# ADVANCED CANOPY ARCHITECTURE MODELING TO IMPROVE PREDICTION OF MAIZE GROWTH

By

Zhongjie Ji

## A DISSERTATION

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

Plant Breeding, Genetics and Biotechnology - Crop and Soil Sciences - Doctor of Philosophy

#### ABSTRACT

Maize (Zea mays L.) is one of the world's most productive crops, benefiting from advancements in agronomic practices and breeding techniques that harness hybrid vigor. Over the past century, maize yields have experienced an eightfold increase, driven by these innovations. As global challenges related to food security intensify with a rapidly growing human population and increased protein consumption, higher breeding goals for maize have become imperative. The maize canopy profoundly impacts plant growth and yield. By using high throughput phenotyping and manual measurements of various leaf and canopy parameters, we found late-season canopy traits significantly impact yield components. The integration of these leaf canopy traits led to the development of a predictive model for yield, achieving an R2 value of 0.483. Innovative leaf angle measurements and a simulation method for leaf curvatures were also introduced. Integrating yield analysis with canopy traits offers critical insights for maize breeding and cultivating high-yield varieties, advancing productivity by leveraging a deeper understanding of canopy dynamics. Through the analysis of multi-year high-throughput phenotyping data, we established a regression method to align Normalized Difference Vegetation Index (NDVI) data with Growing Degree Days (GDD) across diverse environments. Retaining "stay green" ability emerged as a critical factor impacting yield, quantified through our NDVI-GDD curve, exhibiting strong linear correlations with yield at both plot and environmental levels. Leveraging this, we developed physiological and genomic prediction models for yield, demonstrating promising predictive capabilities across practical breeding scenarios. NDVI, as a high-level trait, displayed correlations with several manually measured traits. Our Genome-Wide Association Study (GWAS) using NDVI as an input trait identified significant signals for these correlated traits, suggesting the potential to replace current manual measurements in breeding

programs. Though genomic selection has become commonplace to expedite breeding cycles, accurately forecasting the performance of new genotypes remains a significant challenge. A study evaluating cross-subpopulation genomic prediction within a small tested cross NAM population revealed that the accuracy of cross-subpopulation prediction fell below that of randomly sampled genetics pools. Additionally, prediction accuracy varied considerably among traits and within prediction subpopulations. We delved into these differences, identifying potential explanations. To enhance cross-subpopulation prediction accuracy, we explored the impact of dominant relationship matrices, Gaussian kernel-based relationships, and LD-adjusted methods, which provided limited improvement. The findings highlight the complexity of genomic prediction in diverse breeding scenarios and the need for further research to enhance accuracy and applicability.

#### ACKNOWLEDGMENTS

As I reflect on my PhD journey, I have had both wonderful and challenging experiences. I am the first graduate student in the Thompson Lab and experienced two years of the COVID pandemic. Despite these challenges, I have always found strength and determination in my commitment to my research and the unwavering support of my family and friends.

Firstly, I would like to express my sincere gratitude to my major advisor Dr. Addie Thompson. Your guidance and expertise were invaluable throughout this research journey. I was slow starting this novel research; you provided the patience needed for me to learn and attempt, even when it led to failure or results that fell short. I am deeply thankful for your unwavering support and insightful feedback.

I would also like to thank my thesis committee members, Dr. Daniel Morris, Dr. Eric Olson and Dr. Shinhan Shiu, for their valuable feedback and suggestions. I am grateful for their time and dedication to my work.

In addition, I would like to thank our lab manager Linsey Newton, for her great help and support in field experiments and data collection which built the solid foundation of my research work. My fellow colleagues in the Thompson Lab provided a supportive academic environment that enriched my research and personal development. I am thankful for the stimulating discussions and collaborative efforts we shared.

I would also like to thank the NRT-IMPACT program for their financial support and for providing me with the opportunity to learn computer science. This has been a valuable experience that has helped me in my research. I am also grateful for the opportunity to meet many new friends through this program. I extend heartfelt thanks to Michigan Corn for their steadfast support throughout the years, catalyzing the progress of our work. Additionally, immense appreciation goes out to the G2F collaborators whose collaborative efforts have been pivotal in shaping our shared successes.

I am profoundly grateful to my parents Lianfu Ji and Xiujuan Zhong for their continuous

iv

encouragement and unwavering belief in me throughout the ups and downs of this academic endeavor. Their boundless love and unwavering support have been a constant source of motivation that carried me through the challenges and triumphs of this journey.

## TABLE OF CONTENTS

CHAPTER 1: GENERAL INTRODUCTION REFERENCES	1 12
CHAPTER 2: FROM HAND MEASUREMENTS TO HIGH THROUGHPUT PHENOTYPING:	
UNDERSTANDING MAIZE CANOPY STRUCTURE AND PREDICTING YIELD	18
REFERENCES	42
APPENDIX	46
CHAPTER 3: A CASE STUDY OF CROSS-SUBPOPULATION GENOMIC PREDICTION	52
REFERENCES	72
APPENDIX	74
CHAPTER 4: TIME-SERIES PHENOTYPING REFLECTS THE MAIZE GROWTH AND	
APPLICATION OF YIELD PREDICTION	76
REFERENCES	98
A DDENIDIY	01
	UI

#### **CHAPTER 1: GENERAL INTRODUCTION**

## Maize

Maize (*Zea mays L.*), also known as corn, is in the Poaceae family. It has a monoecious flowering habit with males (staminate with anthers in tassel) and female (pistillate with silk in ear) flowers separate on different positions of the plant. Maize is predominantly cross-pollinated; pollen randomly lands on silks on the ears of adjacent plants or even its own silks. Maize was first domesticated in south Mexico around 9000 years ago (Matsuoka et al., 2002). Balsas teosinte (*Z. mays* ssp. *Parviglumis*) has been determined to be the progenitor of modern maize (Matsuoka et al., 2002).

Due to advances in agronomic practices and breeding to take advantage of heterosis in hybrids, maize yields have increased eight-fold over the last century. Hence, maize has the greatest global production of any crop species with 800 million tons produced worldwide in 2013, accounting for 32% of the total cereal production (Scott and Emery, 2016).

The maize genome exhibits high levels of genetic complexity and diversity among different inbred lines (Yang et al., 2019). Maize has 10 chromosomes (2n=20), whose genome size has expanded dramatically (to 2.3 gigabases) over the last 3 million years with 85% of its genome encoding transposable elements (Schnable et al., 2009). B73 was the first sequenced maize reference genome published in 2009 (Schnable et al., 2009). Next-Generation sequencing technologies led to rapid changes in plant genomics, and another 7 inbred reference genomes have been published, including high-quality whole-genome assembly of B73 (RefGen\_v4) (Jiao et al., 2017), PHN207 (Hirsch et al., 2016) Mo17 (Sun et al., 2018), W22 (Springer et al., 2018), SK (Yang et al., 2019), and the founder of maize Nested association mapping (NAM) founder population (Hufford et al., 2021).

## Phenotyping

Plant phenotyping is the process of quantifying plant traits such as growth, development, tolerance, resistance, architecture, physiology, ecology, yield, and other quantitative parameters (Li et al., 2014). In the past 20 years, with the development of molecular biology and sequencing technologies, genotyping is no longer a barrier in plant breeding. Efficient, accurate, and meaningfully granular phenotyping becomes a bottleneck in plant genetics and breeding. Even seemingly straightforward traits like yield are really an outcome of the combination of several physiological traits. Though yield is an essential trait in breeding, selection directly on yield causes selection on the underlying factors contributing to the measured yield to be indirect, which may slow progress.

Efficient selection of basic traits with biological support and understanding is needed. But selection on extra traits means extra expense for labor and time. Increasing availability of novel technologies enables dynamic analysis and spatially distinct parameter detection that were previously inaccessible (Fiorani and Schurr, 2013). For example, large scale phenotyping and high-level complex traits can be non-invasively captured by images. In addition, measurement accuracy can be improved by reducing random errors, especially human-induced error.

Modern imaging technologies are key for high-throughput phenotyping. But phenotyping does not end with taking pictures of plants. Images are processed to extract measurements with computing algorithms. Image processing is also applied to photon reflectance, absorbance, and transmission. Image technologies include visible light (400-750nm) or RGB; fluorescence, thermal, and spectral imaging; LiDAR and tomographic imaging and Magnetic Resonance Imaging (MRI), PET, and CT.

Visible imaging is the most widely used and cheapest option in plant phenotyping. With sufficient images and advanced algorithms, shoot biomass (Arvidsson et al., 2011), yield traits (Duan et al., 2011), panicle traits (Ikeda et al., 2010), germination rates (Dias et al., 2011)), leaf morphology (Hoyos-Villegas et al., 2014), seed morphology (Joosen et al., 2012) and root architecture (Clark et al., 2011) can be analyzed. Fluorescence imaging mainly reflects photosynthesis which could be used for monitoring the effects and diagnose early stress responses to abiotic and biotic stress before a deterioration in growth can be measured (Lenk et al., 2007; Konishi et al., 2009; Harbinson et al., 2012). Thermal cameras have the ability to capture leaf surface temperature to study plant water status and stomatal conductance, which are related to abiotic and biotic stresses resulting in impaired function of photosynthesis and transpiration (Munns et al., 2010; Zia et al., 2013). Spectral imaging is used to measure green biomass, canopy chlorophyll content, leaf and canopy senescence and plant water status (Li et al., 2014). The absorption bands at different wavelengths can be analyzed to obtain water content, nitrogen content, pigment composition, biomass, and vegetation indices (Schlemmer et al., 2005; Claudio et al., 2006; Mistele and Schmidhalter, 2008).

Phenotyping platforms vary for different research purposes and design. Controlled environment phenotyping in a growth chamber or greenhouse typically involves robotics and an automated imaging system to measure individual plant traits under certain environments. In a field setting, phenotyping platforms fall into two broad categories: ground-based and aerial. Ground-based platforms consist of vehicles and sensors or cameras for medium-scale phenotyping. Unoccupied aerial vehicles carrying sensors and cameras are capable of rapidly characterizing large-scale field experiments.

#### **Canopy trait phenotyping**

The physical arrangement of above-ground vegetation forms the plant canopy. The canopy structure is important because it is directly related to photosynthesis, or the set of reactions that convert solar energy to chemical energy. To achieve high yield, the balance of plant density and productivity of a single plant requires finely tuned canopy architecture. However, it is difficult to comprehensively measure and trace the growth of such a complex structure.

With the development of high-throughput phenotyping, it is possible to measure multiple canopy traits within a short time and with less manual labor (Li et al., 2014), enabling a more comprehensive understanding of canopy architecture. Unoccupied aerial vehicle (UAV) has become a widely-used platform; many studies have investigated phenotypic information from field experiments (Li et al., 2014; Ju and Son, 2018; Mogili and Deepak, 2018). UAVs can be equipped with RGB, hyperspectral and infrared cameras to screen the fields and capture highresolution images for data analysis. Robots are another platform used in phenotyping and can be equipped with cameras, proximal sensors, and LiDAR (Light Detection And Ranging). Using 3D reconstruction, it is possible to build virtual plants with high fidelity for calculating canopy architecture traits (Lin, 2015). Stereo imaging in which two cameras capture the same objects at the same time from slightly different views can be used to capture the 3D images that also include color information (Biskup et al., 2007). There are mature algorithms to calculate plant height, leaf number and LAI (White et al., 2012). Leaf area distribution, leaf angle distribution and plant radius have been measured as novel traits to describe light interception in the canopy (Perez et al., 2019). Physiological canopy traits like canopy temperature and chlorophyll content can be measured by infrared and hyperspectral sensors (Li et al., 2014). Other novel traits can be

discovered through algorithm improvement and research advancement. RGB camera (Li et al., 2016) or LiDAR (Walter et al., 2019) based technologies offer high-resolution images and point clouds that enable the visualization and extraction of multi-dimensional canopy data, for example latent space phenotypes (LSP). Maize LSP traits are useful to describe plant architecture and biomass distribution, with similar heritabilities as traditionally measured traits (Gage et al., 2019). LSP is also considered as an alternative to traditional methods to learn response-to-treatment traits and recover QTL (Ubbens et al., 2020).

