87	0.13	0.11	0.15	0.00	0.00
148	S	TCC	0.15	53	373	219	154	0.14	0.15	0.12	0.00	0.00
148	S	TCG	0.14	38	212	118	94	0.09	0.10	0.08	0.00	0.00
148	S	TCT	0.17	31	101	52	49	-0.03	-0.04	-0.01	0.00	0.00
154	D	GAC	0.37	49	288	157	131	0.10	0.10	0.10	0.01	0.01
154	D	GAT	0.63	55	189	118	71	-0.01	0.02	-0.06	0.01	0.01
244	A	GCA	0.23	30	215	102	113	0.20	0.19	0.21	0.01	0.02
244	A	GCC	0.26	70	495	253	242	0.19	0.20	0.19	0.01	0.02
244	A	GCG	0.33	10	65	26	39	0.18	0.13	0.21	0.01	0.02
244	A	GCT	0.18	94	204	110	94	-0.05	-0.03	-0.08	0.01	0.02
336	V	GTA	0.17	233	360	172	188	0.05	0.06	0.04	0.00	0.00
336	V	GTC	0.2	287	341	149	192	-0.01	-0.02	0.00	0.00	0.00
336	V	GTG	0.35	233	354	158	196	0.05	0.04	0.05	0.00	0.00
336	V	GTT	0.28	190	254	104	150	0.02	-0.01	0.04	0.00	0.00

**Table A 20: AmiE I38V unique beneficial mutations with synonymous codon fitness disparities.**

Position	Mutation	Codon	Codon Frequency	Reference population count	Cumulative selected population count	Selected population count - technical replicate 1	Selected population count - technical replicate 2	Cumulative normalized fitness metric	Technical replicate 1 normalized fitness metric	Technical replicate 2 normalized fitness metric	Variance for synonymous codons of normalized fitness metrics - cumulative	Variance for synonymous codons of normalized fitness metrics - technical replicate 1	Variance for synonymous codons of normalized fitness metrics - technical replicate 2	Shared-unique beneficial mutation bin	
Beneficial for all ->10% increase in growth rate cutoff															
3	I	AAC	0.38	23	104	42	42	0.33	0.55	0.54	0.13	0.11	0.18		
3	I	ATT	0.49	5	12	8	4	-0.03	0.03	-0.14	0.13	0.11	0.18		
3	L	ATG	0.04	24	210	99	111	0.25	0.23	0.26	0.34	0.22	0.11	Beneficial for all - 95% CI cutoff	
3	L	CTG	0.47	12	11	1	0	-0.33	-0.33	-0.33	0.34	0.22	0.11		
3	L	CTT	0.12	43	75	30	45	-0.11	-0.17	-0.07	0.34	0.22	0.11		
3	L	TIA	0.14	13	65	31	34	0.13	0.12	0.15	0.34	0.22	0.11		
4	A	GCA	0.17	15	110	594	555	0.56	0.65	0.65	0.18	0.16	0.16	138V and 1122I shared beneficial - 95% CI cutoff	
4	A	GCC	0.26	38	2513	1277	1236	0.61	0.68	0.61	0.19	0.28	0.16		
4	A	GCG	0.33	16	44	22	22	0.00	0.00	0.00	0.19	0.28	0.16		
4	A	GCT	0.18	12	13	3	10	-0.24	-0.43	-0.13	0.19	0.28	0.16		
4	A	GTC	0.17	17	17	92	850	449	0.53	0.53	0.53	0.18	0.16	0.16	Beneficial for all - 95% CI cutoff
8	A	GCC	0.26	20	11	6	7	-0.43	-0.42	-0.40	0.31	0.29	0.33		
8	A	GCG	0.33	69	2591	1350	1241	-0.51	0.57	0.57	0.31	0.29	0.33		
8	A	GCA	0.18	180	915	915	915	0.94	0.89	0.81	0.30	0.29	0.33		
8	V	GTA	0.17	4	456	229	227	-0.01	-0.01	-0.01	0.02	0.14	0.02	138V unique beneficial - >10% increase in growth rate cutoff	
8	V	GTC	0.2	16	1	0	0	None	None	None	0.74	0.65			
8	V	GTT	0.28	7	143	75	68	0.41	0.41	0.41	0.92	0.74	0.62		
9	L	CTG	0.17	15	3	0	0	None	None	None	0.44	0.44	0.44	138V and WT shared beneficial - >10% increase in growth rate cutoff	
9	L	CTG	0.47	16	4	11	-0.28	-0.47	-0.18	0.44	0.47	0.44			
9	L	CTT	0.12	11	4	3	0	-0.45	-0.45	-0.45	0.44	0.47	0.44		
9	L	CTG	0.13	3	213	144	87	-0.03	-0.03	-0.03	0.44	0.47	0.44		
9	Q	CAG	0.66	8	890	489	401	0.59	0.60	0.50	0.28	0.33	0.24	Beneficial for all - 95% CI cutoff	
9	Q	CAG	0.67	57	826	422	402	-0.17	-0.16	-0.12	0.28	0.33	0.24		
10	A	GCA	0.26	85	230	113	120	0.00	-0.01	0.01	0.29	0.28	0.29		
10	A	GCC	0.36	31	2123	1059	1064	-0.61	0.65	0.65	0.29	0.28	0.29		
10	A	GCG	0.18	18	53	24	27	0.21	0.21	0.20	0.28	0.28	0.29		
10	H	CAT	0.99	99	491	243	248	0.13	0.13	0.13	0.03	0.02	0.05	138V and WT shared beneficial - 95% CI cutoff	
10	H	CAT	0.57	103	185	117	68	-0.11	-0.05	-0.19	0.03	0.07	0.05		
21	A	GCA	0.23	69	273	136	137	0.08	0.08	0.09	0.01	0.01	0.01	138V unique beneficial - 95% CI cutoff	
21	A	GCC	0.26	86	158	79	-0.10	-0.10	-0.10	0.01	0.01	0.01			
21	A	GCA	0.18	72	253	132	123	0.05	0.05	0.05	0.00	0.00	0.00		
21	A	GCT	0.18	63	266	131	155	0.11	0.09	0.13	0.01	0.01	0.01		
21	S	AGG	0.25	61	271	129	142	0.11	0.09	0.12	0.00	0.01	0.00	138V unique beneficial - 95% CI cutoff	
21	S	AGT	0.16	63	233	103	120	0.06	0.04	0.08	0.00	0.01	0.00		
21	S	ATC	0.17	69	247	127	121	0.06	0.05	0.05	0.00	0.01	0.00		
21	S	ATC	0.15	84	224	93	131	-0.01	-0.06	0.03	0.00	0.01	0.00		
21	S	TCG	0.14	59	198	95	103	0.04	0.03	0.06	0.00	0.01	0.00		
21	S	TCT	0.17	29	190	106	84	0.19	0.21	0.17	0.00	0.01	0.00		
21	S	TCA	0.26	88	184	85	85	0.00	0.00	0.00	0.00	0.00	0.00	138V and 1122I shared beneficial - 95% CI cutoff	
69	R	AGG	0.04	50	745	369	376	-0.35	-0.35	-0.36	0.02	0.02	0.03		
69	R	CGA	0.07	102	270	139	131	-0.01	-0.01	-0.02	0.02	0.02	0.02		
69	R	CAG	0.36	123	518	277	261	0.10	0.10	0.10	0.02	0.02	0.02		
69	R	CGT	0.25	121	460	231	229	0.07	0.07	0.07	0.02	0.02	0.03		
72	L	CTA	0.04	34	161	81	80	0.12	0.12	0.12	0.00	0.00	0.00	138V unique beneficial - 95% CI cutoff	
72	L	CTG	0.1	87	307	155	152	0.06	0.06	0.06	0.00	0.00	0.00		
72	L	CTG	0.17	107	107	50	57	-0.01	-0.03	0.01	0.00	0.00	0.00		
72	L	TTC	0.13	85	227	115	112	-0.01	-0.01	-0.01	0.00	0.00	0.00		
72	A	GCA	0.26	65	182	85	99	-0.03	-0.03	-0.03	0.05	0.05	0.06	Beneficial for all - 95% CI cutoff	
72	A	GCC	0.26	77	237	77	38	-0.14	-0.12	-0.09	0.05	0.05	0.06		
93	A	GCG	0.33	41	820	440	380	0.38	0.38	0.39	0.05	0.05	0.06		
93	A	GCT	0.18	22	45	27	18	-0.11	-0.08	-0.15	0.05	0.05	0.06	138V unique beneficial - 95% CI cutoff	
102	C	TGT	0.46	11	4	2	1	-0.14	-0.14	-0.14	0.32	0.32	0.32	138V unique beneficial - 95% CI cutoff	
143	R	AGA	0.07	21	32	17	15	-0.18	-0.18	-0.18	0.02	0.01	0.03	138V unique beneficial - 95% CI cutoff	
143	R	AGG	0.04	20	73	50	23	-0.14	-0.07	-0.24	0.02	0.01	0.03		
143	R	CCT	0.07	40	56	31	25	-0.23	-0.23	-0.23	0.03	0.03	0.03		
143	R	CCT	0.36	154	1055	453	602	0.17	0.12	0.22	0.02	0.01	0.03		
143	R	CGG	0.11	59	112	57	55	-0.13	-0.14	-0.11	0.02	0.01	0.03		
143	R	CGT	0.36	53	90	49	41	-0.16	-0.16	-0.16	0.02	0.01	0.03		
197	C	TGC	0.54	13	3	1	2	-0.04	-0.04	-0.04	0.65	0.90	0.48	138V unique beneficial - >10% increase in growth rate cutoff	
197	C	TGT	0.46	12	199	124	75	-0.40	-0.40	-0.39	0.65	0.90	0.48		
197	N	AGG	0.53	74	91	83	80	-0.01	-0.01	-0.02	0.06	0.06	0.06	138V unique beneficial - >10% increase in growth rate cutoff	
200	N	AGT	0.49	118	1150	617	533	-0.33	-0.35	-0.35	0.06	0.06	0.06		
228	A	GCA	0.23	36	7	4	3	-0.01	-0.01	-0.01	0.29	0.30	0.31	138V unique beneficial - >10% increase in growth rate cutoff	
228	A	GCC	0.26	36	393	186	207	0.36	0.34	0.37	0.29	0.30	0.31		
228	A	GCG	0.33	18	4	2	2	-0.01	-0.01	-0.01	0.29	0.30	0.31		
228	A	GCT	0.17	23	1	1	0	-0.01	-0.01	-0.01	0.29	0.30	0.31		
237	S	AGG	0.25	53	47	24	23	-0.24	-0.23	-0.25	0.06	0.06	0.05	138V unique beneficial - 95% CI cutoff	
237	S	AGT	0.16	84	111	52	59	-0.13	-0.14	-0.11	0.06	0.06	0.05		
237	S	ATC	0.14	54	77	46	41	-0.13	-0.06	-0.09	0.04	0.04	0.05		
237	S	ATC	0.18	80	183	126	126	-0.24	-0.24	-0.24	0.04	0.04	0.05		
237	T	ACA	0.17	24	239	136	103	0.31	0.32	0.29	0.06	0.05	0.07	138V unique beneficial - 95% CI cutoff	
287	T	ACC	0.17	28	37	32	40	0.04	0.04	0.07	0.06	0.05	0.07		
287	T	ACC	0.4	25	72	41	21	16	0.15	0.15	0.15	0.06	0.05	0.07	
287	T	ACC	0.19	6	7	4	3	-0.19	-0.18	-0.23	0.03	0.03	0.07		
288	K	AAG	0.74	168	855	464	391	0.17	0.18	0.16	0.06	0.05	0.05	138V unique beneficial - 95% CI cutoff	
288	K	AAG	0.26	164	214	115	99	-0.16	-0.15	-0.18	0.06	0.05	0.06		
326	F	TTC	0.42	54	288	109	109	0.18	0.22	0.13	0.18	0.17	0.08	138V unique beneficial - >10% increase in growth rate cutoff	
326	F	TTC	0.17	5	4	4	4	0.00	0.00	0.00	0.17	0.17	0.08		
326	R	AGA	0.07	32	43	19	24	-0.15	-0.19	-0.12	0.09	0.10	0.10	138V unique beneficial - 95% CI cutoff	
326	R	AGG	0.04	15	9	6	3	-0.39	-0.31	-0.31	0.09	0.10	0.10		
326	R	CGA	0.07	34	75	42	43	-0.33	-0.32	-0.31	0.09	0.10	0.10		
326	R	CGG	0.36	26	481	228	233	0.23	0.23	0.21	0.09	0.10	0.10		
326	R	CGG	0.11	24	13	5	8	0.03	0.03	0.06	0.09	0.10	0.10		
326	R	CGT	0.36	27	10	4	6	0.06	0.06	0.04	0.09	0.10	0.10		
329	S	AGG	0.25	72	156	68	88	-0.03	-0.07	0.01	0.02	0.03	0.02	Beneficial for all - 95% CI cutoff	
329	S	AGG	0.17	50	55	45	45	0.00	0.00	0.00	0.03	0.03	0.02		
329	S	TCA	0.14	35	97	47	50								

## **APPENDIX B**

### **Chapter 4 supporting information**

**Note B 1: DNA sequences of ΦX174 viral capsid protein gene fragments.** Sequences provided span from the 5' EcoRV to the 3' BamHI site (bolded sequences). Mutagenic regions are highlighted in pink, BsmBI sites are highlighted in green.

>Fragment\_F1

**GATATCTGCAGAATTGCCCTT**CGTCTCATTCAAACGGCCTGTCATCATGGAAAGG  
CGCTGAATTACGGAAAACATTATTAAATGGCGTCGAGCGTCCGGTAAAGCCGCTGA  
ATTGTTCGCGTTACCTCGGTACGCGCAGGAAACACTGACGTTCTACTGACGC  
AGAAGAAAACGTGCGTAAAAATTACGTGCAGAAGGAGTGTAAATTCTAAAGG  
TAAAAAAACGTTCTGGCGCTCGCCCTGGTCGTCCGCAGCCGTTGCGAGGTACTAAAGG  
CAAGCGTAAAGGCCTCGTCTTGGTATGTAGGTGGTCAACAATTAAATTGCAGGG  
GCTTCGGCCCCCTACTTGAGGATAAATTATGTCTAATATTCAAACACTGGCGCCGAGCG  
**TATGCCGCATGACCTTCCCATCTGGCTTCCTGCTGGTCAGATTGGCGTCTTATT**  
ACCATTCAACTACTCCGGTTATCGCTGGCGACTCCTCGAGATGGACGCCGTTGGC  
GCTCTCCGTCTTCTCCATTGCGTCGTGGCCTTGCTATTGACTACTGTAGACATTAAAG  
TACTTTTATGTCCCTCATCGTACGTTATGGTGAACAGTGGATTAAGTTCATGAAG  
**GATGGTGTAAATGCCACTCCTCTCCGACTGT**GAGACGAAGGGCGAATTCCAGCACA  
CTGGCGGCCGTTACTAGT**GGATCC**

>Fragment\_F2

**GATATCTGCAGAATTGCCCTT**CGTCTC GACTGTTAACACTACTGGTTATATTGACC  
ATGCCGCTTTCTTGGCACGATTAACCCGTATACCAATAAAATCCCTAACGATTGTT  
TCAGGGTTATTGAATATCTATAACAACATTTAAAGCGCCGTGGATGCCTGACCG  
TACCGAGGCTAACCCCTAATGAGCTTAATCAAGATGATGCTCGTTATGGTTCCGTTG  
CTGCCATCTCAAAAACATTGGACTGCTCCGCTCCTCCTGAGACTGAGCTTCTCGC  
CAAATGACGACTTCTACCACATCTATTGACATTATGGGTCTGCAAGCTGCTTATGCT  
AATTGACATACTGACCAAGAACGTGATTACTCATGCAGCGTTACCGTGATGTTATT  
TCTTCATTGGAGGTAAAACCTCTTATGACGCTGACAACCGCCTTACTGTCTGAC  
**GCTCTAACTCTGGG**GAGACGAAGGGCGAATTCCAGCACACTGGCGCCGTTACTA  
**GTGGATCC**

>Fragment\_F3

**GATATCTGCAGAATTGCCCTT**CGTCTCCTGGGCATCTGGCTATGATGTTGATGGA  
ACTGACCAAACGTCGTTAGGCCAGTTTCTGGCGTGTCAACAGACCTATAAACAT  
TCTGTGCCGCGTTCTTGTCTGAGCATGGCACTATGTTACTCTTGCCTGTTGAC  
TTTCCGCTACTGCGACTAAAGAGATTCACTAACGCTAAAGGTGCTTGT  
TTATACCGATATTGCTGGCGACCCCTGTTGTATGGCAACTGCCGCCGTGAAATT  
TCTATGAAGGATGTTTCCGTTCTGGTATTGCTTAAGAAGTTAAGATTGCTGAG  
GGTCAGTGGTATCGTTATGCGCCTCGTATGTTCTCCTGCTTATCACCTTGTGAAG  
GCTTCCCATTCAATTAGGAACCGCCTCTGGTATTGCAAGAACGCGTACTTATTG  
CCACCATGATTATGACCAAGTGTTCAGTCCCGTCAAGTGTGAGTGGAAATAGTCA  
GGTAAATTAAATGTGACCGTTATCGCAATCTGCCGACCACTCGCGATTCAATCAT  
GACTTCGTGATAAAAGATTGAGTGTGANNNNNNNNNGTTATAACGCCGAAGC  
GGTAAAAATTAAATTGGCCGCTGAGGGGTTGACCAAGCGAAGC**GAGACGAANG**  
CGAATTCCAGCACACTGGCGCCGTTACTAGT**GGATCC**

>Fragment\_G1

GATATCTGCAGAATTGCCCTT CGTCTCCGAAGCGCGTAGGTTCTGCTTAGGAG  
TTAACATGTTTCAGACTTTATTCTCGCCATAATTCAAACCTTTCTGATAAGC  
TGGTTCTCACTCTGTTACTCCAGCTTCCGGCACCTGTTACAGACACCTAAAGC  
TACATCGTCAACGTTATTTGATAGTTGACGGTTAATGCTGGTAATGGTGGTTT  
CTTCATTGCATTAGATGGATACATCTGTCAACGCCGCTAATCAGGTTCTGTTG  
GTGCTGATATTGCTTTGATGCCGACCCTAAATTGCTGTTGGTCGCTTGA  
GTCTTCTCGGTTCCGACTACCCTCCCAGCTGCCTATGATGTTATCCTTGAATGGT  
CGCCATGATGGTGGTTATTATACCGTCAAGGACTGTGACTATTGACGTCCTCCC  
CGTACGGAGACGAAGGGCGAATTCCAGCACACTGGCGGCCGTTACTAGT**GGATCC**

>Fragment\_G2

GATATCTGCAGAATTGCCCTT CGTCTCCGTACGCCGGCAATAATGTTATGTTGG  
TTTCATGGTTGGTCTAACCTTACCGCTACTAAATGCCCGGGATTGGTTCGCTGAAT  
CAGGTTATTAAAGAGATTATTGTCTCCAGCCACTTAAGTGAGGTGATTATGTTG  
GTGCTATTGCTGGCGGTATTGCTCTGCTCTGCTGGTGGCGCCATGTCTAAATTGTT  
TGGAGGCGGTAAAAAGCCGCCTCCGGTGGCATTCAAGGTGATGTGCTTGCTACCG  
ATAACAATACTGTAGGCATGGGTGATGCTGGTATTAAATCTGCCATTCAAGGCTCTA  
ATGTTCTAACCTGATGAGGCCCGCCCTAGTTTGTCTGGTGTATGGCTAAAGC  
TGGTAAGGACTTCTGAAGGTACGTTGCAGGCTGGCACTTCTGCCGTTCTGATAA  
GTTGCTTGATTGGTGGACTTGGCAAGTCTGCCGCTGATAAAGGAAAGGATAC  
TCG**GAGACGAAGGGCGAATTCCAGCACACTGGCGGCCGTTACTAGTGGATCC**

**Note B 2: Amino acid sequences of ΦX174 mutagenized viral capsid protein gene fragments.** Only mutagenized regions of gene fragments shown.