#### Genetics

Understanding the contribution of genetics to phenotypic variation is critical to plant geneticists and breeders. Quantitative genetics introduced statistical methods and molecular markers to detect associations between regions of the genome and traits of interest. Quantitative genetics was initially applied in animal breeding. In 1940, maize breeders introduced basic quantitative genetics theory in breeding programs (Carena et al., 2010). Since then, quantitative genetics has played an important role in maize breeding (Carena et al., 2010). Quantitative trait locus (QTL) mapping has been successfully applied in plant breeding for about 40 years (Tanksley et al., 1982). Enormous numbers of QTL have been identified across the plant kingdom which can be used as QTL-associated markers to select individuals likely to express the desired traits. This maker-assisted selection (MAS) works well for single traits, such as disease resistance and abiotic stress tolerance, which have QTL major effects. But for more complex traits, QTL have limited power for detecting the many small effect loci contributing to important traits which may also be highly affected by environmental factors.

## **Genomic selection & prediction**

Genomic prediction uses genetic marker effects to predict the phenotypic values of traits. The method that uses genomic prediction to make selection for desired traits is called genomic selection. Compared to using several significant trait-associated makers to build up the prediction model, genomic selection considers marker effects across the whole genome simultaneously for prediction of performance and candidate selection (Crossa et al., 2017). Genomic selection aims to enhance genetic gain which excludes the environmental effects. Genomic selected has been shown to be a powerful tool in animal breeding and now has been introduced in plant breeding with the help of high-density genomic markers. In genomic prediction or selection, the population consists of a training set that is both phenotyped and genotyped, and a testing set that is genotyped but not be phenotyped (Crossa et al., 2017). Genome-wide marker effects are computed by fitting phenotypic data and molecular marker from the training set, then applied in the testing set to estimate each line's Genomic Estimated Breeding Value, or GEBV (Meuwissen et al., 2001). GEBV is used to replace the traditionally measured phenotypic data and can be calculated using only genomic markers. Hence, significant time, money and labor on phenotyping can be saved by genomic selection per cycle. Unlike phenotypic selection, genotypic selection can also be conducted in off-season nurseries, shortening the selection cycle and allowing multiple cycles per year.

Based on statistical methods, genomic prediction can be divided into parametric based (penalized approach and Bayesian approach) and non-parametric based approaches. Ridge regression best linear unbiased prediction (RR-BLUP), genomic best linear unbiased prediction (GBLUP), Bayesian A, B, C models and non-linear machine learning models have been applied in genome regression models. There are several factors that affect the accuracy of genomic

prediction, including the population size and genetic diversity and its relationship to the breeding population (Crossa et al., 2014), as well as the heritability of traits under selection. For simple traits with high heritability, genomic prediction accuracy is high, and will easily obtain the ideal genotype and phenotype (Charmet et al., 2014). For complex traits with low heritability, the accuracy of genomic prediction is relatively low, but it still can significantly decrease the cost of the breeding process. The genomic prediction accuracy of complex traits determined by large number of markers are more sensitive to heritability and population size.

Following the first applications in livestock breeding, plant breeders and geneticists have successfully applied genomic selection to breed for a variety of traits. Based on ranking from genomic selection, the top 100 rice hybrids show 16% increased yield compared to the average of all potential hybrids (Xu et al., 2014). The predictive abilities of malting traits in spring barely ranged 0.14 to 0.58 and winter barley ranged between 0.40-0.80 in cross validation (Schmidt et al., 2016). And genomic selection showed higher accuracy on maize under drought stress than traditional selection (Beyene et al., 2015). Prediction based on genomic data also is an efficiency way to enhance the haploid induction rate in maize, compared with phenotypic selection (Almeida et al., 2020).

#### **GxE** interaction

Genotype and the environment interaction ( $G \times E$ ) is the joint effect of genetics and environments. The role  $G \times E$  plays in plant phenotypic variation can even be significantly greater than the sum of the genetic and environmental effects themselves (Heath and Nelson, 2002). G x E, also known as phenotypic plasticity, manifests as a change in phenotype in different environments (El-Soda et al., 2014).

 $G \times E$  variation is thought to arise from a combination of overdominance, pleiotropy, epistasis, linkage, and epigenetic causes. The seemingly stochastic effects of  $G \times E$  may cause divergent responses in a population when introduced to a new environment which could be positive, negative, or null effects.  $G \times E$  may contribute to an advantage in fitness, adaptation and evolution in new environments (Kusmec et al., 2018). But more often, plants exhibit worse performance in novel environments than those to which they are adapted. Understanding of  $G \times E$  is important for domestication of novel germplasm, selection on plasticity, and optimizing agricultural management. Since maize has been adapted to be productive and cultivated worldwide with various habitats and environmental conditions, it provides a perfect model to study  $G \times E$  variation across a wide range of environments.

Multi-environment linear mixed models can handle correlation between environments via GBLUP and have been used to predict phenotypic performance using pedigree and molecular markers. This molecular marker and pedigree GBLUP model was used to quantify  $G \times E$  in wheat genomic prediction. Inclusion of G x E for wheat grain yield led to a 17-34% higher prediction accuracy. Similarity, integration of environmental covariates and crop modeling into genomic prediction of winter wheat led to better performance but less variability than the genomic data-based model (Heslot et al., 2014).

#### **Crop growth modeling**

With advances in plant physiology, soil sciences and micrometeorology, crop growth models (CGM) can also be used to predict plant phenotypic traits (Technow et al., 2015). Some critical yet difficult-to-measure emergent traits from CGMs such as resource capture, utilization of fertilizer, water and solar use efficiency, and allocation among plant organs (also called harvest index) are used as plant breeding traits or targets. CGMs use mathematical equations to

quantify the interaction of biological and environmental factors in a real dynamic soil-plantenvironment system which can simulate trajectories of the whole crop life cycle. However, as CGMs are process-based models often tuned to a discrete number of varieties, plant breeders face barriers to genotype-based prediction of plant phenotypes. If this limitation is overcome, CGMs have the potential to tackle the genotype to phenotype prediction problems in plant breeding. CGMs have a unique advantage in their characterization of environmental factors which could play an essential role in predicting complex traits heavily affected by genotype × environment × management interactions.

In early stages, plant scientists and breeders used CGMs to understand the physiology of adaptative traits in QTL studies. Later, epistasis became an important consideration in trait formation. Because CGMs include nonlinear relationships among phenotypic traits, they have potential to quantify the epistatic effects in genomic prediction models by integrating with biological knowledge (Technow et al., 2015).

To define an individual variety within a CGM, it is necessary to input a set of parameters that define how that variety grows and develops throughout its life cycle. However, the scale of phenotyping for vital physiological traits is insufficient to provide sound and adequate parameters to fully define a variety within a CGM. With advances in imaging technologies and matched algorithms, phenotyping may no longer be a bottleneck in applicability of CGMs, as plenty of data can be collected and generated in a short time. Conversely, a good and efficient CGM can reduce the need for phenotyping once complex traits can be accurately estimated via integrated biological models informed by low-level traits and environmental data. Progress in both CGMs and phenotyping will lead to a better understanding of plant physiology, as well as easier and more efficient selection in plant breeding.

The prediction yield in maize is determined by kernel number and size, with crop, plant or ear growth rate indicating the health status of the plant (Vega et al., 2001; Echarte et al., 2004). The impact of water, nitrogen and density stress on kernel traits is most significant around anthesis, which is incorporated in the simulation of maize yield in CGM (Soufizadeh et al., 2018). CGMs are a powerful tool to explore yield potential and conduct gap analysis to optimize management of sorghum varieties in a dryland production system under current climate and future climate conditions of Australia (Hammer et al., 2014, 2020). Soufizadeh (2018) used the comprehensive APSIM maize model to accurately predict the N dynamics on an organ and crop level scale across a range of genotypes, environments, management strategies, and their interaction.

#### **Genomes to Fields Project**

The maize Genomes to Field (G2F) project, launched in 2014 and continuing each year, attempts to understand G × E variation. This project has measured phenotypic traits in hybrid maize across the North American corn belt with wide geographical and climatic diversity. The G2F project is an umbrella interdisciplinary project which connects genetics, genomics, plant physiology, agronomy, climatology, geography, and crop modeling with the analytical tools derived from computational sciences, statistics, and engineering (AlKhalifah et al., 2018). During the growth season, G2F collaborators collected large-scale genotypic, phenotypic, environmental and metadata datasets (McFarland et al., 2020).

#### Importance

With the explosion of human population, deterioration of the environment, the conflict of rapidly decreasing arable farmland and dramatically increased demand of food, agriculture production has been challenged. In addition, increasing meat consumption in developing

countries and biofuel usage in developed countries make this more compelling. Hence, high and stable yield are two fundamental but urgent targets for plant breeding and agriculture. Maize, as one of most productive crops, with the largest cultivated area will play an important role in the crisis of food and agriculture. Leaf and canopy traits are vital to capture the solar energy to conduct photosynthesis. Due to limitation of labors and costs, leaf and canopy traits are not easy to comprehensively measure. The relationship of non-leaf phenotypic traits and leaf photosynthetic organs is worth exploration. High through-put phenotyping tools can be applied to record and quantify the leaf and canopy. Such data can be used to develop new types of crop growth models or as a supplement to modify existing crop growth models. Because of the complexity of leaf and canopy characteristics, its overall architecture cannot be well explained by several major QTL. Genomic prediction enables estimation and prediction of these traits with small effect alleles. Leaf and canopy traits are sensitive to environmental factors, so this project also provides valuable data for GxE research. Understanding leaf and canopy traits would provide a direction to find an ideotype of maize plants, and to optimize plant density in field production.

Global climate change is a fact that will challenge agricultural production and food security. This multi-location and multi-environmental GxE study provide extensive data to interpret the impacts of environmental and GxE effects. CGM with environmental data and associated genetic information could lead to more accurate prediction into new climate scenarios, and what genotypes will thrive. Plant breeders will have the ability to mitigate the detrimental effects of climate before they occur. Environmentally customized breeding and management could be a promising way to improve yield potential.

#### REFERENCES

- AlKhalifah, N., D.A. Campbell, C.M. Falcon, J.M. Gardiner, N.D. Miller, et al. 2018. Maize Genomes to Fields: 2014 and 2015 Field Season Genotype, Phenotype, Environment, and Inbred Ear Image Datasets. *BMC Research Notes 11*(1): 452. https://doi.org/10.1186/s13104-018-3508-1
- Almeida, V.C., H.U. Trentin, U.K. Frei, and T. Lübberstedt. 2020. Genomic prediction of maternal haploid induction rate in maize. *The Plant Genome* 13(1): e20014. https://doi.org/10.1002/tpg2.20014.
- Arvidsson, S., P. Pérez-Rodríguez, and B. Mueller-Roeber. 2011. A growth phenotyping pipeline for Arabidopsis thaliana integrating image analysis and rosette area modeling for robust quantification of genotype effects. *New Phytologist 191*(3): 895–907. https://doi.org/10.1111/j.1469-8137.2011.03756.x.
- BEADLE, G.W. 1939. Teosinte and the origin of Maize. Journal of Heredity 30(6): 245–247.
- Beyene, Y., K. Semagn, S. Mugo, A. Tarekegne, R. Babu, et al. 2015. Genetic Gains in Grain Yield Through Genomic Selection in Eight Bi-parental Maize Populations under Drought Stress. *Crop Science* 55(1): 154–163. https://doi.org/10.2135/cropsci2014.07.0460.
- Biskup, B., H. Scharr, U. Schurr, and U. rascher. 2007. A stereo imaging system for measuring structural parameters of plant canopies. *Plant, Cell & Environment 30*(10): 1299–1308. https://doi.org/10.1111/j.1365-3040.2007.01702.x.
- Burgueño, J., G. de los Campos, K. Weigel, and J. Crossa. 2012. Genomic prediction of breeding values when modeling genotype × environment interaction using pedigree and dense molecular markers. *Crop Science* 52(2): 707–719. https://doi.org/10.2135/cropsci2011.06.0299.
- Carena, M.J., A.R. Hallauer, and J.B. Miranda Filho. 2010. Quantitative genetics in maize breeding. *Springer New York, New York*, NY.
- Charmet, G., E. Storlie, F.X. Oury, V. Laurent, D. Beghin, et al. 2014. Genome-wide prediction of three important traits in bread wheat. *Molecular Breeding* 34(4): 1843–1852. https://doi.org/10.1007/s11032-014-0143-y.
- Clark, R.T., R.B. MacCurdy, J.K. Jung, J.E. Shaff, S.R. McCouch, et al. 2011. Threedimensional root phenotyping with a novel imaging and software platform. *Plant Physiol*. 156(2): 455. https://doi.org/10.1104/pp.110.169102.
- Claudio, H.C., Y. Cheng, D.A. Fuentes, J.A. Gamon, H. Luo, et al. 2006. Monitoring drought effects on vegetation water content and fluxes in chaparral with the 970 nm water band index. *Remote Sensing of Environment 103*(3): 304–311. https://doi.org/10.1016/j.rse.2005.07.015.