>Fragment\_F1

MSNIQTGAERMPHDLSHLGFLAGQIGRLITISTPVIAGDSFEMDAVGALRLSPLRRGLAI  
DSTVDIFTFYVPHRHVYGEQWIKFMKDGVNATPL

>Fragment\_F2

NTTGYIDHAAFLGTINPDTNKIPKHLFQGYLNIYNNYFKAPWMPDRTEANPNELNQDDA  
RYGFRCCHLKNIWTAPLPETELSRQMTTSTSIDIMGLQAYANLHTDQERDYFMQRY  
RDVISSFGGKTSYDADNRPLLVMRSN

>Fragment\_F3

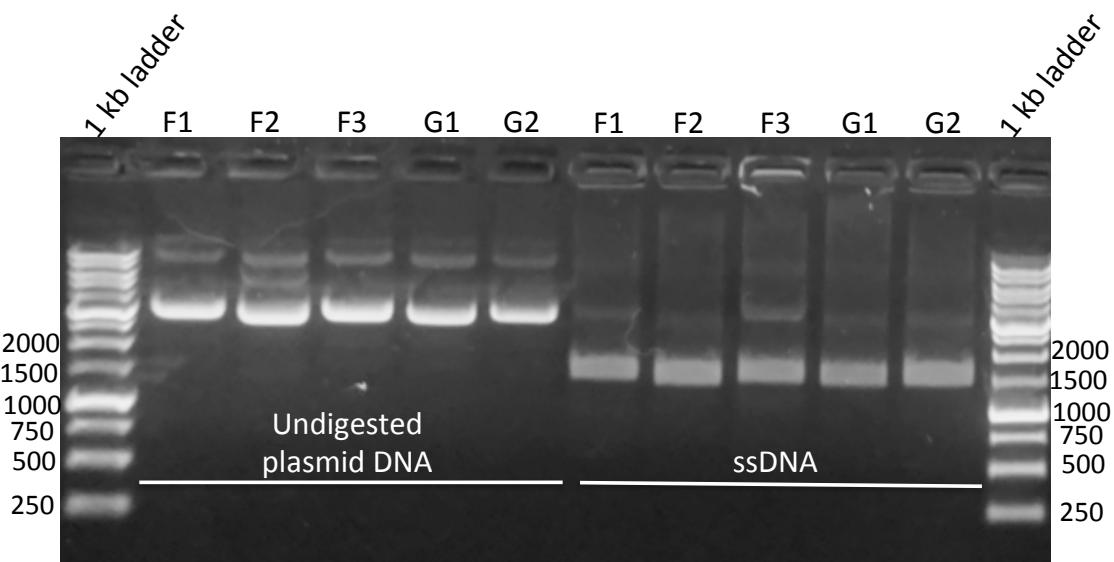
SGYDVGTDQTSLGQFSGRVQQTYKHSVPRFFVPEHGTMFTLALVRFPPATKEIQYLN  
AKGALTYTDIAGDPVLYGNLPPREISMKDVFRSGDSSKKFKIAEGQWYRYAPSYVSPAY  
HLLEGFPFIQEPPSGDLQERVLIRHHDYDQCFQSVQLLQWNSQVKFNVTVYRNLPTRD  
SIMTS

>Fragment\_G1

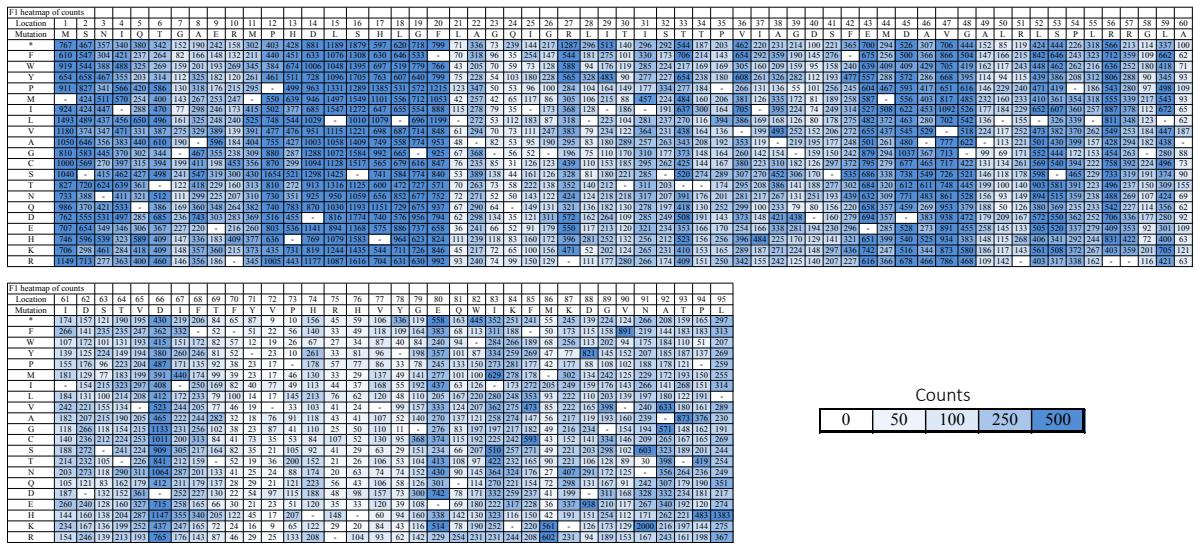
FQTFISRHNNSNFFSDKLVLTSPASSAPVLQTPKATSSTLYFDSLTVNAGNGGLHCIQM  
DTSVNAANQVVSVGADIAFDADPKFFACLVRFESSIONPTTLPTAYDVYPLNGRHDGGY  
YTVKDCVTIDVLP

>Fragment\_G2

PGNNVYVGFMVWSNFTATKCRGLVSLNQVIKEIICLQPLK



**Figure B 1: Introduction of nicking sites into shuttle vectors containing viral genes.**  
Verification of the introduction of the BbvCI nicking site into the shuttle vector containing the viral genes is shown by the generation of ssDNA as in the NSM protocol. Samples were run on a 1% agarose gel with SYBR™ Safe DNA gel stain (Invitrogen) added before casting, the ladder used is the 1 kb DNA ladder from GoldBio.



**Figure B 2: F1 heatmap of counts.**

Counts

**Figure B 3: F2 tile 1 heatmap of counts.**

**Figure B 4: F2 tile 2 heatmap of counts.**



ITL 1 heatmap of counts																																																																																																																																																																																											
Location	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305																																																																																																																																
Mutations	S	G	T	D	V	G	T	D	Q	S	L	G	F	G	S	R	G	V	Q	I	K	H	S	V	P	R	F	V	P	E	H	G	T	M	F	T	L	A	V	R	F	P	E	A	T	K	E	Y	L	N	A																																																																																																																																								
F	135	161	170	171	170	169	167	166	165	165	164	163	162	161	160	159	158	157	156	155	154	154	153	152	151	150	149	148	147	146	145	144	143	142	141	140	139	138	137	136	135	134	133	132	131	130	129	128	127	126	125	124	123	122	121	120	119	118	117	116	115	114	113	112	111	110	109	108	107	106	105	104	103	102	101	100	99	98	97	96	95	94	93	92	91	90	89	88	87	86	85	84	83	82	81	80	79	78	77	76	75	74	73	72	71	70	69	68	67	66	65	64	63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0											
F	135	161	170	171	170	169	167	166	165	165	164	163	162	160	159	158	157	156	155	154	153	152	151	150	149	148	147	146	145	144	143	142	141	140	139	138	137	136	135	134	133	132	131	130	129	128	127	126	125	124	123	122	121	120	119	118	117	116	115	114	113	112	111	110	109	108	107	106	105	104	103	102	101	100	99	98	97	96	95	94	93	92	91	90	89	88	87	86	85	84	83	82	81	80	79	78	77	76	75	74	73	72	71	70	69	68	67	66	65	64	63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0													
W	107	204	141	88	79	68	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																																																																																																																									
W	163	215	174	122	98	75	69	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																																																																																																																									
P	161	258	81	105	97	94	81	78	74	73	72	71	70	69	68	67	66	65	64	63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																																																																																																								
M	125	243	111	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305									
M	140	258	138	86	92	90	88	86	84	82	80	78	76	74	72	70	68	66	64	62	60	58	56	54	52	50	48	46	44	42	40	38	36	34	32	30	28	26	24	22	20	18	16	14	12	10	8	6	4	2	0	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305																																																																												
V	114	257	102	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305
V	114	257	102	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305
S	140	258	138	86	92	90	88	86	84	82	80	78	76	74	72	70	68	66	64	62	60	58	56	54	52	50	48	46	44	42	40	38	36	34	32	30	28	26	24	22	20	18	16	14	12	10	8	6	4	2	0	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305																																																																												
H	174	264	170	93	91	89	87	85	83	82	80	78	76	74	72	70	68	66	64	62	60	58	56	54	52	50	48	46	44	42	40	38	36	34	32	30	28	26	24	22	20	18	16	14	12	10	8	6	4																																																																																																																																										

Mutation	K	G	A	L	T	Y	T	D	G	P	V	Y	G	N	I	P	R	E	I	S	M	K	F	D	B			
* <sup>1</sup>	104	101	102	115	113	129	102	200	192	90	47	140	98	87	111	336	177	74	219	70	207	100	69	102	72	21		
F	105	78	129	168	77	101	99	56	150	193	85	79	29	78	68	61	303	61	102	26	116	111	77	240	65	101	104	
Y	110	78	129	101	85	105	120	169	202	208	214	38	21	324	75	73	130	177	179	93	94	104	108	322	77	125	100	133
P	64	78	121	136	84	98	152	74	145	191	163	88	35	160	50	112	87	297	—	136	99	54	369	94	104	135	58	
M	116	69	69	101	77	96	104	124	102	95	49	111	38	88	130	151	160	75	129	125	238	106	88	114	112	205		
I	114	91	106	122	135	135	65	100	92	259	210	116	55	151	84	51	164	188	124	65	118	87	227	284	61	84	174	
L	90	68	125	122	122	122	102	125	150	101	31	111	79	100	107	701	233	190	166	270	216	112	160	129	176	114		
A	73	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123		
A	73	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123	123		
G	113	—	161	88	60	80	146	122	122	122	122	69	33	130	57	—	121	134	143	57	110	213	62	126	205	106		
C	6	160	147	140	93	77	74	84	137	239	298	63	19	44	159	162	156	125	265	55	561	104	91	338	62	149	226	
S	104	94	196	180	153	74	181	181	181	370	370	104	26	54	267	161	186	187	995	398	254	268	82	84	86	113	129	
S	65	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67	67		
N	179	70	117	96	79	208	205	210	210	210	210	210	210	210	210	210	210	210	210	210	210	210	210	210	210	210		
O	70	80	84	149	71	62	115	157	111	224	166	113	56	155	48	101	261	298	208	105	123	260	120	72	122	139		
D	91	164	151	149	91	154	93	151	121	338	288	—	55	101	167	535	152	298	207	83	103	160	94	251	97	34	701	
E	198	124	105	120	80	70	186	169	147	247	163	21	128	38	140	281	164	161	49	67	215	37	111	268	185	104	212	
E	112	79	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142		
E	112	79	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142	142		
R	131	74	89	111	91	113	63	103	103	203	203	93	23	51	152	20	18	114	178	177	152	177	183	191	197	110	105	