- Crossa J, Campos G de los, Pérez P, et al. Prediction of genetic values of quantitative traits in plant breeding using pedigree and molecular markers. Genetics. 2010;186:713–724.
- Crossa, J., P. Pérez, J. Hickey, J. Burgueño, L. Ornella, et al. 2014. Genomic prediction in CIMMYT maize and wheat breeding programs. *Heredity* 112(1): 48–60. https://doi.org/10.1038/hdy.2013.16.
- Crossa, J., P. Pérez-Rodríguez, J. Cuevas, O. Montesinos-López, D. Jarquín, et al. 2017. Genomic selection in plant breeding: methods, models, and perspectives. *Trends in Plant Science* 22(11): 961–975. https://doi.org/10.1016/j.tplants.2017.08.011.
- Dias, P.M.B., S. Brunel-Muguet, C. Dürr, T. Huguet, D. Demilly, et al. 2011. QTL analysis of seed germination and pre-emergence growth at extreme temperatures in Medicago truncatula. *Theor Appl Genet 122*(2): 429–444. https://doi.org/10.1007/s00122-010-1458-7.
- Doebley, J., and L. Lukens. 1998. Transcriptional regulators and the evolution of plant form. *Plant Cell 10*(7): 1075. https://doi.org/10.1105/tpc.10.7.1075.
- Duan, L., W. Yang, C. Huang, and Q. Liu. 2011. A novel machine-vision-based facility for the automatic evaluation of yield-related traits in rice. *Plant Methods* 7(1): 44. https://doi.org/10.1186/1746-4811-7-44.
- Dwyer, L.M., and D.W. Stewart. 1986. Leaf area development in field-grown maize1. *Agronomy Journal* 78(2): 334–343. https://doi.org/10.2134/agronj1986.00021962007800020024x.
- Echarte, M., A. Conchillo, D. Ansorena, and I. Astiasarán. 2004. Evaluation of the nutritional aspects and cholesterol oxidation products ff pork liver and fish patés. *Food Chemistry* 86(1): 47–53. https://doi.org/10.1016/j.foodchem.2003.08.027.
- Eubanks, M. 1995. A cross between two maize relatives: Tripsacum Dactyloides Andzea Diploperennis (Poaceae). *Economic Botany* 49(2): 172–182. https://doi.org/10.1007/BF02862921.
- Fiorani, F., and U. Schurr. 2013. Future scenarios for plant phenotyping. *Annu. Rev. Plant Biol.* 64(1): 267–291. https://doi.org/10.1146/annurev-arplant-050312-120137.
- Francis, C.A., J.N. Rutger, and A.F.E. Palmer. 1969. A rapid method for plant leaf area estimation in maize (Zea mays L.)1. *Crop Science* 9(5): cropsci1969.0011183X000900050005x.
- Gage, J.L., E. Richards, N. Lepak, N. Kaczmar, C. Soman, et al. 2019. In-field whole-plant maize architecture characterized by subcanopy rovers and latent space phenotyping. *The Plant Phenome Journal* 2(1): 190011. https://doi.org/10.2135/tppj2019.07.0011.
- Hammer, G.L., G. McLean, S. Chapman, B. Zheng, A. Doherty, et al. 2014. crop design for specific adaptation in variable dryland production environments. *Crop Pasture Sci.* 65(7): 614–626.

- Hammer, Graeme.L., G. McLean, E. van Oosterom, S. Chapman, B. Zheng, et al. 2020. Designing crops for adaptation to the drought and high-temperature risks anticipated in future climates. *Crop Science* 60(2): 605–621. https://doi.org/10.1002/csc2.20110.
- Harbinson, J., A.E. Prinzenberg, W. Kruijer, and M.G. Aarts. 2012. High throughput screening with chlorophyll fluorescence imaging and its use in crop improvement. *Current Opinion in Biotechnology 23*(2): 221–226. https://doi.org/10.1016/j.copbio.2011.10.006.
- Heath, A.C., and E.C. Nelson. 2002. Effects of the interaction between genotype and environment. Research into the genetic epidemiology of alcohol dependence. *Alcohol Res Health* 26(3): 193–201.
- Heslot, N., D. Akdemir, M.E. Sorrells, and J.-L. Jannink. 2014. Integrating environmental covariates and crop modeling into the genomic selection framework to predict genotype by environment interactions. *Theor Appl Genet 127*(2): 463–480. https://doi.org/10.1007/s00122-013-2231-5.
- Hirsch, C.N., C.D. Hirsch, A.B. Brohammer, M.J. Bowman, I. Soifer, et al. 2016. Draft Assembly of Elite Inbred Line PH207 Provides insights into genomic and transcriptome diversity in maize. *Plant Cell* 28(11): 2700. https://doi.org/10.1105/tpc.16.00353.
- Hoyos-Villegas, V., J.H. Houx, S.K. Singh, and F.B. Fritschi. 2014. Ground-based digital imaging as a tool to assess soybean growth and yield. *Crop Science* 54(4): 1756–1768. https://doi.org/10.2135/cropsci2013.08.0540.
- Hufford, M.B., A.S. Seetharam, M.R. Woodhouse, K.M. Chougule, S. Ou, et al. 2021. *De novo* assembly, annotation, and comparative analysis of 26 diverse maize genomes. *bioRxiv*: 2021.01.14.426684. https://doi.org/10.1101/2021.01.14.426684.
- Ikeda, M., Y. Hirose, T. Takashi, Y. Shibata, T. Yamamura, et al. 2010. Analysis of rice panicle traits and detection of qtls using an image analyzing method. *Breeding Science* 60(1): 55–64. https://doi.org/10.1270/jsbbs.60.55.
- Iltis, H.H. 1983. From Teosinte to Maize: The catastrophic sexual transmutation. science 222(4626): 886. https://doi/org/10.1126/Science.222.4626.886.
- Jarquín, D., J. Crossa, X. Lacaze, P. Du Cheyron, J. Daucourt, et al. 2014. A reaction norm model for genomic selection using high-dimensional genomic and environmental data. *Theor Appl Genet* 127(3): 595–607. https://doi.org/10.1007/s00122-013-2243-1.
- Jiao, Y., P. Peluso, J. Shi, T. Liang, M.C. Stitzer, et al. 2017. Improved maize reference genome with single-molecule technologies. *Nature* 546(7659): 524–527. https://doi.org/10.1038/nature22971.
- Joosen, R.V.L., D. Arends, L.A.J. Willems, W. Ligterink, R.C. Jansen, et al. 2012. Visualizing the genetic landscape of Arabidopsis seed performance. *Plant Physiol.* 158(2): 570. https://doi.org/10.1104/pp.111.186676.

- Ju, C., and H.I. Son. 2018. Multiple UAV systems for agricultural applications: Control, Implementation, and Evaluation. Electronics 7(9). https://doi.org/10.3390/electronics7090162.
- Konishi, A., A. Eguchi, F. Hosoi, and K. Omasa. 2009. 3d monitoring spatiotemporal effects ff herbicide on a whole plant using combined range and chlorophyll a fluorescence imaging. *Funct. Plant Biol.* 36(11): 874–879.
- Kusmec, A., N. de Leon, and P.S. Schnable. 2018. Harnessing phenotypic plasticity to improve maize yields. *Frontiers in Plant Science* 9: 1377. https://doi.org.10.3389/fpls.2018.01377.
- Lenk, S., L. Chaerle, E.E. Pfündel, G. Langsdorf, D. Hagenbeek, et al. 2007. Multispectral fluorescence and reflectance imaging at the leaf level and its possible applications. *Journal of Experimental Botany* 58(4): 807–814. https://doi.org/10.1093/jxb/erl207.
- Li, W., Z. Niu, H. Chen, D. Li, M. Wu, et al. 2016. Remote estimation ff canopy height and aboveground biomass ff maize using high-resolution stereo images from a low-cost unmanned aerial vehicle system. *Ecological Indicators* 67: 637–648. https://doi.org/10.1016/j.ecolind.2016.03.036.
- Li, L., Q. Zhang, and D. Huang. 2014. A review of imaging techniques for plant phenotyping. *Sensors 14*(11). https://doi.org/10.3390/s141120078.
- Lin, Y. 2015. LiDAR: An important tool for next-generation phenotyping technology of high potential for plant phenomics? *Computers and Electronics in Agriculture 119*: 61–73. https://doi.org/10.1016/j.compag.2015.10.011.
- Mangelsdorf PC, Reeves RG 1939. The origin of Indian corn and its relatives. Texas Agricultural Experiment Station Bulletin No. 574. College Station, TX, USA: Texas Agricultural Experimental Station.
- Matsuoka, Y., Y. Vigouroux, M.M. Goodman, J. Sanchez G., E. Buckler, et al. 2002. A single domestication for maize shown by multilocus microsatellite genotyping. *Proc Natl Acad Sci USA 99*(9): 6080. https://doi.org/10.1073/pnas.052125199.
- McFarland, B.A., N. AlKhalifah, M. Bohn, J. Bubert, E.S. Buckler, et al. 2020. Maize genomes to fields (g2f): 2014–2017 field seasons: genotype, phenotype, climatic, soil, and inbred ear image datasets. bmc research notes 13(1): 71. https://doi.org10.1186/s13104-020-4922-8.
- Meuwissen, T.H.E., B.J. Hayes, and M.E. Goddard. 2001. Prediction of total genetic value using genome-wide dense marker maps. *Genetics* 157(4): 1819–1829. https://doi.org/10.1093/genetics/157.4.1819.
- Mistele, B., and U. Schmidhalter. 2008. Spectral measurements of the total aerial n and biomass dry weight in maize using a quadrilateral-view optic. *Field Crops Research 106*(1): 94–103. https://doi.org/: 10.1016/j.fcr.2007.11.002.

- Mogili, U.R., and B.B.V.L. Deepak. 2018. Review on application of drone systems in precision agriculture. *Procedia Computer Science* 133: 502–509. https://doi.org/10.1016/j.procs.2018.07.063.
- Munns, R., R.A. James, X.R.R. Sirault, R.T. Furbank, and H.G. Jones. 2010. New phenotyping methods for screening wheat and barley for beneficial responses to water deficit. *Journal of Experimental Botany* 61(13): 3499–3507. https://doi.org/10.1093/jxb/erq199.
- Perez, R.P.A., C. Fournier, L. Cabrera-Bosquet, S. Artzet, C. Pradal, et al. 2019. Changes in the vertical distribution of leaf area enhanced light interception efficiency in maize over generations of selection. *Plant, Cell & Environment* 42(7): 2105–2119. https://doi.org/10.1111/pce.13539.
- Schlemmer, M.R., D.D. Francis, J.F. Shanahan, and J.S. Schepers. 2005. Remotely measuring chlorophyll content in corn leaves with differing nitrogen levels and relative water content. *Agronomy Journal* 97(1): 106–112. https://doi.org/0.2134/agronj2005.0106.
- Schmidt, M., S. Kollers, A. Maasberg-Prelle, J. Großer, B. Schinkel, et al. 2016. Prediction of malting quality traits in barley based on genome-wide marker data to assess the potential of genomic selection. *Theor Appl Genet 129*(2): 203–213. https://doi.org/10.1007/s00122-015-2639-1.
- Schnable, P.S., D. Ware, R.S. Fulton, J.C. Stein, F. Wei, et al. 2009. The B73 maize genome: complexity, diversity, and dynamics. *Science* 326(5956): 1112. https://doi.org/10.1126/science.1178534.
- Soufizadeh, S., E. Munaro, G. McLean, A. Massignam, E.J. van Oosterom, et al. 2018. Modelling the nitrogen dynamics of maize crops – Enhancing the APSIM maize model. *European Journal of Agronomy 100*: 118–131. https://doi.org/10.1016/j.eja.2017.12.007.
- Springer, N.M., S.N. Anderson, C.M. Andorf, K.R. Ahern, F. Bai, et al. 2018. The maize W22 genome provides a foundation for functional genomics and transposon biology. *Nature Genetics* 50(9): 1282–1288. https://doi.org/10.1038/s41588-018-0158-0.
- Sun, S., Y. Zhou, J. Chen, J. Shi, H. Zhao, et al. 2018. Extensive intraspecific gene order and gene structural variations between Mo17 and other maize genomes. *Nature Genetics* 50(9): 1289–1295. https://doi.org/10.1038/s41588-018-0182-0.
- Tanksley, S.D., H. Medina-Filho, and C.M. Rick. 1982. Use of naturally-occurring enzyme variation to detect and map genes controlling quantitative traits in an interspecific backcross of tomato. *Heredity* 49(1): 11–25. https://doi.org/10.1038/hdy.1982.61.
- Technow, F., C.D. Messina, L.R. Totir, and M. Cooper. 2015. Integrating crop growth models with whole genome prediction through approximate Bayesian computation. *PLOS ONE 10*(6): e0130855. https://doi.org/10.1371/journal.pone.0130855.

- Ubbens, J., M. Cieslak, P. Prusinkiewicz, I. Parkin, J. Ebersbach, et al. 2020. Latent space phenotyping: automatic image-based phenotyping for treatment studies. *Plant Phenomics* 2020: 5801869. https://doi.org/10.34133/2020/5801869.
- Vega, C.R.C., F.H. Andrade, and V.O. Sadras. 2001. Reproductive partitioning and seed set efficiency in soybean, sunflower and maize. *Field Crops Research* 72(3): 163–175. https://doi.org/10.1016/S0378-4290(01)00172-1.
- Walter, J.D.C., J. Edwards, G. McDonald, and H. Kuchel. 2019. Estimating biomass and canopy height with lidar for field crop breeding. *Frontiers in Plant Science 10*: 1145. https://doi.org/10.3389/fpls.2019.01145.
- White, J.W., P. Andrade-Sanchez, M.A. Gore, K.F. Bronson, T.A. Coffelt, et al. 2012. Fieldbased phenomics for plant genetics research. *Field Crops Research* 133: 101–112. https://doi.org/10.1016/j.fcr.2012.04.003.
- Xu, S., D. Zhu, and Q. Zhang. 2014. Predicting hybrid performance in rice using genomic best linear unbiased prediction. *Proc Natl Acad Sci* USA 111(34): 12456. https://doi.org/10.1073/pnas.1413750111.
- Yang, N., J. Liu, Q. Gao, S. Gui, L. Chen, et al. 2019. Genome assembly of a tropical maize inbred line provides insights into structural variation and crop improvement. *Nature Genetics* 51(6): 1052–1059. https://doi.org/10.1038/s41588-019-0427-6.
- Zia, S., G. Romano, W. Spreer, C. Sanchez, J. Cairns, et al. 2013. Infrared thermal imaging as a rapid tool for identifying water-stress tolerant maize genotypes of different phenology. *Journal of Agronomy and Crop Science* 199(2): 75–84. https://doi.org/10.1111/j.1439-037X.2012.00537.x.

## CHAPTER 2: FROM HAND MEASUREMENTS TO HIGH THROUGHPUT PHENOTYPING: UNDERSTANDING MAIZE CANOPY STRUCTURE AND PREDICTING YIELD

#### Abstract

The maize canopy plays a crucial role in determining plant growth, development, and ultimately, crop yield. Understanding the complex interplay between canopy traits and yield is vital for effective maize breeding and crop management strategies. This study comprehensively examined various canopy traits, including hand-measured height, flowering characteristics, leaf growth parameters, leaf angles, leaf size, stand count, lodging, and time series vegetation indices (VIs), as well as canopy cover rate obtained from Unoccupied Aerial Vehicles (UAVs). The analysis focused on investigating the relationship between canopy traits and yield. The results revealed that late-season canopy traits played a more significant role in yield components. A promising yield prediction model, utilizing the analyzed canopy traits, achieved an R<sup>2</sup> value of 0.483. Additionally, the study introduced new leaf angle measurements and proposed a novel simulation method for leaf curvatures based on these measurements. The integration of yield and canopy analysis provides valuable insights for maize breeding and guidance for achieving ideal maize types and high-yield cultivars. These findings contribute to advancing maize breeding efforts and enhancing crop productivity through improved understanding and utilization of canopy traits.

## Introduction

Due to the rapidly expanding population, there is a significant gap between the increasing demand for food and the current supply. Maize, being one of most widely planted and productive grain crops, has shown great potential in the past five decades, resulting in improved yields. Now, maize is expected to play an even more crucial role in ensuring future food security.

The maize plant exhibits distinct morphological characteristics, including its large, elongated leaves that grow in an alternate pattern along the stem. The canopy plays a crucial role in maize productivity by capturing and utilizing solar radiation (T. Liu et al., 2018) and directly impacting plant density, which is a key determinant of yield (Duncan, 1971; Tokatlidis & Koutroubas, 2004). Canopy structure is highly complex, as it is influenced by multiple leaf traits, such as leaf number, size, ear leaf position, curvature, and inclination. Moreover, intricate physiological interactions exist among these leaf traits. For example, there tends to be a negative relationship between individual leaf size and leaf numbers, and the ear leaf position is typically at around 2/3 of the total leaves (Birch et al., 1997).

Canopy structure is influenced not only by genetic factors but also by environmental conditions (Chen et al., 2019). Research has shown that higher plant densities can induce changes in leaf and canopy morphology (Sher et al., 2017). Dense plant populations often exhibit smaller and thinner leaves, more upright leaf angles at the individual plant level, and a greater overall leaf area index (LAI) (J. Li et al., 2018; Pendleton et al., 1968; WILSON et al., 1998). This compact canopy configuration enhances light interception and utilization within the canopy, leading to increased photosynthetic rates and higher yields (Du et al., 2021). However, excessive plant density can result in competition for light, lodging, negatively impacting crop yield (Sher et al., 2018; H. Zhang et al., 2022). Achieving high crop yields to meet food demands requires implementing a high plant density with appropriate canopy structure (Sher et al., 2016; Tokatlidis et al., 2011).

Meanwhile, throughout the growing season, the maize canopy is dynamically changing as it progresses through various growth stages and eventually withers during maturity. With the advancement of high-throughput phenotyping, it has become possible to efficiently measure

multiple canopy traits within a short time period with reduced manual labor (L. Li et al., 2014), leading to a more comprehensive understanding of canopy architecture. Unoccupied aerial vehicles (UAVs) have become widely used platforms, and numerous studies have investigated the phenotypic information obtained from field experiments using UAVs. Some UAV-based spectral vegetation indices (VIs) can reflect the status of the canopy and are utilized for yield prediction. The Normalized Difference Vegetation Index (NDVI) is among the most commonly used indices to quantify vegetation growth and health (Ozyavuz et al., 2015). The NDVI value indicates higher vegetation density and photosynthetic activity (Gamon et al., 1995). While NDVI has been used to estimate LAI in crops under low LAI conditions, it exhibits limitations when estimating high LAI due to saturation issues (Hashimoto et al., 2019; Wang et al., 2005) (paper list https://www.mdpi.com/2072-4292/12/12/1979). A robust relationship has been observed between NDVI values and biomass accumulation in maize [11]. UAV-based NDVI has also demonstrated promising results in predicting yield (Tamás et al., 2023) and above ground biomass (Meiyan et al., 2022).

Various other multispectral vegetation indices have been developed and applied for canopy research and yield prediction in different crop species. The Green Chlorophyll Vegetation Index (GCVI) and RedEdge Chlorophyll Vegetation Index (RECVI) Index have been reported to remain sensitive in high LAI estimation, while other normalized difference VIs tend to saturate more quickly (Nguy-Robertson et al., 2012; Viña et al., 2011). Wu et al. (Wu et al., 2019) used time-series Blue Normalized Difference Vegetation Index (BNDVI) to estimate the grain yield of 25 maize hybrids and achieved an average correlation coefficient (r) = 0.51 using a partial least squares regression (PLSR) model. Danilevicz et al. (Danilevicz et al., 2021) successfully predicted the yield of hybrid maize by combining multispectral vegetation indices

and genotype information under different field management, with an accuracy of R2 = 0.73. Fan et al assessed maize grain yield using a survey of 200 continuous bands from 400–1000 nm and ridge regression, achieving the highest r = 0.54.

To improve our understanding of maize leaf and canopy traits and facilitate leaf canopy optimization in maize breeding, this study aims to (1) investigate the interactions among a set of leaf canopy traits, (2) utilize high-throughput phenotyping to evaluate canopy changes throughout the growing season, and 3) understand the impact of the canopy on yield and develop predictive models for yield using leaf canopy traits.

#### **Materials & Methods**

The field experiments were conducted in East Lansing, Michigan, as part of the Genomes to Fields initiative (www.genomes2fields.org). The Genomes to Fields initiative is a collaborative project conducted across multiple locations, aiming to understand the interactions between maize genotypes and environments to accurately predict phenotypic traits in different conditions. The hybrid materials used in this study were developed from a PHW65 reference population, including doubled haploids derived from PHW65/PHN11, PHW65/Mo44, and PHW65/MoG crosses, which were then test-crossed to PHT69. The Michigan location consisted of 500 two-row plots. Because of germination, traffic and machine harvest problem, 468 plots were used in this study. Since this study was focused on examining the relationship between leaf and canopy phenotypes at the plot level, the analysis and discussion primarily revolved around phenotypic rather than genotypic aspects.

To collect traits in the field, we utilized the Field Book app (phenoapps.org) on cellphones following the Genomes to Fields standard operating procedure (https://www.genomes2fields.org/resources/). For each plot, we recorded the anthesis date,

which corresponds to the date when 50% of the plants' anthers were visible on over half of the main tassel spike. Similarly, the silking date was determined as the date when 50% of the plants' silks emerged in each individual plot. Before flowering, we counted green snap as the number of plants broken between the ground level and the top ear node. At harvest time, we assessed stand count (number of plants), stalk lodging (number of plants broken at ground level and top ear node, and root lodging (number of plants with substantially leaning stems). We selected two representative plants per plot – one each from the left and right rows – to measure plant-scale traits. For these two selected plants, we measured the distance from the ground soil line to the ligule of the flag leaf as the plant height, and from the ground soil line to the top ear-bearing node as the ear height. To accurately determine the ear leaf number and total leaf number, we marked the specific leaf number weekly during the plant vegetative stage. Based on these marked leaves, we counted the leaf number located at the top ear node as the ear leaf number, and the flag leaf below the tassel as total leaf number. Additionally, we recorded the leaf number of the uppermost ligule (top collared leaf) at three timepoints during V4-V12 growth stages, shown in Table 2.1.