## Counts

0    50    100    250    500

**Figure B 5: F3 tile 1 heatmap of counts.**

**Figure B 6: F3 tile 2 heatmap of counts.**

Counts

**Figure B 7: G1 tile 1 heatmap of counts.**

Counts																			
	0	50	100	250	500														
I <sub>2</sub> (2) heatmap of counts																			
Location	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133
Mutation	R	H	D	G	S	Y	T	V	K	K	D	C	V	T	I	D	V	L	P
*	19	123	96	200	139	334	349	372	223	98	437	218	75	75	315	193	266	231	
F	5	127	65	122	207	171	271	299	378	162	70	331	256	95	125	190	164	384	223
S	21	126	111	117	120	121	122	123	124	125	126	127	128	129	130	131	132	133	307
Y	10	226	131	68	154	160	244	310	183	135	636	244	240	121	250	179	262	232	
P	50	123	107	55	209	131	283	368	373	217	86	344	243	152	49	216	168	372	-
M	22	99	72	68	214	127	270	401	314	235	105	332	364	119	110	194	193	235	222
I	14	102	92	72	164	120	265	526	469	172	64	353	239	140	100	279	275	327	223
N	38	121	111	101	120	121	122	123	124	125	126	127	128	129	130	131	132	133	276
V	21	86	121	133	266	105	304	375	201	145	100	300	140	100	100	100	100	100	193
A	27	98	108	81	280	145	295	523	516	179	130	321	337	269	92	344	317	205	274
G	47	121	111	-	146	413	389	651	184	150	327	271	85	85	862	235	232	239	
C	265	121	72	59	237	203	369	416	439	234	113	113	113	113	97	306	196	279	218
H	113	121	111	101	120	121	122	123	124	125	126	127	128	129	130	131	132	133	252
F	52	154	103	63	168	131	158	182	102	234	107	410	271	71	148	332	239	292	310
N	32	182	118	83	183	131	334	440	417	247	162	302	231	185	135	376	232	287	264
S	34	149	83	89	182	145	355	369	350	229	104	304	293	145	72	301	182	249	262
D	14	146	83	82	122	807	154	334	425	426	207	-	340	255	108	85	301	308	234
E	8	114	83	82	122	807	154	334	425	426	207	-	340	255	108	85	301	308	234
H	144	133	68	102	215	363	363	371	231	92	92	92	92	92	274	187	232	297	
K	13	126	111	87	175	107	253	363	372	-	99	320	271	135	105	277	223	265	260
R	-	251	55	335	183	299	449	392	251	117	421	181	118	63	339	171	461	238	

**Figure B 8: G1 tile 2 heatmap of counts.**

G2 heatmap of counts		Counts																																						
Location	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
Mutation	P	G	N	N	V	Y	V	G	F	M	V	W	S	N	F	T	A	T	K	C	R	G	L	V	S	L	N	Q	V	I	K	E	I	C	L	Q	P	E	K	
*	880	725	785	330	686	1018	769	660	607	348	82	241	504	163	839	422	599	645	651	706	390	844	360	1176	486	457	879	756	980	829	1454	576	458	232	172	174	240	164	252	522
F	623	548	782	243	585	1000	767	771	-	372	147	182	477	156	-	378	444	713	421	607	311	627	311	1175	464	304	871	573	987	842	1018	448	683	525	154	183	112	72	293	386
W	724	708	612	276	788	934	697	1100	592	508	86	-	362	135	890	391	579	700	519	713	281	645	327	841	601	445	850	908	845	691	1128	1060	383	215	524	135	278	135	203	467
Y	733	604	895	380	720	-	793	666	554	425	109	215	493	175	984	439	563	679	551	665	385	665	322	1140	537	309	1121	661	1050	757	1218	409	536	426	202	140	120	152	261	367
P	-	757	689	277	776	882	822	950	610	513	92	293	568	184	964	496	543	853	529	564	329	638	378	1078	715	702	936	1093	1073	751	1263	1015	442	256	118	326	361	-	292	520
M	878	700	894	250	765	946	743	865	544	-	107	259	400	184	1001	390	525	823	627	692	352	703	365	1090	655	585	1110	899	1094	1022	1175	1039	834	566	164	213	291	158	255	527
I	716	648	825	454	722	970	830	753	550	660	158	181	405	142	1072	544	594	594	554	561	344	658	278	1217	446	388	1147	663	1163	-	1406	468	-	-	146	248	122	66	323	485
L	910	771	723	358	682	962	752	839	642	668	101	243	466	155	1009	529	538	751	467	609	352	789	-	1143	704	-	860	1192	960	834	1252	985	488	222	136	-	511	218	-	515
V	776	664	743	243	-	803	-	1196	622	700	1	311	360	149	870	425	737	709	570	660	389	795	411	-	568	586	963	735	-	841	1177	1605	511	345	167	195	273	122	315	526
A	988	763	740	316	956	931	838	1292	658	611	142	281	428	175	1054	547	-	846	537	676	331	666	270	1152	740	494	983	785	1136	911	1265	1648	482	266	201	160	267	177	201	538
G	614	-	655	249	762	915	790	-	673	518	128	265	331	211	880	473	674	690	629	757	386	-	320	1132	550	369	901	881	1182	822	1265	1245	472	184	223	241	194	116	255	464
C	839	1133	792	346	799	1032	797	944	673	550	106	337	451	217	942	517	573	794	546	-	513	649	396	1238	538	421	932	835	1067	903	1258	841	495	281	-	251	221	160	251	474
S	811	1001	884	352	792	914	793	878	758	486	93	243	-	260	926	613	916	779	543	760	492	865	319	1324	-	362	1133	723	1140	931	1540	899	513	348	190	275	291	197	270	542
T	1015	892	882	1487	749	1093	794	854	699	469	74	208	530	248	1051	-	581	-	521	637	433	745	314	1234	532	297	1139	813	96	755	1333	609	474	249	146	197	144	166	256	472
N	935	809	-	-	747	1160	841	816	683	384	96	133	480	-	1006	548	551	783	580	566	383	775	295	1170	522	334	-	931	1077	845	1661	603	501	285	139	252	136	106	249	474
Q	1125	808	308	320	818	872	739	847	534	567	83	265	481	167	892	512	532	785	549	618	347	722	322	1086	693	684	1021	-	1061	831	1292	1229	381	258	136	301	-	233	309	494
D	691	988	907	368	768	1105	958	1141	665	479	169	230	369	196	953	428	775	767	620	617	344	719	272	1236	559	328	1076	796	1233	865	1393	1195	568	411	174	160	186	104	286	460
E	713	701	767	326	940	1009	836	1010	613	427	108	196	436	183	944	421	647	781	598	721	431	911	286	1191	514	330	1025	741	1072	878	1544	-	479	334	181	201	180	174	224	480
H	905	718	776	393	795	995	866	815	586	479	130	161	481	176	1060	400	530	710	535	569	478	786	252	1198	445	520	1154	1023	1065	784	1097	477	621	316	111	292	419	179	368	413
K	783	786	913	338	750	941	749	799	630	405	96	215	400	218	866	440	577	670	-	587	416	872	366	1093	511	283	1140	752	917	832	-	572	487	267	115	213	176	137	255	-
R	848	907	842	380	809	883	807	859	668	459	116	216	502	223	969	540	655	783	501	813	-	838	277	1076	641	546	947	1033	1125	930	1429	839	508	260	213	504	353	200	388	498

Counts
0
50
100
250
500

**Figure B 9: G2 heatmap of counts.**

**Table B 1: Mutant library preparation summary.** Summary table of the transformants required for sufficient library coverage, transformants obtained during comprehensive mutant library preparation by nicking scanning mutagenesis, and the fold excess of the number of transformants required for coverage.

Gene mutated	F1	F2	F3	G1	G2
Number of Residues	95	144	182	132	40
Transformants obtained following NSM	620,000	890,000	370,000	620,000	640,000
Required transformants for 99.9% coverage of possible library	13781	20889	26401	19148	5803
Fold excess over amount required for coverage	45	43	14	32	110

**Table B 2: Primers for incorporating BbvCI nicking sites into the pCR2.1-topo shuttle vector.** Blue text is the overlap region with the shuttle vector, red is the KpnI site, green is the BbvCI site.

Nick incorporation fwd
GCTCTACG <b>GGTACC</b> GCTGAGGGAGCTCGGATCCACTAGTAACG
Nick incorporation rvs
GCTCTAACG <b>GGTACCAAG</b> CTGGCGTAATCATGGTC

**Table B 3: Inner and outer primers for PCR reactions for Illumina sequencing.** Red indicates overhang regions for attaching Illumina adapter primers (inner PCR primers) or overhangs for attaching to inner PCR product (outer PCR primers), black is the overlap region in the gene or the barcode, blue is the Illumina adapter.

Inner PCR Primers	Sequence (5' to 3')
Fragment F1 Fwd	gttcagagtctacagtccga <ins>cc</ins> c <ins>tactt</ins> aggataaaatt
Fragment F2 Tile 1 Fwd	gttcagagtctacagtccgac <ins>atc</ins> gactca <ins>tatagg</ins> gc <ins>aa</ins>
Fragment F2 Tile 2 Fwd	gttcagagtctacagtccgac <ins>atc</ins> gacttaatcaagat <ins>gtat</ins> gct
Fragment F3 Tile 1 Fwd	gttcagagtctacagtccgac <ins>atc</ins> cg <ins>tctct</ins> ctggca
Fragment F3 Tile 2 Fwd	gttcagagtctacagtccgac <ins>atc</ins> gaaggat <ins>tttcc</ins> gt
Fragment G1 Tile 1 Fwd	gttcagagtctacagtccgac <ins>atc</ins> gaattggcc <ins>ctct</ins> tag
Fragment G1 Tile 2 Fwd	gttcagagtctacagtccgac <ins>atc</ins> ggttaat <ins>gttgt</ins> aatgg
Fragment G2 Fwd	gttcagagtctacagtccgac <ins>atc</ins> ggcc <ins>cctct</ins> tagatgca
Fragment F1 Rvs	<ins>ccttggcacccgagaattcca</ins> ttcg <ins>tctcac</ins> agt <ins>cgg</ins>
Fragment F2 Tile 1 Rvs	<ins>ccttggcacccgagaattcca</ins> caac <ins>cgaaaccataacg</ins>
Fragment F2 Tile 2 Rvs	<ins>ccttggcacccgagaattcca</ins> ttcg <ins>tctccc</ins> agag
Fragment F3 Tile 1 Rvs	<ins>ccttggcacccgagaattcca</ins> tct <ins>tagacgaat</ins> caccaga
Fragment F3 Tile 2 Rvs	<ins>ccttggcacccgagaattcca</ins> tc <ins>acactcaat</ins> ctttatca
Fragment G1 Tile 1 Rvs	<ins>ccttggcacccgagaattcca</ins> tg <ins>caatgaagaaaacca</ins>
Fragment G1 Tile 2 Rvs	<ins>ccttggcacccgagaattcca</ins> tt <ins>cgtctccgtac</ins>
Fragment G2 Rvs	<ins>ccttggcacccgagaattcca</ins> tt <ins>gaccgcctcca</ins>
<b>Illumina outer primer adapter</b>	aat <ins>gatacggcgaccaccgagat</ins> ctac <ins>ac</ins> gttcagagtctacagtccga
<b>Illumina outer PCR adapters and barcodes</b>	
RPI31 (Fragment F1)	caagcagaagacggcatac <ins>gagat</ins> ATCGTG <ins>gtgactggagttccttggcacccg</ins> <ins>agaattcca</ins>
RPI15 (Fragment F2 Tile 1)	caagcagaagacggcatac <ins>gagat</ins> TGACAT <ins>gtgactggagttccttggcacccg</ins> <ins>agaattcca</ins>
RPI16 (Fragment F2 Tile 2)	caagcagaagacggcatac <ins>gagat</ins> GGACGG <ins>gtgactggagttccttggcaccc</ins> <ins>gagaattcca</ins>
RPI17 (Fragment F3 Tile 1)	caagcagaagacggcatac <ins>gagat</ins> CTCTAC <ins>gtgactggagttccttggcacccg</ins> <ins>agaattcca</ins>
RPI18 (Fragment F3 Tile 2)	caagcagaagacggcatac <ins>gagat</ins> GC <ins>GGAC</ins> <ins>gtgactggagttccttggcaccc</ins> <ins>gagaattcca</ins>
RPI19 (Fragment G1 Tile 1)	caagcagaagacggcatac <ins>gagat</ins> TTTCAC <ins>gtgactggagttccttggcacccg</ins> <ins>agaattcca</ins>
RPI20 (Fragment G1 Tile 2)	caagcagaagacggcatac <ins>gagat</ins> GGCCAC <ins>gtgactggagttccttggcaccc</ins> <ins>gagaattcca</ins>
RPI21 (Fragment G2)	caagcagaagacggcatac <ins>gagat</ins> CGAAAC <ins>gtgactggagttccttggcaccc</ins> <ins>gagaattcca</ins>