A WatchDog weather station (Spectrum Technologies, Aurora, IL, USA) was installed in the field to monitor the weather conditions. The dates recorded during the study were converted to growing degree days (GDD), which provide a meaningful measure for explaining physiological changes compared to simply counting days after planting. GDD was calculated on a daily basis using the following formula: GDD = (Tmax + Tmin) / 2 - Tbase, where Tmax is the maximum temperature, Tmin is the minimum temperature, and Tbase is set at 50 °F as the base temperature. Any temperature exceeding 86 °F was capped at Tmax = 50 °F before calculating the GDD. Likewise, the lower limit temperature was set at 10°C. Leaf initiation rate was

calculated as the number of GDDs needed to generate each collared leaf during vegetative growth by using the leaf number counts along with the GDD values for the dates of their collection.

Ear leaf area was measured in the field on two representative plants either by scanning or by measuring their length and width. Leaf area can be estimated using the equation 0.75 x leaf length x width (Francis et al., 1969). By knowing the ear leaf (largest leaf) area, the largest leaf number, and the total leaf number, the total leaf area can be calculated using the equation (Dwyer & Stewart, 1986):  $Y = Y_0 \exp[-b(X - X_0)^2 + c(X - X_0)^3]$ , where X is the number of an individual leaf,  $X_0$  is the number of the largest leaf, Y is the leaf area of an individual leaf,  $Y_0$  is the largest leaf area,  $b = -0.009 - (\exp -0.2 * n)$ , and  $c = 0.0006 - (\exp -0.43 * n)$ , where n is the total leaf number. The estimated total leaf area index (LAI) per plot was calculated using the equation: TLA x (stand count – green snap) / plot area.

Ear Leaf Angle and Simulation:

Traditionally, leaf angle was defined as the stem angle (Figure 2.1B), but this approach overlooks changes in angle and curvature, which have physiological significance. To address this, we developed a tool capable of simultaneously measuring angle and distance (Figure 2.1 A). In addition to the stem angle, we measured two additional points: the highest point on the leaf and the tip of the leaf (Figure 2.1 C). In the field, we measured two leaves on the two marked plants per plot: one below the ear and one above the ear (L1 and L2 respectively). For some plants, the end angle exceeded the measurement limit (>90 degrees), so it was recorded as 90 degrees and classified into a new category (endangle1 and endangle2 for L1 and L2).



Figure 2.1: Tools, methods, and simulation for leaf angle. A: Customized tool for simultaneously measuring leaf angle and length. B: Examples for traditional leaf angle measurements C: Examples for our novel leaf angle and length measurements. D&E: Two examples of leaf curvature simulations.

Based on the measured leaf traits, polynomial functions were used to simulate the leaf curvatures. We observed that no single polynomial curve would reflect all the measured traits from stem to end. We thus modeled each blade as the continuous union of two separate polynomial curves. The first one modeled the blade from stem to its apex  $f(x) = A_1(x - B)^{N1} + C$ ,

while the second one modeled the blade from its apex to its end  $g(x) = A_2(x - B)^{N^2} + C$  (Figure A2.1). This second curve was ignored for leaves whose end point was behind their apex, meaning that the end of the leaf was curled inwards. The parameters B and C correspond to the location of the apex, which is the same for both curves f and g. All the parameters were computed following basic trigonometrical manipulations based on the measured lengths and angles. To avoid curves that looked too distorted, we limited the possible values of N1 and N2 to be between 1 and 2.25. Once the polynomial curve parameters were determined, the curve length of f(x), g(x), and their union was computed by approximating them as piecewise continuous linear functions (Figure 2.2).



Figure 2.2: Simulated leaf curvature for leaves. A: Leaf curvature from base to apex point simulated by a polynomial function. B: Leaf curvature from apex to tip point simulated by a second polynomial function. C: A combined leaf curvature simulation from the two parts. The color in the figures represents the curvature length, with lighter shades indicating longer leaves.

## **Drone Data Acquisition and Processing**

Aerial survey flights equipped with RGB and multispectral cameras were conducted at eight timepoints throughout the entire growing season, ensuring clear and windless conditions. Table 2.1 presents the corresponding flight time and GDD information. To ensure high-quality aerial image data, manual weeding was performed to maintain a clean field. The multispectral sensor used in this study consisted of five bands: blue (475 nm with 32 nm FWHM), green (560 nm with 27 nm FWHM), red (668 nm with 14 nm FWHM), rededge (717 nm with 12 nm

FWHM), and near infrared (842 nm with 57 nm FWHM) (MicaSense Inc., Seattle, WA, USA;

http://www.micasense.com/)

Date	GDD	Measurements
5/27/19	-	Machine planting
6/24/19	305.95	Stand count
7/1/19	480.8	Full expanded leaf count
7/2/19	508.35	Drone: RGB+ Multispec
7/10/19	698.1	Full expanded leaf count
7/15/19	810.2	Drone: RGB+ Multispec
7/17/19	865.7	Full expanded leaf count
7/24/19	1028.05	Green Snap
7/28/19	1124.05	Drone: RGB+ Multispec
8/12/19	1438.9	Drone: RGB+ Multispec
9/2/19	1833.95	Drone: RGB+ Multispec
9/11/19	1977.65	Drone: RGB+ Multispec
10/7/19	2329.25	Drone: RGB+ Multispec
11/4/19	2437.15	Stalk and Root Lodging

Table 2.1: Date and GDD for each measurement.

UAV flight paths were guided by GPS-autopilot using pre-programmed grid survey patterns across the field plots. The collected UAV images were processed using Pix4Dmapper to generate orthophotos for both RGB and individual reflectance bands. Initially, a rough polygon shape file with labeled plot information was created using UAStools (Anderson & Murray, 2020) in R. Subsequently, in ArcMAP, the polygon boundaries were refined based on high-resolution RGB orthophotos.

For each flight survey, plant pixels were extracted using a customized threshold applied to the RGB orthophotos. Based on the single band reflectance value observed from each flight, NIR was relatively independent, the rest of four bands are highly correlated and correlated (Figure A2.2). The five multispectral band orthophotos from each flight were used to calculate the Normalized Difference Vegetation Index (NDVI) and Green Chlorophyll Vegetation Index (GCVI) maps. Reflectance values for different bands and vegetation indices were extracted by averaging all valid pixels within each bounded plot using the defined polygons. The ratio of pure plant pixels to total pixels within a plot was considered the canopy cover rate.

## **Models for Prediction**

In this study, we trained two commonly used and interpretable machine learning regression models for yield prediction. The first model is Ridge regression, a linear regression technique that addresses the issue of multicollinearity in the data. It incorporates a penalty parameter to control the amount of shrinkage applied to the coefficients. The second model is Partial Least Squares Regression (PLSR), a statistical method used for predictive modeling in multivariate data analysis. PLSR is capable of considering the correlation between predictor variables and the response variable, making it particularly suitable for chemometrics and spectroscopy applications. PLSR effectively handles multicollinearity among predictor variables. The yield prediction models were developed at five levels 1) using only traditional handmeasured traits, 2) incorporating hand-measured traits and comprehensive leaf angle measurements, 3) including hand-measured traits and comprehensive leaf angle measurements, and vegetation indices.

The performance of the models was evaluated using the coefficient of determination (R2) and mean square error (MSE), comparing the predicted yield to the actual harvest yield. Additionally, a five-fold cross validation approach (80% training set and 20% test set) was employed to assess the robustness of the models. All input variables in the models were scaled at the same level which allows us to compare the coefficients in ridge regression to quantify the

effects of input variables. And variable importance projection (VIP) in PLSR were used to quantify the effects of the input variables on the predictions.

#### Results

Figure 2.3 displays the pattern of canopy changes. The canopy cover increased rapidly from the first to the second fight and continued to grow slowly until the third flight (at 1124.05 GDD), which was close to the anthesis date. The canopy remained closed during the fourth flight. However, by the fifth flight, a decrease in canopy cover was observed in most plots, with even more noticeable variance.

The results of the study revealed a rising trend in canopy-covered NDVI from the first to the third survey, which corresponded to an increase in canopy cover (Figure 2.3 A&B, Table A2.1). However, immediately after reaching its peak during the fourth survey, a decline in NDVI was observed. This decrease was attributed to the withering of bottom leaves despite the canopy cover remaining closed at that time. The NDVI continued to decrease until the last flight. Canopy-covered GCVI exhibited a similar pattern to NDVI but demonstrated greater sensitivity and variance in the last two flights (Figure A2.3).

Three fully expanded leaf counts were conducted at time points of 480.8, 698.1, and 865.7 GDD (Table 1). Variations in fully expanded leaf numbers were observed among different plots for each survey, which led to variations in leaf initiation slope and intercept (Figure 2.3C).

The initial two drone flight surveys were conducted at time points of 580.35 and 810.2 GDD. The correlation between leaf initiation regression and GDD changes during the first two flights was 0.33, indicating that the observed changes in NDVI during the first flights partially captured information on leaf initiation. The overall NVDI of the plot was also influenced by

stand count and other leaf characteristics such as leaf size and angles, which were not measured early in the season.



Figure 2.3: Time series of canopy and leaf characteristics per plot throughout the growing season, over growing degree days (GDD). A: plot mean of NDVI. B: RGB based canopy cover area. C: Fully expanded leaf number; red points are the real recorded expanded leaf numbers for three different time points, while the line is a regression of expanded leaf number over GDD.

## **Leaf Simulation**

The leaf stem angle, mid angle, mid length, end angle, and end length exhibited a wide range of variance. Due to gravity, the order of angle rank was stem angle, mid angle and end angle. The leaf above ear (L2) had a larger stem angle, but a smaller mid angle compared to the leaf below the ear (L1) (Figure 2.4A). For length, the end length was greater than the mid length, and L1 had significantly longer mid length and end length than L2 (Figure 2.4B).

Based on the measured leaf traits for L1 and L2, polynomial functions were used to simulate the leaf curvatures (Figure 2.2), and the leaf length was also integrated from these functions (Figure 2.2). Comparing measured ear leaf length and simulated leaf length, L1 had the longest leaf length, with the largest variance observed. The ear leaf length was in the middle, while L2 had the shortest leaf length (Figure 2.4C).



Figure 2.4: Box plots of measured L1 and L2 leaf measurements, A: three leaf angles, B: length, C: the whole integrated leaf length form polynomial function and measured the ear leaf length.

## Interactions

Hand measured traits showed different distributions (Figure A2.4). Figure 2.5 displays the correlation between leaf and canopy traits. Leaf Area Index (LAI), Total Leaf Number (TLN), Total Leaf Area (TLA), root lodging, stalk lodging, anthesis GDD, silking GDD, Plant Height (PHT), and Ear Height (EHT) were found to be positively correlated with each other and clustered together. The leaf initiation slope exhibited a positive correlation with PHT, TLN, and TLA, but a negative correlation with anthesis GDD and silking GDD. Stand count, used in LAI calculation, showed a negative correlation with TLA and plant height but a positive correlation with stalk lodging.

L1 stem angle was positively correlated with stalk lodging, while L1 mid angle was negatively correlated with EHT. L1 mid length showed a positive correlation with TLN, TLA, PHT, and leaf initiation slope, but a negative correlation with stalk lodging. L1 end angle exhibited a positive correlation with stalk lodging, while L1 end length was negatively correlated with stand count, silking GDD, and stalk lodging, but positively correlated with TLN, TLA, and leaf initiation slope (Figure 2.5).


Figure 2.5: Correlation and clustering of all hand measured traits.

L2 stem angle was negatively correlated with TLA but positively correlated with stalk lodging. L2 mid angle was positively correlated with stand count and leaf initiation intercept, but negatively correlated with TLA and PHT. L2 mid length was positively correlated with leaf initiation slope and TLA, but negatively correlated with ELN, stand count, TLA, anthesis GDD, silking GDD, stalk lodging, and EHT. L2 end angle was positively correlated with stand count, stalk lodging, and EHT, but negatively correlated with leaf initiation slope. L2 end length was negatively correlated with ELN, stalk lodging, and root lodging, but positively correlated with leaf initiation slope (Figure 2.5).

## **Spectral Data**

The NDVI values from different survey dates showed a positive correlation, with neighboring flights displaying higher correlation values. Specifically, the NDVI values from 7/2 and 7/15 were positively correlated with TLA, LAI, stand count, leaf initiation rate, plant L1 and L2 end length, PHT, TLA, and EHT, but negatively correlated with flowering time (Figure 2.6). NDVI on 7/28 was positively correlated with LAI, stand count, stalk lodging, PHT, TLA, EHT, anthesis and silking GDD, TLN, root lodging, and L1 mid length. Similarly, NDVI on 8/12 was positively correlated with LAI, stand count, stalk lodging, PHT, TLA, EHT, anthesis and silking GDD, TLN, root lodging, L1 mid length, and L2 mid length. NDVI on 9/2 was positively correlated with LAI, PHT, TLA, EHT, anthesis and silking GDD, TLN, root lodging, L1 mid length, and L2 mid length. Finally, NDVI in 9/11 displayed a similar correlation pattern to 9/2, except for a negative correlation with stand count (Figure 2.6). Compared to NDVI, the GCVI exhibited a similar correlation pattern with hand measured traits (Figure 2.5). The only difference was some cluster pattern switching and correlation value fluctuation, but no essential change.

Hand measured traits and ridge regressions were used to fit VIs from six flight surveys. The adjust  $R^2$  showed that the NDVI from 7/02, 7/15, and 7/28 can be explained by 0.58, 0.61, 0.49 of the hand measured traits. In these three time points, stand count contributed the most, followed by plant height (Figure 2.7). Hand measured traits only accounted for a very smaller percentile of total variance in 8/12, 9/02, and 9/11. Compared to the first three vegetation indices surveys, stand count did not play a dominant role; plant height and stalk lodging were significantly influenced vegetation indices, but still made a relatively small contribution (Table A2.3). GCVI also gave a similar result, but slightly lower in adjusted  $R^2$  (Figure A2.6).



Figure 2.6: Correlation heat map of NDVI and all hand measurements including the leaf angle traits.



Figure 2.7: The adjusted coefficient of determination of lasso regression models using all hand measured traits to predict different dates of NDVI.

## **Relationship with Yield**

From the correlation heat map, yield was positively corelated with LAI, plant height, leaf initiation slope, and L1 and L2 mid length, but negatively corelated with anthesis GDD, silking GDD, stalking lodging, L2 stem angle, and mid angle. Yield was positively correlated with VIs from every flight survey (Figure 2.8). The correlation pattern of yield for NDVI and GCVI was similar (Figure 2.5). The highest correlated flight was NDVI on Sept. 11. The correlation between yield and GCVI was consistently higher than the correlation between yield and NDVI for each individual flight across the season.



Figure 2.8: Correlation heat map of yield and all hand measurements including the leaf angle traits and NDVI.

# **Yield Prediction**

By using all non-angle traits as input in models to predict yield, ridge regression and PLSR achieved  $R^2$  values of 0.350 and 0.286, respectively (Table 2.3). By adding leaf stem angle, the prediction  $R^2$  increased to 0.356 for ridge regression and 0.280 for PLSR. Incorporating comprehensive leaf angle measurements, the  $R^2$  values improved to 0.390 and 0.280 by using ridge regression, and PLSR. When time series NDVI and GCVI were separately used as inputs in the prediction model, the  $R^2$  values were 0.318 and 0.320 for NDVI using ridge regression and PLSR, respectively. GCVI showed a similar performance with  $R^2$  values of 0.308 and 0.311 (Table 2.3). Combining vegetation indices and hand-measured traits resulted in better predictions than using either indices or hand-measured traits alone. The  $R^2$  values reached 0.483 and 0.478 when combining hand-measured traits and NDVI using ridge regression and PLSR (Table 2.3). When NDVI was switched to GCVI, the  $R^2$  values were 0.465 and 0.464.



Figure 2.9: Sorted mean of variable importance projection (VIP) of all input features (NDVI) in 5-fold PLSR yield prediction.

Table 2.2: Sorted Top 5 Absolute coefficient of input variables in ridge regression of different dates' NDVI.

Rank	7/02	7/15	7/28	8/12	9/02	9/11
1	Stand Count	Stand Count	Stand Count	PHT	PHT	PHT
2	slope	slope	PHT	slope	slope	slope
3	Intercept	PHT	Intercept	Intercept	EHT	Intercept
4	EHT	Intercept	slope	ELN	Intercept	ELN
5	ELN	EHT	EHT	EHT	S Lodging	L2 M Length

Table 2.3: Coefficient of determination of different input variable sets in ridge regression and PLSR for yield prediction.

	Ridge Regression		PLSR	
Inputs:	$R^2$	MSE	$\mathbb{R}^2$	MSE
Non angle hand measurements	0.286	1.140	0.286	1.140
	(0.084)	(0.134)	(0.084)	(0.135)
Stem angle hand measurements	0.280	1.150	0.280	1.150
	(0.108)	(0.172)	(0.107)	(0.172)
All angle Hand measurements	0.278	1.154	0.280	1.150
	(0.122)	(0.202)	(0.107)	(0.172)
NDVIs	0.319	1.089	0.320	1.088
	(0.160)	(0.268)	(0.164)	(0.278)
NDVIs all angle hand measurements	0.484	0.826	0.479	0.834
	(0.121)	(0.207)	(0.128)	(0.220)
GCVIs	0.304	1.13	0.307	1.11
	(0.152)	(0.257)	(0.158)	(0.270)
GCVIs all angle hand measurements	0.458	0.868	0.457	0.870
	(0.124)	(0.216)	(0.126)	(0.222)

Value in () is the standard variation across the five folds.

### Discussion

#### **Vertical Canopy Traits**

Plant height (PHT) is considered a fundamental trait, and previous studies have shown that it is influenced by various factors (Ghimire & Timsina, 2015; Silva et al., 2016). Consistent with these findings, our study also found a positive correlation between PHT and traits such as ELN, TLN, EHT, TLA, and LAI. This group of height-related traits need to be carefully balanced in plant breeding. On one hand, a higher PHT can lead to increased leaf number and total leaf area, which promotes photosynthesis and solar radiation use efficiency (Edmeades & Daynard, 1979). On the other hand, a higher PHT can also result in lodging, which negatively affects yield (Q. Zhang et al., 2014). In this study, the positive effects of increased leaf area from PHT outweighed the negative effects of lodging induced by PHT.

### Leaf angle and plant density

In this experiment, we planted a fixed number of seeds in each plot and measured stand count, green snap count, and stalk lodging at different time points (Table 1). Stand count, which was measured earliest which was less affected by environmental factors, as plant density can induce spatial competition and affect light availability. Although excessively high plant density can lead to lodging and green snap, these effects were minimal in our study. Under high density, plants tended to have smaller angles, which is consistent with the positive correlation between stalk lodging and L1 and L2 stem angle. The high canopy cover and vegetation indices observed on 7/28, as well as the variance of stand count, confirmed the plasticity of leaf angle in response to available plants. Compared to traditional stem angle measurements, the mid-length was found to be more correlated with yield. A longer leaf mid-length can indicate less stalk lodging but larger TLA, both of which are beneficial for yield. However, the end length was not as strongly

correlated with yield. Combining simulated leaf curvature and correlation analysis, we can conclude that a stiff leaf (small angle and mid-length) is favorable for TLA and yield, and the descending part of the leaf is as important as the ascending part which is agreed by previous studies (G. Liu et al., 2017; Mason & Zuber, 1976).

Since L1 and L2 were measured on the same plant, their angles and lengths were closely related. However, L1 generally had a larger angle compared to L2, while the mid-length of L2 was larger than that of L1. Furthermore, although the total length of L1 was slightly greater than that of L2, this could be attributed to the senescence sequence or angle distribution for improved photosynthetic efficiency (Niinemets, 2007). Additionally, there was some ambiguity in determining the position of L2, as some plants had more than one ear. When comparing the length of simulated L1 and L2, as well as the measured ear leaf length, it was observed that L1 exhibited the longest length. This length distribution was consistent with previous studies, L1 position were generally closer to 2/3 of total leaf number (Birch et al., 1997).

LAI represents the leaf area per unit field area and combines the vertical and horizonal dimensions of the canopy. In this study, the total possible LAI for the season was estimated using a power function that considered the leaf number, leaf area, and the number of plants in the field. Although this LAI value could not have been observed on any specific date throughout the season as it does not consider senescence, it represents the overall potential leaf area within the canopy. This estimate of total LAI was positively correlated with yield.

Positive correlations have been reported between maize yield and stalk strength in previous studies (Singh, 1970). A strong stalk provides benefits such as improved lodging resistance and efficient assimilate transport.

## **Time Series of Canopy**

During the early stages, NDVI primarily reflected canopy expansion. Therefore, the NDVI values from 7/02 and 7/15 captured information about leaf initiation slope and stand count. Additionally, the negative correlation between early vegetation indices and flowering indicated that rapidly growing plants in early stages required fewer GDD to reach flowering, highlighting the efficiency of growth. Flowering is a key step in the transition from vegetation growth to the reproductive stage (Iqbal et al., 2017; Wellmer & Riechmann, 2010). The transition in correlation between flowering GDD and vegetation indices occurred around 7/28, which corresponds to when most plots had finished flowering. This reflects the senescence of lower-position leaves in plots that flowered early. Leaf angle measurements were conducted around the same time as the late-season vegetation indices. Stiff leaves with small angles and mid-lengths led to higher vegetation indices.

NDVI and GCVI showed positive correlations with yield throughout the season. The variance of NDVI and GCVI was higher in the late season, corresponding to stronger correlation with yield compared to the early season. The variance of NDVI and GCVI reflects the greenness variability of canopy (Burke & Lobell, 2017; Mkhabela et al., 2005). Higher NDVI and GCVI values in late season, closer to maturity, were associated with genotypes that tended to have higher yields. In contrast, flowering GDD showed a negative correlation with yield. Based on the flowering and late season drone vegetation indices, a longer reproductive period had more potential to achieve high yield. In other words, the ideal variety would flower early, but stay green and have a long grain-filling period. This information could be used for breeders to make efficient selections.

Stand count and stalk lodging directly determined the number of harvestable plants in each plot, hence the positive and negative relationships found between stand count and stalk lodging with yield. Root lodging and green snap had little effect on yield, though they also can determine the number of harvestable plants. This inconsistency may be root lodging and green snap tend to occur in specific weather conditions, especially heavy rains and wind (Gardiner et al., 2016), that were not observed in this growing season.

#### **Yield Prediction and Relationships**

To mitigate the negative effects of multicollinearity in the time-series vegetation indices, we selected regression algorithms that can handle this issue. The addition of leaf angle data had little impact on the prediction accuracy. In this study, plant density (stand count, stalk lodging) had the most significant influence on leaf angle. Compared to one-time in-field hand-measured traits, multi-temporal aerial VIs measurements provided more insights into canopy dynamics. VIs values obtained after flowering showed a higher correlation with yield compared to those obtained before flowering. The highest yield correlation was observed on the September 11th flight. The prediction accuracy using vegetation indices alone was similar to that using handmeasured traits alone. However, combing the hand-measured traits and vegetation indices resulted in superior yield prediction performance, indicating complementary information between these two data sources. Analysis of VIP of PLSR and coefficients in ridge regression, stalk lodging and stand count also played important roles in yield prediction model (Figure 2.9, Figure A2.7, Table A2.4). Stand count was highly correlated with the early VIs, but stalk lodging showed less correlation even later in the season. There are two possible explanations for this phenomenon. First, the stalk lodging occurred very late in the growing season (Robertson et al., 2016), which was not covered by flight surveys in this study; second, as green tissue turned

yellow, VIs lost their ability to detect differences in plant status (Gu et al., 2008). Upon comparing the VI on 7/15 and 7/28, it was observed that there was a saturation present in the NDVI on 7/28, whereas no such striation was observed in GCVI. But the NDVI model achieved a slightly better yield prediction accuracy. This finding suggests that the saturation of NDVI during the growth peak has minimal impact on the accuracy of the model in predicting yield. The stable and sensitive nature of NDVI contributes to this advantage. More representativeness of agronomy traits was found in NDVI in correlation analysis and variance analysis in model fitting, compared to GCVI.