## **APPENDIX C**

**Purification of the TROP2 extracellular domain from a stable insect cell line using  
ammonium sulfate precipitation**

## **Abstract**

The tumor associated calcium signal transducer 2 (TROP2) is an oncogenic transmembrane protein that is overexpressed in aggressive and late stage cancers. Here we present a facile method for purifying approximately 3 mg/L quantities of the extracellular domain of TROP2 from a stable *Drosophila* Schneider's (S2) cell line. The secreted oncogene is purified by ammonium sulfate fractionation, immobilized metal affinity chromatography, and size exclusion chromatography. The folding fidelity of the highly purified product was confirmed with a binding assay that tests for a unique conformational epitope in TROP2.

## **Introduction**

The tumor associated calcium signal transducer 2 (TROP2) is a 36 kDa type-I transmembrane glycoprotein associated with late stage and aggressively metastatic tumors<sup>1-3</sup>. TROP2 is basally expressed in human trophoblasts when the fetal tissue joins into the maternal circulation<sup>4</sup>. Ectopic expression of TROP2 in cancer cells results in increased proliferation and tumorigenicity<sup>5-8</sup>, while gene silencing of TROP2 decreases these same characteristics<sup>5,7,9</sup>. Following an activation event that is currently unknown, TROP2 undergoes regulated intramembrane proteolysis<sup>10</sup> at two distinct sites to release the extracellular domain (TROP2Ex, residues 27-274) and the intracellular domain (TROP2Ic, residues 298-323). TROP2Ic migrates to the nucleus where it increases the expression of genes associated with cell proliferation<sup>10-11</sup>. The liberated extracellular domain interacts with signaling proteins, though with which ones and how is incompletely understood<sup>8,12-13</sup>.

Previous biophysical analysis has found that TROP2Ex can assume a dimeric conformation similar to its paralog epithelial cell adhesion molecule (EpCAM)<sup>14-16</sup>. Interestingly, the glycosylation state of the protein seems to dictate whether the monomeric (favored when fully glycosylated) or dimeric (favored when unglycosylated) form predominates at equilibrium<sup>15</sup>. Previous work studying the migration of metastatic cancer cells found that the oligomeric state of TROP2 might govern its function<sup>9</sup>. Whereas binding of TROP2 with the bivalent m7E6 monoclonal antibody (mAb)<sup>17</sup> resulted in significant inhibition of migration of MDA-MB-231 cancer cells while not impacting toxicity or cell proliferation, binding with the monovalent m7E6 Fab did not significantly alter migration. These results suggest that the dimer to monomer transition of TROP2 may play a role in its activation and mechanism of action.

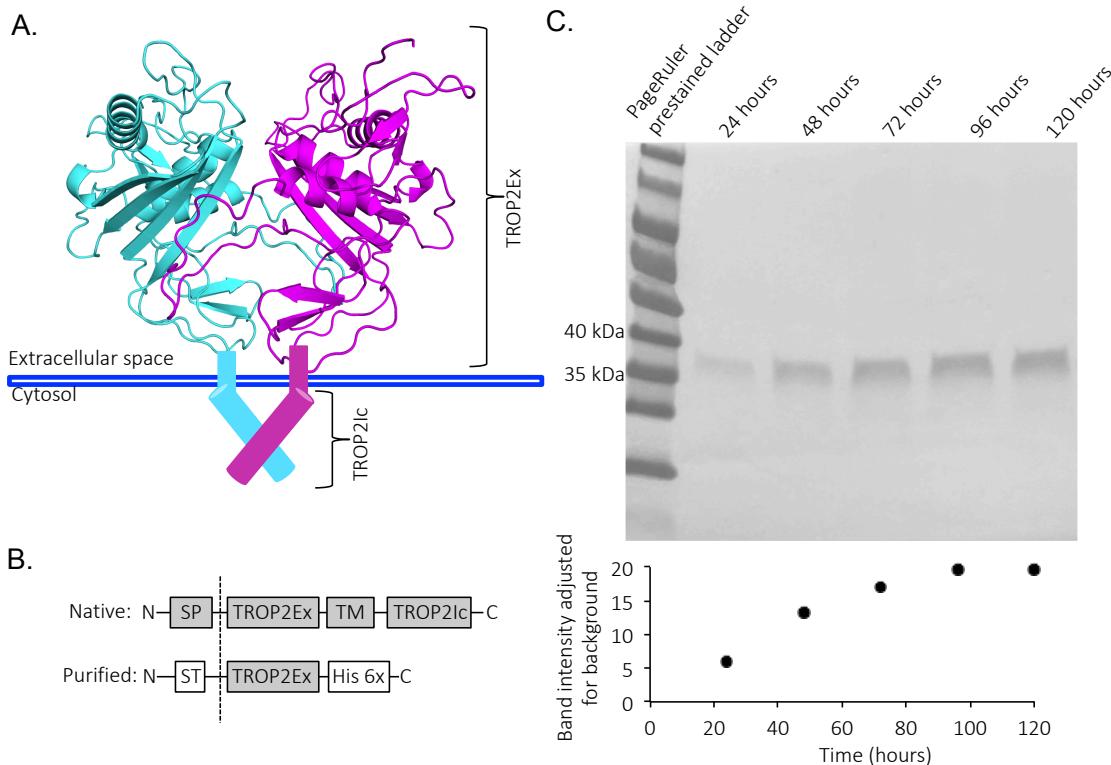
Here we provide a novel and facile approach for the purification of properly folded TROP2Ex from insect cells. We used a stably transformed *Drosophila* Schnieder's (S2) cell line to secrete the protein. The secreted protein is removed from the culture media with a classical but underutilized ammonium sulfate fractionation method. Next, the protein is further purified with immobilized metal affinity chromatography and size exclusion chromatography. We obtain yields above 3 mg/L of induction culture with excellent purity. The correct folding of the protein has been verified with flow cytometry binding assays using a single chain variable fragment (scFv) of the m7E6 mAb displayed on the yeast surface, which recognizes a conformational epitope in TROP2Ex. The purified TROP2Ex can be used in a variety of *in vivo* and *in vitro* experiments to help understand how TROP2 is activated. Additionally, the yeast displayed m7E6 scFv system presented here can be used to map the paratope involved in the TROP2 m7E6 binding interaction.

## Results and discussion

### *Expression of TROP2Ex from a stable cell line*

As TROP2Ex contains six potential disulfide bonds and multiple N-linked glycosylation sites, we selected *Drosophila* Schneider 2 (S2) cells for our expression host

because of their potential ability to accurately maintain these post-translational modifications, their low house-keeping requirements, ease of transfection, the advantages of protein secretion, and the ease of forming stably transfected cell lines. The TROP2 gene was truncated to residues



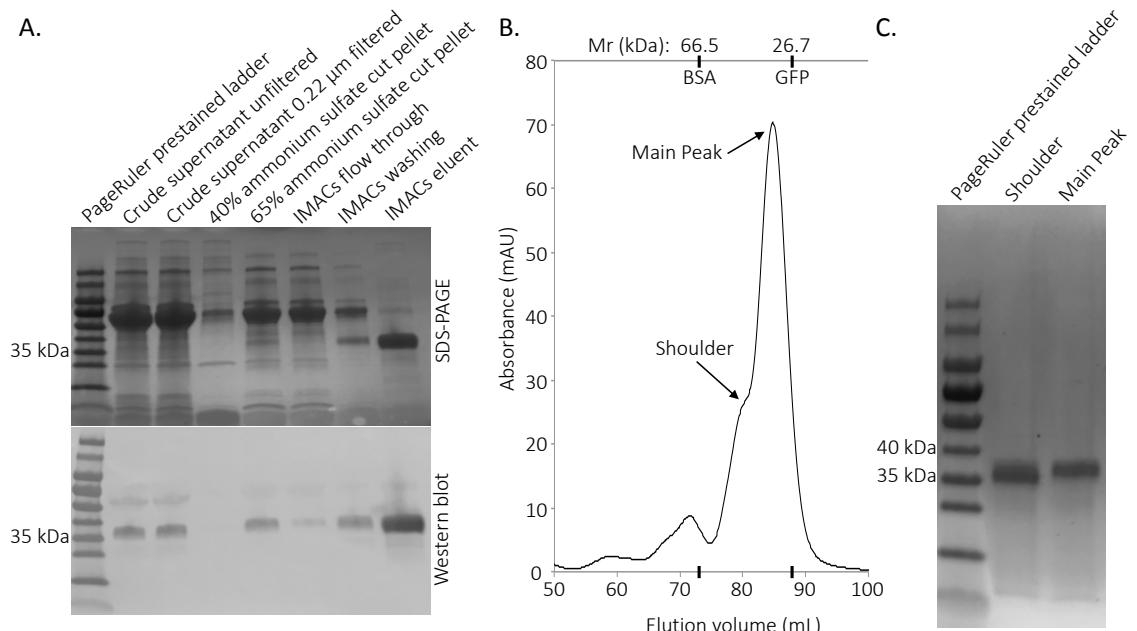
**Figure C 1: TROP2Ex and its expression.** A. Homology model of the dimeric TROP2Ex domains generated using the paralog EPCAM<sup>14,16</sup> (pdb: 4MZV). This model was previously published<sup>9</sup>. B. Diagram of the native and purified form of TROP2, with the dashed line representing the truncation point in the secretory pathway, SP = signal peptide, ST = secretion tag, TM = transmembrane domain, TROP2Ic = TROP2 intracellular domain. C. Western blot testing for TROP2Ex in equivalent volumes of the stable S2 cell supernatant following initiation of induction and band intensity adjusted for background is plotted against time; bands were quantified with ImageJ by standard protocols.

27-254 (TROP2Ex) for expression and secretion in S2 cells (**Figure C 1A-B**). The native secretion signal peptide sequence (SP) was replaced with the secretion tag needed for expression in S2 cells (ST) (**Figure C 1B**). A His 6x tag was added to the C-terminal end of the truncated TROP2Ex gene (**Figure C 1B**) to facilitate purification.