Most yield prediction studies using vegetation indices have focused on different treatment conditions, such as different nitrogen levels, water availability (well-watered or drought conditions), and different sowing dates. In these scenarios, the majority of observed trait variance comes from the treatment, increasing the prediction accuracy (Herrmann et al., 2020; Wu et al., 2019). By contrast, in our study, all data were collected from a single field environment, where plot-level yield variance was primarily influenced by similar background genotypes and random microenvironments in the field. Several studies used tens of VIs derived by several reflectance bands in machine learning models, but the results are not very promising (Fan et al., 2022). The similar prediction accuracy from with only using timeseries single VI. Except for NIR, the reflectance of red, green, blue and rededge were highly clustered and correlated in all flight surveys. VIs derived from several bands as input for machine learning model could only provide redundant information. A very slight difference of NDVI and GCVI in correlation analysis and yield prediction also confirmed this assumption. Furthermore, the varieties used in this study were highly uniform, indicating our mixed canopy prediction approach could be more reliable under real breeding scenarios to select the superior genotypes

with high yield. The high prediction accuracy for yield underscores the importance of canopy traits. The inclusion of hand measured traits in our best yield prediction models highlights the need for more high-throughput phenotyping methods to capture latent canopy traits.

## Conclusion

The maize leaf canopy is composed of a set of interrelated traits that reflect the physiological growth and development strategies of maize. These complex interrelationships maintain a dynamic balance of these traits. The strong correlation between yield and canopy traits emphasizes the importance of understanding leaf and canopy characteristics in the decision-making process of maize breeding. Our study provides evidence that plot-based high-throughput canopy phenotyping can capture some of the hand-measured leaf traits at the plant level. High-throughput phenotyping has great potential for estimating critical agronomic traits in large maize breeding populations or as a high-dimensional integrative phenotypic trait that reflects the true status of the canopy. By using leaf-scale to canopy-scale traits and combining time-series vegetation indices, we achieved promising accuracy in yield prediction. This approach not only benefits yield estimation but also provides evidence for selecting multiple traits in the breeding process.

#### REFERENCES

- Anderson, S. L., & Murray, S. C. (2020). R/UAStools: Plotshpcreate: Create multi-polygon shapefiles for extraction of research plot scale agriculture remote sensing data. *Frontiers in Plant Science*, 11, 511768.
- Birch, C. J., Hammer, G. L., & Rickert, K. G. (1997). Improved methods for predicting individual leaf area and leaf senescence in maize (Zea mays). *Australian Journal of Agricultural Research*, 49(2), 249–262.
- Burke, M., & Lobell, D. B. (2017). Satellite-based assessment of yield variation and its determinants in smallholder African systems. *Proceedings of the National Academy of Sciences*, *114*(9), 2189–2194.
- Chen, T.-W., Cabrera-Bosquet, L., Alvarez Prado, S., Perez, R., Artzet, S., Pradal, C., Coupel-Ledru, A., Fournier, C., & Tardieu, F. (2019). Genetic and environmental dissection of biomass accumulation in multi-genotype maize canopies. *Journal of Experimental Botany*, 70(9), 2523–2534. https://doi.org/10.1093/jxb/ery309
- Danilevicz, M. F., Bayer, P. E., Boussaid, F., Bennamoun, M., & Edwards, D. (2021). Maize yield prediction at an early developmental stage using multispectral images and genotype data for preliminary hybrid selection. *Remote Sensing*, *13*(19), 3976.
- Du, X., Wang, Z., Lei, W., & Kong, L. (2021). Increased planting density combined with reduced nitrogen rate to achieve high yield in maize. *Scientific Reports*, 11(1), 358.
- Duncan, W. G. (1971). Leaf Angles, Leaf Area, and Canopy Photosynthesis1. *Crop Science*, *11*(4), cropsci1971.0011183X001100040006x. https://doi.org/10.2135/cropsci1971.0011183X001100040006x
- Dwyer, L. M., & Stewart, D. W. (1986). Leaf Area Development in Field-Grown Maize1. Agronomy Journal, 78(2), 334–343. https://doi.org/10.2134/agronj1986.00021962007800020024x
- Edmeades, G., & Daynard, T. (1979). The development of plant-to-plant variability in maize at different planting densities. *Canadian Journal of Plant Science*, *59*(3), 561–576.
- Fan, J., Zhou, J., Wang, B., de Leon, N., Kaeppler, S. M., Lima, D. C., & Zhang, Z. (2022). Estimation of Maize Yield and Flowering Time Using Multi-Temporal UAV-Based Hyperspectral Data. *Remote Sensing*, 14(13), 3052.
- Francis, C. A., Rutger, J. N., & Palmer, A. F. E. (1969). A Rapid Method for Plant Leaf Area Estimation in Maize (Zea mays L.)1. *Crop Science*, 9(5), cropsci1969.0011183X000900050005x. https://doi.org/10.2135/cropsci1969.0011183X000900050005x

- Gamon, J. A., Field, C. B., Goulden, M. L., Griffin, K. L., Hartley, A. E., Joel, G., Penuelas, J., & Valentini, R. (1995). Relationships between NDVI, canopy structure, and photosynthesis in three Californian vegetation types. *Ecological Applications*, 5(1), 28–41.
- Gardiner, B., Berry, P., & Moulia, B. (2016). Wind impacts on plant growth, mechanics and damage. *Plant Science*, 245, 94–118.
- Ghimire, B., & Timsina, D. (2015). Analysis of yield and yield attributing traits of maize genotypes in Chitwan, Nepal. *World Journal of Agricultural Research*, *3*(5), 153–162.
- Gu, Y., Hunt, E., Wardlow, B., Basara, J. B., Brown, J. F., & Verdin, J. P. (2008). Evaluation of MODIS NDVI and NDWI for vegetation drought monitoring using Oklahoma Mesonet soil moisture data. *Geophysical Research Letters*, 35(22).
- Hashimoto, N., Saito, Y., Maki, M., & Homma, K. (2019). Simulation of reflectance and vegetation indices for unmanned aerial vehicle (UAV) monitoring of paddy fields. *Remote Sensing*, 11(18), 2119.
- Herrmann, I., Bdolach, E., Montekyo, Y., Rachmilevitch, S., Townsend, P. A., & Karnieli, A. (2020). Assessment of maize yield and phenology by drone-mounted superspectral camera. *Precision Agriculture*, 21, 51–76.
- Iqbal, N., Khan, N. A., Ferrante, A., Trivellini, A., Francini, A., & Khan, M. (2017). Ethylene role in plant growth, development and senescence: Interaction with other phytohormones. *Frontiers in Plant Science*, 8, 475.
- Li, J., Xie, R. Z., Wang, K. R., Hou, P., Ming, B., Zhang, G. Q., Liu, G. Z., Wu, M., Yang, Z. S., & Li, S. K. (2018). Response of canopy structure, light interception and grain yield to plant density in maize. *The Journal of Agricultural Science*, 156(6), 785–794. Cambridge Core. https://doi.org/10.1017/S0021859618000692
- Li, L., Zhang, Q., & Huang, D. (2014). A Review of Imaging Techniques for Plant Phenotyping. Sensors, 14(11), 20078–20111. https://doi.org/10.3390/s141120078
- Liu, G., Hou, P., Xie, R., Ming, B., Wang, K., Xu, W., Liu, W., Yang, Y., & Li, S. (2017). Canopy characteristics of high-yield maize with yield potential of 22.5 Mg ha– 1. *Field Crops Research*, 213, 221–230.
- Liu, T., Chen, J., Wang, Z., Wu, X., Wu, X., Ding, R., Han, Q., Cai, T., & Jia, Z. (2018). Ridge and furrow planting pattern optimizes canopy structure of summer maize and obtains higher grain yield. *Field Crops Research*, 219, 242–249. https://doi.org/10.1016/j.fcr.2018.02.012
- Mason, L., & Zuber, M. S. (1976). Diallel analysis of maize for leaf angle, leaf area, yield, and yield components 1. *Crop Science*, *16*(5), 693–696.

- Meiyan, S., Mengyuan, S., Qizhou, D., Xiaohong, Y., Baoguo, L., & Yuntao, M. (2022). Estimating the maize above-ground biomass by constructing the tridimensional concept model based on UAV-based digital and multi-spectral images. *Field Crops Research*, 282, 108491. https://doi.org/10.1016/j.fcr.2022.108491
- Mkhabela, M. S., Mkhabela, M. S., & Mashinini, N. N. (2005). Early maize yield forecasting in the four agro-ecological regions of Swaziland using NDVI data derived from NOAA's-AVHRR. *Agricultural and Forest Meteorology*, *129*(1–2), 1–9.
- Nguy-Robertson, A., Gitelson, A., Peng, Y., Viña, A., Arkebauer, T., & Rundquist, D. (2012). Green leaf area index estimation in maize and soybean: Combining vegetation indices to achieve maximal sensitivity. *Agronomy Journal*, *104*(5), 1336–1347.
- Niinemets, U. (2007). Photosynthesis and resource distribution through plant canopies. *Plant, Cell & Environment*, *30*(9), 1052–1071.
- Ozyavuz, M., Bilgili, B., & Salici, A. (2015). Determination of vegetation changes with NDVI method. *Journal of Environmental Protection and Ecology*, *16*(1), 264–273.
- Pendleton, J., Smith, G., Winter, S., & Johnston, T. (1968). Field investigations of the relationships of leaf angle in corn (zea mays l.) to grain yield and apparent photosynthesis 1. Agronomy Journal, 60(4), 422–424.
- Robertson, D. J., Lee, S. Y., Julias, M., & Cook, D. D. (2016). Maize stalk lodging: Flexural stiffness predicts strength. *Crop Science*, *56*(4), 1711–1718.
- Sher, A., He, L., Zhang, S., Li, J., & Song, Y. (2016). Analysis and characterisation of interplant competition on maize canopy morphology for modelling. 189–193.
- Sher, A., Khan, A., Ashraf, U., Liu, H. H., & Li, J. C. (2018). Characterization of the effect of increased plant density on canopy morphology and stalk lodging risk. *Frontiers in Plant Science*, 9, 1047.
- Sher, A., Khan, A., Cai, L. J., Ahmad, M. I., Asharf, U., & Jamoro, S. A. (2017). Response of maize grown under high plant density; performance, issues and management-a critical review. Adv. Crop Sci. Technol, 5(3), 1–8.
- Silva, T. N., Moro, G. V., Moro, F. V., Santos, D. M. M. dos, & Buzinaro, R. (2016). Correlation and path analysis of agronomic and morphological traits in maize. *Revista Ciência Agronômica*, 47, 351–357.
- Singh, T. (1970). Association between certain stalk traits related to lodging and grain yield in maize (Zea mays L.). *Euphytica*, *19*(3), 394–397.
- Tamás, A., Kovács, E., Horváth, É., Juhász, C., Radócz, L., Rátonyi, T., & Ragán, P. (2023). Assessment of NDVI Dynamics of Maize (Zea mays L.) and Its Relation to Grain Yield in a Polyfactorial Experiment Based on Remote Sensing. *Agriculture*, 13(3). https://doi.org/10.3390/agriculture13030689