A stable S2 cell line that secretes TROP2Ex was successfully generated as described in the **Materials and Methods** (**Figure C 1C**). The secreted TROP2Ex runs larger than the predicted molecular weight - ~28.8 kDa – because N-linked glycosylations can occur at up four distinct sites after translation<sup>15</sup>. The robustness of the cell line was verified through multiple successful TROP2Ex purifications from the same lineage over the course of months (data not shown). After generating the stable S2 cell line we next increased the culture volumes possible by transitioning from static to shaking cultures. The addition of a non-ionic surfactant the cells maintained healthy morphologies in the shaking cultures as assessed by confocal microscopy. The optimal induction time for secreted TROP2Ex titers was performed the shaking cultures, with a 96-hour induction period striking the best balance between titers and volumetric productivity (**Figure C 1C**). The shaking cultures were successfully scaled to 300 mL and induced for 96 hours using copper.

#### *Purification of TROP2Ex*

TROP2Ex was purified using ammonium sulfate fractionation, IMACs, and chromatography. Fractionation with ammonium sulfate at 40% saturation, followed by fractionation at 65% saturation efficiently removed some of the contaminating proteins and contaminant copper and allowed for resuspension of the secreted TROP2Ex in IMAC buffer (**Figure C 2A**). Batch IMAC resulted in nearly pure TROP2Ex (**Figure C 2A**) and size exclusion chromatography (SEC) using a Superdex 200 column equilibrated with 1x PBS at 4°C removed the remaining impurities (**Figure C 2B-C**). TROP2Ex eluted from SEC as a peak with a shoulder between 78-90 mL elution volume (**Figure C 2B-C**), suggesting that both monomeric and dimeric TROP2Ex exist under elution conditions.

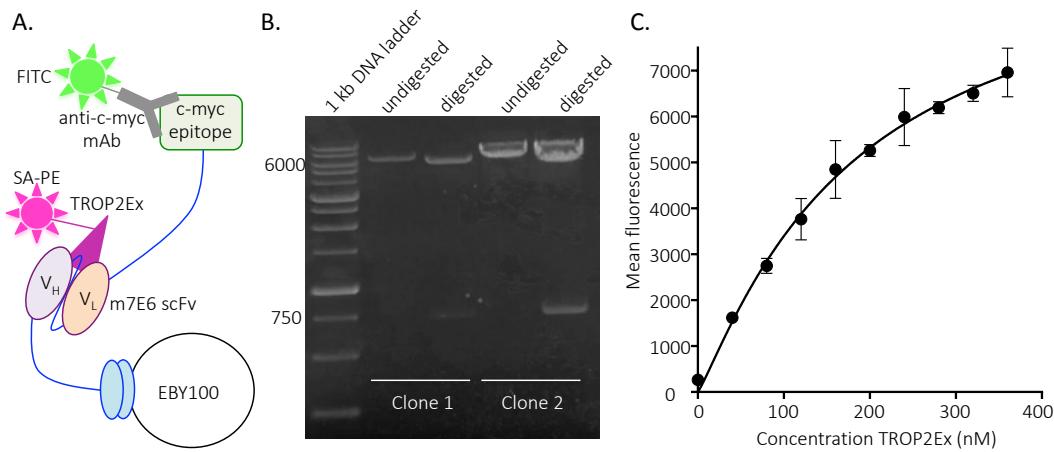


**Figure C 2: Purification of TROP2Ex from insect cell cultures.** A. Images of an SDS-PAGE gel and a Western blot containing identical samples resulting from the purification of TROP2Ex. B. FPLC-SEC chromatogram of the IMACs eluent, column was equilibrated with 1x PBS at 4°C, main peak and respective shoulder at ~78 – 90 mL are TROP2Ex. C. SDS-PAGE image of the shoulder (80 mL elution volume) and main peak (85 mL elution volume).

Interestingly, when the purified TROP2Ex main peak is run again with the same SEC conditions as before both the main peak and shoulder are resolved (data not shown). All subsequent analysis corresponds only to TROP2Ex from the main peak (**Figure C 2C**). Under conditions used in these experiments TROP2Ex (unglycosylated mass = 28.8 kDa) elutes at ~85 mL (**Figure C 2B**), while BSA (66.5 kDa) elutes at ~73 mL, and GFP (26.7 kDa) elutes at ~88 mL off of the Superdex 200 column in identical conditions. This elution pattern indicates that our purified TROP2Ex is predominantly monomeric. This finding is consistent with previous results finding that only a small fraction of the glycosylated TROP2Ex is in the dimeric form when at equilibrium<sup>15</sup>. In total, we obtain yields of 3.6 mg/L of purified monomeric TROP2Ex from the stable cell line.

### Validation of TROP2Ex folding

To verify that TROP2Ex was in a properly folded conformation, we assessed its binding to a yeast displayed scFv derived from the high affinity anti-TROP2 mAb m7E6<sup>17</sup> (**Figure C 3A**). m7E6 recognizes a conformational epitope



**Figure C 3: Verification of the proper assembly of the m7E6 scFv yeast display construct and of binding to TROP2Ex.** A. A diagram of the yeast display system, the binding interactions detected, and the respective conjugate fluorophores: SA-PE = streptavidin R-phycoerythrin conjugate. B. A DNA gel containing restriction enzyme analysis products for verifying the assembly of the m7E6 scFv pETconNK construct via homologous recombination. Following digestion the DNA was run on a 1% agarose gel prestained with ethidium bromide at 120 V for 50 minutes, the 1 kb DNA ladder is from GoldBio. Following electrophoresis the gel was imaged using a Safe Imager™ Blue Light Transilluminator. C. Clonal titration of yeast displayed m7E6 scFv with TROP2Ex, n = 3, error bars = 1 S.D.

located in the membrane distal region of TROP2Ex<sup>9</sup>. Overhang PCR and homologous recombination in yeast were used to assemble to the anti-TROP2 m7E6 scFv yeast display construct as described in the **Materials and Methods** section (**Figure C 3B**). Clonal titrations resulted in a binding curve that spanned a >25-fold change in signal intensity and provide a preliminary dissociation constant ( $K_D$ ) of  $180 \pm 54$  nM (**Figure C 3C**). There are multiple possible reasons the titrations didn't reach saturation: we didn't expand the titrations to high enough concentrations, the yeast displayed form of m7E6 might not be stable, or perhaps the

binding reaction conditions are not fully optimized for this pairing. The obtained  $K_D$  is ~4 times larger than the ~44 nM  $K_D$  reported for the binding of this pairing where TROP2Ex is displayed on the surface of yeast and recombinant m7E6 fragments of antibodies (fabs) are applied<sup>9</sup>. The large difference in the  $K_D$  could be the result of the titrations failing to reach saturation, different glycosylation patterns on either protein, differences introduced into the CDRs or framework of the antibody fragments resulting from the truncation of the framework region, or it is possible that the orientation of the  $V_H$  and  $V_L$  domains of the scFv aren't optimal. This pairing binds with a dissociation constant in the hundreds of nanomolar, indicating that the purified TROP2Ex contains the same conformational epitope as the native TROP2Ex and is likely correctly folded.

## Conclusion

Here we presented a facile method for purifying the extracellular domain of TROP2 from *Drosophila* Schneider's (S2) insect cells. The cell lines generated are stable and can be seeded for TROP2Ex production once recovered from storage. Removing the need for retransformation allows for an easy scale up of the induction cultures. Additionally, we have presented a method for purifying TROP2Ex using ammonium sulfate fractionation, IMACs, and chromatography. This method allows for the purification of >3 mg/L of properly folded TROP2Ex from an induction culture in a single day.

## Materials and methods:

### *Materials*

Antibiotics were purchased from GoldBio and ThermoFisher. All enzymes were from New England Biolabs (NEB), and all other chemicals were purchased from Sigma-Aldrich. The

TROP2Ex gene block, and all primers were purchased from Integrated DNA technologies (IDT); the m7E6 full IgG gene which was purchased from GenScript. Flow cytometry probes were purchased from Miltenyi Biotec (anti-c-myc FITC) and ThermoFisher (streptavidin R-phycoerythrin conjugate). Other materials used have their respective manufacturers listed in the text.

### *Plasmid Construction*

The gene for TROP2Ex (Kozac-Bip-TROP2Ex) was designed as follows following a 5' to 3' convention: a KpnI restriction site, followed by a Kozac sequence and a Bip secretion tag (**Supporting text, Note S1**), then the TROP2Ex gene (residues 27-274), with an AgeI restriction site. Kozac-Bip-TROP2Ex was codon optimized for expression in *D. melanogaster* and was ordered as a single gene block from IDT (**Supporting text, Note S1 and S3**). The Kozac-Bip-TROP2Ex gene was subcloned into the pMT/V5-His B plasmid (Invitrogen) using KpnI and AgeI restriction enzymes and standard cloning methods to produce pMT-Bip-TROP2Ex containing an in-frame C-terminal 6x His tag. The final construct was transformed into *E. coli* XL-1 Blue and a glycerol cell stock was generated. A scratch of the freezer stock was used to inoculate 50 mL LB + 50  $\mu$  g/mL carbenicillin which was grown overnight at 37°C. DNA was Midi-prepped (Qiagen) from 25 mL of the overnight culture, and the purified DNA was resuspended in nuclease free H<sub>2</sub>O and frozen at -20°C.

The m7E6 scFv<sup>17</sup> (**Supporting text, Note S2 and S4**) comprises (5' to 3') V<sub>H</sub> residues 1-118, a flexible (GGGGS)<sub>3</sub> linker, and V<sub>L</sub> residues 1-111. Using the full m7E6 IgG gene that was previously purchased from Genscript<sup>17</sup> as template, overhang PCR was used to assemble the scFv, introduce the linker, and introduce homology regions from the pETconNK yeast display

plasmid<sup>18</sup>. Primers used for assembly are in **Note S5** of the supporting text. The PCR products were cleaned and concentrated using a Monarch® PCR & DNA Cleanup Kit (NEB) and eluted into nuclease free H<sub>2</sub>O. The scFv gene was assembled with pETconNK vector cut with NdeI and XhoI using homologous recombination into *S. cerevisiae* EBY100 cells. For homologous recombination the V<sub>H</sub> and V<sub>L</sub> inserts were combined with the cut pETconNK plasmid at a ratio of 30 fmol : 1 fmol (respective insert : cut vector) in a final volume of 8 µL. All of the DNA mixture was transformed into EBY100 chemically competent yeast using standard protocols. Transformation mixtures were plated on agar plates containing Synthetic Complete media supplemented with 2% (w/v) dextrose (SDCAA agar) and grown at 30°C for 2-3 days. Colonies were picked and grown in 1 mL of SDCAA and 1x penicillin/streptomycin (SDCAA media) at 30°C with shaking at 250x rpm for ~16 hours. The 1 mL cultures were then used to inoculate 50 mL of SDCAA media that was grown for generation of freezer stocks essentially as described in Klesmith et al.<sup>18</sup>. Both plasmids used in this work are available from Addgene: m7E6 scFv pETConNK (ID: 125731), and pMT-Bip-TROP2Ex (ID: 125730).

#### *Transformation of S2 Drosophila Cells and the Generation of the Stable Cell Line*

The S2 *Drosophila* cells were cultured, transformed, and a stable cell line was generated essentially as described by the manufacturer (Invitrogen). Briefly, 3 mL/well of Schneider's S2 *Drosophila* Media (ThermoFisher) supplemented with 10% heat-inactivated fetal bovine serum (full S2 media) was seeded with 1x10<sup>6</sup> cells/mL in a 6-well plate (Fisher Scientific). Cells were incubated at 28°C until a density of 2-4x10<sup>6</sup> cells/mL was reached, after approximately 16 hours. After reaching density the cells were transformed with 19 µg pMT-Bip-TROP2Ex plasmid and 1 µg pCoPuro plasmid<sup>19</sup> via standard calcium phosphate transfection. Calcium phosphate DNA

precipitates were prepared using the Calcium Phosphate Transfection Kit (Invitrogen) as described by the manufacturer. The DNA precipitate solution was added to respective wells such that 5 µg total DNA was applied. Additions were performed in a dropwise fashion with gentle swirling of the culture after each addition. After addition of the DNA precipitates the cultures were incubated at 28°C overnight. The next morning the transformation cultures were spun at 440 xg for 3 minutes and the transfection media was removed via aspiration. Cells were resuspended in 3 mL of full S2 media and spun as before. This media was aspirated and the washing procedure was repeated once more. Cells were resuspended in 3 mL full S2 media, placed back into the original well and incubated at 28°C for 48 hours. Puromycin toxicity was assessed as described in the *Drosophila* S2 Cell Manual (Invitrogen). It was determined that 2 µg/mL puromycin was the optimal concentration for selection of stable S2 cells harboring the pMT-Bip-TROP2Ex and pCoPuro plasmids.

Following the recovery of the transformation cultures, a single culture was brought to a final concentration of inducer - 500 µM CuSO<sub>4</sub> - and incubated for 48 hours at 28°C. After 48 hours the culture was spun at 3,200 xg for 10 min to pellet the cells. The supernatant was filtered through a 0.22 µM PES filter and stored at 4°C. A portion of the aliquot was assessed via Western blotting with standard protocols. Briefly, the sample was denatured in Laemmli's buffer supplemented to 1.5% β-mercaptoethanol at 1x for 10 minutes at 98°C. The denatured protein was run on a 4-20% Mini-PROTEAN® TGXTM precast gels (Bio-Rad) for 1 hour at 120 V, and the PageRuler Prestained Protein Ladder (10-180 kDa, Thermo-Fisher) was used as a ladder. Once separated on the gel the proteins were transferred to a nitrocellulose membrane using the iBlot™ Gel Transfer Device with iBlot™ Transfer Stacks as per the manufacturer instructions (ThermoFisher). The primary antibody was Anti-6x His tag® horseradish peroxidase conjugate

antibody (abcam, ab1187), and the blot was resolved using the Pierce™ DAB Substrate Kit as described by the manufacturer (ThermoFisher).

The generation of the stable cell line was performed in parallel to the initial testing of TROP2Ex production. Recovered cells were spun at 440 xg for 3 minutes the supernatant was aspirated, the pellet was resuspended in 3 mL of full S2 media + 2 µg/mL puromycin and placed back into the original well. Every 3-4 days the culture was pelleted, the supernatant removed via aspiration, gently resuspended in 3 mL of fresh full S2 media with 2 µg/mL puromycin, placed back into the original well and grown at 28°C. This practice was continued until the cells had grown a density of 6-20x10<sup>6</sup> cells/mL. Once at density the cells were then passaged 1:2 into full S2 media + 2 µg/mL puromycin to 6 mL final volume. At this point the culture was moved into a Corning® 25 cm<sup>2</sup> culture flask with a vented cap (Sigma-Aldrich) and grown at 28°C. Once the cells reached a density of 6-20x10<sup>6</sup> cells/mL they were passaged 1:2 into 12 mL full S2 media with puromycin in a Corning® 75 cm<sup>2</sup> culture flask with a vented cap (Sigma-Aldrich) and grown at 28°C. When the cells were passaged into the 75 cm<sup>2</sup> flask a 6 well plate was also seeded with 3 mL full S2 media + 2 µg/mL puromycin with 1x10<sup>6</sup> cells/mL and TROP2Ex expression was tested as previously described. The stable cell line was cultured in the 75 cm<sup>2</sup> flask until the cells grew to 1-2x10<sup>7</sup> cells/mL, at which time freezer stocks were prepared.

To prepare freezer stocks the cultures were spun at 440 xg for 3 minutes and the conditioned media was collected in a 50 mL falcon tube via careful pipetting. Cells were resuspended in 10 mL room temperature sterile 1x PBS (10 mM Na<sub>2</sub>HPO<sub>4</sub> + 2 mM KH<sub>2</sub>PO<sub>4</sub> + 2.7 mM KCl + 137 mM NaCl, pH 7.4) and spun as before, the resulting supernatant was removed via aspiration. Freezing buffer was prepared by combining 9 mL conditioned media with 9 mL full S2 media and 2 mL 100% DMSO. The cells were gently resuspended in the room

temperature freezing buffer at  $1.1 \times 10^7$  cells/mL and 1 mL aliquots were dispensed into cryogenic storage vials. A freezer box was lined with paper towel and the samples were stored therein, this box was placed into a Styrofoam container and this was stored at -80°C. After 24 hours at -80°C the vials were moved into liquid nitrogen storage. The remaining culture was continually expanded via 1:2 dilutions into fresh full S2 media with puromycin as before. Once at 30 mL final volume, cells were passaged 1:10 into full S2 media with puromycin every ~6-7 days.

*Production and Purification of TROP2Ex:*

In a sterile 250 mL Erlenmeyer flask, 50 mL dynamic full S2 media (full S2 media + 2 µg/mL puromycin + 0.1% Pluronic® F68 (Sigma-Aldrich) + 1x penicillin/streptomycin) was seeded with stable S2 cells (pMT-Bip-TROP2Ex-pCoPuro) at  $1 \times 10^6$  cells/mL, this was grown at 28°C with shaking at 110 rpm. This culture was grown to a density of ~ $2-4 \times 10^6$  cells/mL, once there CuSO<sub>4</sub> was added to a final concentration 500 µM. Extractions were performed every 24 hours for 120 hours following initiation of TROP2Ex expression. Collected samples were analyzed via Western blotting as described previously.

For larger scale production 300 mL of dynamic full S2 media was seeded and grown as described previously in a 2 L baffled culture flask. Once at ~ $2-4 \times 10^6$  cells/mL (~24-48 hours) the culture was induced CuSO<sub>4</sub> as before for 96 hours at 28°C with shaking at 110 rpm. After induction the culture was spun at 3000 xg for 30 minutes at 4°C, the supernatant was saved and the pellet discarded. The volume of the supernatant was measured with a graduated cylinder and then filtered through a 0.22 µm PES filter. After filtration, the supernatant was placed into a sterile 1L beaker with a stir bar, placed into an ice bath on a stir plate, and brought to 0.1 mM PMSF. With continuous gentle stirring solid ammonium sulfate was slowly added to the

supernatant until the solution was 40% of the saturation limit for ammonium sulfate at 0°C. The mass of solid ammonium sulfate to add was calculated using a web-based calculator (<http://www.encorbio.com/protocols/AM-SO4.htm>). After addition of ammonium sulfate the beaker was covered and stirred for 3 hours in the ice bath. After the equilibration, the 40% saturated ammonium sulfate supernatant was spun at 10,000 xg for 15 minutes at 4°C. Following centrifugation the supernatant and precipitate were separated, the pellet was discarded, and the supernatant was moved back into the ice bath. With gentle stirring the 40% saturated supernatant had solid ammonium sulfate slowly added until the solution was 65% of the saturation limit for ammonium sulfate at 0°C . The 65% ammonium sulfate saturated solution was covered and stirred for 1 hour in the ice bath. After the second equilibration the solution was spun at 15,000 xg for 15 minutes at 4°C. The supernatant was discarded, and the pellets were resuspend in a total of 40 mL 1x PBS + 20 mM imidazole, and pooled together.

Batch immobilized metal affinity chromatography (IMAC) was performed using a 1 mL bed volume of Talon® resin (Takara Bio) pre-equilibrated with 1x PBS + 20 mM imidazole. A 3-hour adherence reaction was performed by rocking the samples with the resin in an ice bath. The slurry was pelleted at 700 xg for 2 minutes and the supernatant was discarded. The slurry was then resuspended in 10 mL of 1x PBS with 20 mM imidazole with gentle pipetting, transferred to a 15 mL conical tube, and incubated for 5 minutes with rocking in the ice bath. The resin was pelleted as described above and the supernatant was discarded, this wash step was repeated four more times. To elute TROP2Ex the pelleted resin was resuspended with 1 mL of 1x PBS + 400 mM imidazole with gentle pipetting and incubated for 5 minutes with rocking in an ice bath. The slurry was spun as before and the supernatant was collected and stored in ice, this was repeated

five more times. IMACs purified TROP2Ex was concentrated from ~6 mL to ~2.5 mL with Amicon® Ultra-4 10,000 NMWL centrifugal filters and filtered through a 0.22 µm PES filter.

The IMAC eluent was further purified by FPLC-SEC using a HiLoad 16/600 Superdex 200 pg size exclusion column equilibrated with 1x PBS with a mobile phase flow rate of 1 mL/min. The main peak eluted at 85 mL with a shoulder that eluted at 80 mL, and the main peak and shoulder were not pooled. The collected fractions were filter through a 0.22 µm PES filter and stored at 4°C. Successful purification of the protein was confirmed by Western blotting and the purity of the final product was confirmed using SDS-PAGE stained with SimplyBlue SafeStain as per the manufacturer instructions (ThermoFisher). The purified protein was quantified using both Bradford assays with standard protocols, and by the Edelhoch method with a predicted molar extinction coefficient at A<sub>280</sub> of 16680 M<sup>-1</sup> cm<sup>-1</sup>.

#### *Yeast display of the m7E6 scFv and clonal titration*

Yeast surface display of the m7E6 scFv with EBY100 cells was optimized with the optimal displaying conditions discovered to be 30°C for 16-18 hours with shaking at 250 rpm in a 1 mL culture. TROP2Ex was chemically biotinylated using NHS-Ester chemistry and clonal titrations were performed exactly as described by Medina-Cucurella and Whitehead<sup>20</sup>.

#### **Acknowledgements**

We would like to thank Dr. C. Chan for allowing us to use her laboratory infrastructure and Dr. A. Oak for her training and help with insect cell culturing.

## Supporting text

**Note S1. DNA sequence of the TROP2Ex gene block inserted into the pMT vector.** *D. melanogaster* optimized TROP2Ex gene (underlined text), Kozac sequence (blue text), Bip secretion tag (orange text), KpnI and AgeI restriction sites are bolded and italicized.

***GGTAC***ACCATGGGA***ATGAAAGTTATGCATATTACTGGCCGTGCGCCTTGTTGGC***  
***CTCTCGCTCGGGCATACGGCAGCCCAGGATAATTGTACTGTCCAACGAATAAAATG***  
ACGGTATGCTCCCCGGACGGACCAGGCCGTAGGTGTCAATGTGGGCACTCGGAAG  
TGGAATGGCTGTTGACTGTTCCACACTCACGAGTAAATGTTGCTCCTGAAGGCAAG  
GATGAGCGCTCCTAAAAACGCTCGAACACTCTCGTAGGCCAAGTGAGCACCGCCTGG  
TCGACAACGACGGCTTGTACGATCCCATTGCGACCCAGAGGGCAGGTTAAAGCT  
CGCCAGTGCAACCAAACGTCGGTCTGCTGGTGTAAACTCCGTGGGTACGACGG  
ACCGACAAAGGTGATTGTCCTCGCTGTGACGAACCTCGTCCGACGCCATCACATA  
CTCATTGATCTGAGGCATCGTCCCACCGCTGGAGCGTTCAATCACTCGGACTTGGAC  
GCCGAAC TGCGCCGATTGTTCCGAGAACGATATAGGTTGCATCCAAAGTTGTTGCT  
GCGGTACATTACGAGCAACCAACCATAACAGATAGAATTGAGGCAAATACCAAGTCA  
AAAAGCGGCGGGAGATGTAGATATTGGCGATGCAGCCTATTATTCGAAAGGGATA  
TTAAAGGCGAATCGTTGTTCCAAGGTGAGGCGGACTCGACCTCCGGGTTAGGGC  
GAGCCATTGCAGGTCGAAAGGACTCTGATCTATTACCTCGATGAAATACCTCCGAAG  
TTAGCATGAAGAGGTTG***ACCGGT***

**Note S2. DNA sequence of the m7E6 scFv.** The sequence for the V<sub>H</sub> domain is blue, the sequence of the linker is underlined, the V<sub>L</sub> sequence is lower case.

CAGGTCCAGCTGAAGGAAAGCGGCCCCGGCCTGGTGGCCCCTCCCAGTCTCTGAG  
CATCACCTGCACAGTGAGCGGCTTCTCCCTGACCTTACGGCGTGCACGGGTGCG  
GCAGCCTCCTGGCAAGGGCCTGGAGTGCTGGCTGGCGTGTGGACCGGGCTCCA  
CAGATTATAACTCTGCCCTGATGAGCCGGCTGTCCATCAACAAGGACAATTCCAAGT  
CTCAGGTGTTCTGAAGATGAATAGCCTGCAGACCGACGATACAGCCATGTACTATT  
GTGCCCCGGACGGCGATTACGACAGATATACCATGGATTACTGGGGCCAGGGCACC  
**AGCGTGACAGTGAGC**GGTGGAGGGGTTCAAGGCGGGGTGGAAGCGGTGGAGGGG  
GTAGCagtgatatcgctgactcagtcgtccctgttcaactggcgtaggcctggccagagagctacaatctctgcagagcctccaag  
tctgtgaggcacccggctacagctatgcactggtaccagcagaagccagggccagccccctaagctgtatctggcctcaacct  
gaaagccggcgtgcgtcggtctctggcagccgctccggcacagacttaccctgaatattcaccctgaggaggaggatgccgc  
cacatactattccagcactccagggagctgccctacacattccggcggcggcaccaagctggagatcaag

**Note S3. Secreted TROP2Ex amino acid sequence.**

HTAAQDNCTCPNKMTVCSPDGGGRQCRALGSGMAVDCSTLTSKCLLKARMSAP  
KNARTLVRPSEHALVDNDGLYDPDCDPEGRFKARQCNCNTSVCWCVNSVVRRTDKGD  
LSLRCDELVRTHILIDLRHRPTAGAFNHSSDLDAELRRLFRERYRLHPKFVAAVHYEQPT  
IQIELRQNTSQKAAGDVDIGDAAYYFERDIKGESLFQGRGGLDLRVRGEPLQVERTLIYY  
LDEIPPKFSMKRLTGHHHHHH

**Note S4. Amino acid sequence of the m7E6 scFv.**

QVQLKESGPGLVAPSQSLISITCTVSGFSLTSYGVHWVRQPPGKGLEWLGVIWTGGSTDY  
NSALMSRLSINKDNSKSQVFLKMNSLQTDDTAMYCARQGDYDRYTMDYWQGQGTSV  
TVSGGGGSGGGGSGGGSSDIVLTQSPASLA VSLGQRATISCRASKSVSTSGYSYMHWY  
QQKPGQPKLIYLASNLESGVPARFSGSGSGTDFTLNIHPVEEDAATYYCQHSRELPY  
TFGGGTKEIK

**Note S5. Primers for overhang PCR for forming the m7E6 scFv and introduction of the cut sites for cloning into pETconNK.**

V<sub>H</sub> forward:

5'- AGCGGAGGC GGAGGGTCGGCTAGCCATATGCAGGTCCAGCTGAAGGAAAG -3'

V<sub>H</sub> reverse:

5'-CTCCACCGCTTCCACCCCCGCCTGAACCCCCCTCCACCGCTCACTGTCACGCTGGT -  
3'

V<sub>L</sub> forward:

5'-

AGGCGGGGGTGGAAGCGGTGGAGGGGGTAGCAGTGATATCGTGCTGACTCAGTCCC  
-3'

V<sub>L</sub> reverse:

5'- GAAATAAGCTTTGTT CGGATCCGCC CCGCTGAGCTTGATCTCCAGCTGGTGC -  
3'

## **REFERENCES**

## REFERENCES

1. Fong D, Moser P, Krammel C, Gostner JM, Margreiter R, Mitterer M, Gastl G, Spizzo G: **High expression of TROP2 correlates with poor prognosis in pancreatic cancer.** *Brit. J. Cancer* 2008, 99:1290–1295.
2. Wu M, Liu L, Chan C: **Identification of novel targets for breast cancer by exploring gene switches on a genome scale.** *BMC Genomics* 2011, 12:547.
3. McDougall ARA, Tolcos M, Hooper SB, Cole TJ, Wallace MJ: **Trop2: from development to disease.** *Dev. Dynam.* 2015, 244:99–109.
4. Lipinski M, Parks DR, Rouse RV, Herzenberg LA: **Human trophoblast cell-surface antigens defined by monoclonal antibodies.** *Proc. Natl. Acad. Sci. USA* 1981, 78:5147–5150.
5. Wang J, Day R, Dong Y, Weintraub SJ, Michel L: **Identification of Trop-2 as an oncogene and an attractive therapeutic target in colon cancers.** *Mol. Cancer Ther.* 2008, 7:280–285.
6. Cubas R, Zhang S, Li M, Chen C, Yao Q: **Trop2 expression contributes to tumor pathogenesis by activating the ERK MAPK pathway.** *Mol. Cancer* 2010, 9:253.
7. Trerotola M, Cantanelli P, Guerra E, Tripaldi R, Aloisi AL, Bonasera V, Lattanzio R, Lange R, Weidle UH, Piantelli M, Alberti S: **Upregulation of Trop-2 quantitatively stimulates human cancer growth.** *Oncogene* 2013, 32:222–233.
8. Trerotola M, Jernigan DL, Liu Q, Siddiqui J, Fatatis A, Languino LR: **Trop-2 Promotes Prostate Cancer Metastasis By Modulating  $\beta$ 1 Integrin Functions.** *Cancer Res.* 2013, 73:3155–3167.
9. Kowalsky CA, Faber MS, Nath A, Dann HE, Kelly VW, Liu L, Shanker P, Wagner EK, Maynard JA, Chan C, Whitehead TA: **Rapid fine conformational epitope mapping using comprehensive mutagenesis and deep sequencing.** *J. Biol. Chem.* 2015, 290:26457–26470.
10. Stoyanova T, Goldstein AS, Cai H, Drake JM, Huang J, Witte ON: **Regulated proteolysis of Trop2 drives epithelial hyperplasia and stem cell self-renewal via  $\beta$ -catenin signaling.** *Genes Dev.* 2012, 26:2271–2285.
11. Pavšič M, Ilc G, Vidmar T, Plavec J, Lenarčič B: **The cytosolic tail of the tumor marker protein Trop2 – a structural switch triggered by phosphorylation,** *Sci. Rep.* 2015, 5:10324.
12. Trerotola M, Li J, Alberti S, Languino LR: **Trop-2 inhibits prostate cancer cell adhesion to fibronectin through the  $\beta$ 1 integrin-RACK1 axis.** *J. Cell Physiol.* 2012, 227:3670–3677.

13. Trerotola M, Ganguly KK, Fazli L, Fedele C, Lu H, Dutta A, Liu Q, De Angelis T, Riddell LW, Riobo NA, Gleave ME, Zoubeidi A, Pestell RG, Altieri DC, Languino LR: **Trop-2 is up-regulated in invasive prostate cancer and displaces FAK from focal contacts.** *Oncotarget* 2015, 6:14318-14328.
14. Litvinov SV, Velders MP, Bakker HA, Fleuren GJ, Warnaar SO: **Ep-CAM: a human epithelial antigen is a homophilic cell-cell adhesion molecule.** *J. Cell Biol.* 1994, 125:437-446.
15. Vidmar T, Pavšič M, Lenarčić B: **Biochemical and preliminary X-ray characterization of the tumor-associated calcium signal transducer 2 (Trop2) ectodomain.** *Protein Express Purif.* 2013, 91:69-76.
16. Pavšič M, Gunčar G, Djinović-Carugo K, Lenarčić B: **Crystal structure and its bearing towards an understanding of key biological functions of EpCAM.** *Nat. Commun.* 2014, 5:4764.
17. Liu S, Ho W, Strop P, Dorywalska MG, Rajpal A, Shelton DL, Tran T: **Antibodies specific for TROP-2 and their uses.** *US Patent Application* 2013.
18. Klesmith JR, Bacik J, Wrenbeck EE, Michalczyk R, Whitehead TA: **Trade-offs between enzyme fitness and solubility illuminated by deep mutational scanning.** *Proc. Natl. Acad. Sci. USA* 2017, 114:2265-2270.
19. Iwaki T, Figuera M, Ploplis VA, Castellino FJ: **Rapid selection of *Drosophila* S2 cells with the puromycin resistance gene.** *BioTechniques* 2003, 35:482-486.
20. Medina-Cucurella AV, Whitehead TA: **Characterizing protein-protein interactions using deep sequencing coupled to yeast surface display.** *Methods in Mol. Biol.* 2018, 1764:101-121.