- Tokatlidis, I. S., Has, V., Melidis, V., Has, I., Mylonas, I., Evgenidis, G., Copandean, A., Ninou, E., & Fasoula, V. A. (2011). Maize hybrids less dependent on high plant densities improve resource-use efficiency in rainfed and irrigated conditions. *Field Crops Research*, 120(3), 345–351. https://doi.org/10.1016/j.fcr.2010.11.006
- Tokatlidis, I. S., & Koutroubas, S. D. (2004). A review of maize hybrids' dependence on high plant populations and its implications for crop yield stability. *Field Crops Research*, 88(2), 103–114. https://doi.org/10.1016/j.fcr.2003.11.013
- Verhulst, N., Govaerts, B., Nelissen, V., Sayre, K. D., Crossa, J., Raes, D., & Deckers, J. (2011). The effect of tillage, crop rotation and residue management on maize and wheat growth and development evaluated with an optical sensor. *Field Crops Research*, 120(1), 58–67. https://doi.org/10.1016/j.fcr.2010.08.012
- Viña, A., Gitelson, A. A., Nguy-Robertson, A. L., & Peng, Y. (2011). Comparison of different vegetation indices for the remote assessment of green leaf area index of crops. *Remote Sensing of Environment*, 115(12), 3468–3478.
- Wang, Q., Adiku, S., Tenhunen, J., & Granier, A. (2005). On the relationship of NDVI with leaf area index in a deciduous forest site. *Remote Sensing of Environment*, 94(2), 244–255.
- Wellmer, F., & Riechmann, J. L. (2010). Gene networks controlling the initiation of flower development. *Trends in Genetics*, 26(12), 519–527.
- WILSON, T. D., BROOK, R. M., & TOMLINSON, H. F. (1998). INTERACTIONS BETWEEN NÉRÉ (PARKIA BIGLOBOSA) AND UNDER-PLANTED SORGHUM IN A PARKLAND SYSTEM IN BURKINA FASO. *Experimental Agriculture*, 34(1), 85–99. Cambridge Core. https://doi.org/10.1017/S0014479798001069
- Wu, G., Miller, N. D., De Leon, N., Kaeppler, S. M., & Spalding, E. P. (2019). Predicting Zea mays flowering time, yield, and kernel dimensions by analyzing aerial images. *Frontiers in Plant Science*, 10, 1251.
- Zhang, H., Zhang, C., Sun, P., Jiang, X., Xu, G., & Yang, J. (2022). Optimizing planting density and nitrogen application to enhance profit and nitrogen use of summer maize in Huanghuaihai region of China. *Scientific Reports*, 12(1), 2704.
- Zhang, Q., Zhang, L., Evers, J., van der Werf, W., Zhang, W., & Duan, L. (2014). Maize yield and quality in response to plant density and application of a novel plant growth regulator. *Field Crops Research*, 164, 82–89.

# APPENDIX



Figure A2.1: The polynomial function used to fit the leaf curvature.



Figure A2.2: The correlation and clustering of single bands for different flight dates. A: 7/02; B: 7/15: C: 7/28; D: 8/12; E: 9/02; F: 9/11.



Figure A2.3: Time series of plot mean vegetative indices over GDD for six flights. A: Whole plot GCVI, B: Canopy covered NDVI, C: Canopy covered GCVI.



Figure A2.4: Density distribution of all hand measured traits, with dashed line showing mean value.



Figure A2.5: Correlation heat map of GCVI and all hand measurements including the leaf angle traits.



Figure A2.6: The adjusted coefficient of determination using all hand measurand traits and lasso regression for different dates of NDVI.



Figure A2.7: Sorted mean of variable importance projection (VIP) of all input features (GCVI) in 5-fold PLSR yield prediction.

	7/02	7/15	7/28	8/12	9/02	9/11
r	0.6692186	0.8146454	0.6641018	0.4837311	0.5484422	0.7864576
Table A2.	.2: The correlation	of potential t	otal leaf area	index and NE	OVI.	
	702	715	728	812	902	911
r	0.4981308	0.5953713	0.5805949	0.250449	0.1699696	0.1777904

Table A2.1: The correlation of canopy cover and NDVI.

Rank	7/02	7/15	7/28	8/12	9/02	9/11
1	Stand Count	Stand Count	Stand Count	PHT	PHT	PHT
2	slope	slope	PHT	slope	slope	slope
3	Intercept	PHT	Intercept	Intercept	EHT	Intercept
4	EHT	Intercept	slope	ELN	Intercept	ELN
5	ELN	EHT	EHT	EHT	S Lodging	L2 M Length

Table A2.3: Sorted Top 5 absolute coefficient of input variables in ridge regression of different dates' NDVI.

Table A2.4: Mean coefficient of input variables in 5-fold ridge regression of all-inclusive model (NDVI or GCVI).

Variables	NDVI	GCVI
7/02	-0.035	-0.006
7/15	0.347	0.210
7/28	-0.089	-0.228
8/12	0.330	0.265
9/02	-0.252	0.034
9/11	0.575	0.440
TLA	0.010	-0.002
Stand.Count	0.279	0.442
AnthesisGDD	-0.185	-0.118
SilkingGDD	-0.131	-0.105
GreenSnap	-0.050	-0.036
Stalk Lodging	-0.488	-0.435
Root Lodging	-0.165	-0.153
PHT	0.240	0.301
EHT	0.013	0.003
slope	-0.119	-0.165
Intercept	-0.057	-0.078
PlantL1.Stem.Angle	-0.032	-0.020
PlantL1.Mid.Angle	0.137	0.100
PlantL1.Mid.Length	-0.050	-0.044
PlantL1.End.Angle	0.001	0.030
PlantL1.End.Length	-0.084	-0.034
PlantL2.Stem.Angle	0.012	-0.016
PlantL2.Mid.Angle	-0.042	-0.002
PlantL2.Mid.Length	0.006	-0.023
PlantL2.End.Angle	-0.053	-0.083
PlantL2.End.Length	0.020	0.019
endangle1	-0.002	0.003
endangle2	-0.071	-0.075

# CHAPTER 3: A CASE STUDY OF CROSS-SUBPOPULATION GENOMIC PREDICTION

### Abstract

Genomic selection has been used in breeding programs to speed up the breeding cycle. But predicting the performance of newly introduced genotypes or population is still very challenging. In this paper we utilized a small tested cross NAM population from Genomes to Field project to conduct a case study for cross-subpopulation genomic prediction. The prediction accuracy of cross-subpopulation prediction was worse than randomly sampling from all genetics pools, and the prediction accuracy within cross-subpopulation scenarios varied by trait and subpopulation. We analyzed these differences and found possible explanations. A dominance relationship matrix, Gaussian kernel-based relationship, and LD adjusted methods were used to explore their effect on prediction accuracy in cross-subpopulation scenarios.

## Introduction

Maize is a member of the grass family and is native to Mexico. It plays an essential role in global agriculture, serving as a vital source of food, livestock feed, and bioenergy. With the progress of breeding technology, maize yields have increased from 30 bushels in 1930s to around 180 bushels per acre in 2020 (Schlenker 2020). Thanks to advancements in genotyping technologies and computational methods, the breeding method has evolved from selecting individuals based on specific, well-understood genetic markers linked to a known trait to utilizing comprehensive genome-wide datasets. This shift enables the prediction of genetic values for multiple traits, proving particularly effective for complex, polygenic traits. The approach leverages extensive genomic datasets to forecast an individual's genetic predisposition or potential for a specific trait.

Various statistical genomic prediction models have been developed and applied in

breeding. The Genomic Best Linear Unbiased Prediction (G-BLUP) model uses genome-wide makers to estimate a sample covariance matrix (Meuwissen *et al.* 2001; Habier *et al.* 2007; VanRaden 2008). The relationship covariance matrix can be additive, dominant, and Euclidean distance based. Fitting regressions with whole-genome markers to estimate marker effects is an alternative approach in genomic prediction, leveraging comprehensive genomic information for accurate trait predictions. In marker based genomic prediction, ridge regression, lasso regression, and several Bayesian regression approaches with varying prior density assumptions (BayesA, Bayes B, BayesC, etc.) have been widely used (Isik *et al.* 2017; Robertsen *et al.* 2019).

Genomic prediction has emerged as a widely utilized tool in multiple fields including agriculture, specifically in the improvement of complex traits in maize. Its effectiveness lies in its ability to utilize genotypic data for selection purposes, regardless of the season or timing, even before direct phenotypic selection can take place. Recent studies have demonstrated the successful application of genomic prediction in enhancing traits like yield and disease resistance in maize with high efficacy compared with traditional breeding methods (Jannink *et al.* 2010; Lorenz 2013). Additionally, a study highlighted that genomic prediction can boost the genetic gain in yield by as much as 30% within maize breeding programs (Crossa *et al.* 2017).

In genomic prediction, it is often advisable to reduce the total number of markers utilized, both for the sake of computational efficiency and to avoid overfitting. Linkage disequilibrium (LD) is the non-random association of alleles at different loci within a given population (Song *et al.* 2021). One approach to pruning uses linkage disequilibrium (LD) patterns to select representative and non-redundant genetic markers. LD pruning is commonly used in genomic studies to reduce marker redundancy and computational complexity. In GWAS, LD pruning has been used to improve the statistical power and reduce the risk of overfitting. However, LD

pruning can also remove the true associated SNPs and increase the type ii error rate in GWAS. In genomic prediction, there are few studies focused explicitly on the effects of LD pruning (Ye *et al.* 2019; Song and Hu 2022).

Elite germplasm that has experienced both adaptation and selection processes has shown less diversity and higher LD compared to the whole maize genetic pool (Maccaferri *et al.* 2005). Germplasm introgression is critical to plant breeding and genetic improvement, as it enriches valuable genetic diversity and introduces new beneficial traits. Because of the narrowness of current elite genetic pools, genomic prediction has inferior performance when predicting new germplasm introgressions than predicting within the pool. In this study, we utilized a small NAM population from the 2018 Genome to Fields (G2F) project for cross-subpopulation genomic prediction as a case study to provide insight into genomic prediction in germplasm introgression.

### **Materials and Methods**

#### Population

In the 2018 G2F project, hybrid materials were developed from a PHW65 NAM population which has three founder families, including double haploids derived from PHW65/PHN11, PHW65/Mo44, & PHW65/MoG, then test crossed with PHT69. PHW65, PHN11, and PHT69 are expired plant variety protection (ex-PVP) lines from Pioneer Hybrid seed company; Mo44 and MoG are public lines developed in Missouri.

Because of seed limitations and to increase representation of the full set across sites, the field was planted as a modified randomized complete block design. Some hybrids were replicated twice at a site, while some were only assigned with one rep. All inbred parents represented in the 2018 G2F hybrids were sequenced at ~5x coverage. Pools of 24 samples were sequenced on a HiSeq X Ten lane at Novogene. 4.2 million genic SNPs with B73 AGPv4

coordinates were called via Practical Haplotype Graph (information provided by Genomes to Fields collaborative). SNPs were imputed in BEAGLE and roughly cleaned by using MAF<5%. The SNPrelate package in R was used to conduct the LD pruning on the preprocessed SNPs. Final thresholds of LD pruning were set to 0.99, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2 and 0.1.

#### Field data measurements

In this study, we used field data from the 2018 Michigan and Wisconsin locations. Michigan's site had 2 blocks with 250 plots each, while Wisconsin site had 2 blocks with 400 plots each. In these two locations, the total number of genotypes were 460, which consisted of 168 MoG families, 185 PHN11 families, 90 Mo44 families and 18 yellow stripe (check) hybrids (Figure 3.1).

In the field, all data collection was based on Genomes to Fields standard operating procedures (https://www.genomes2fields.org/resources/). For each plot, we noted anthesis date as the date at which 50% of the plants' anthers were exuded on over half of the main tassel spike. Similarly, the silking date is when 50% of plants' silks have emerged within each plot. We counted the number of plants broken between the ground level and the top ear node before flowering (green snap). We choose two representative plants in each plot (one in each of the two rows) to measure the plant-scale traits. For these two plants, we measured the distance from ground soil line to ligule of the flag leaf as the plant height, and from the ground soil line to the top ear bearing node as ear height. The height from ear to flag leaf height we designated as the differential height. The yield was reported as the real harvested weight for a single plot without any conversion.

### **Genomic prediction**

Because of the unbalanced experimental design, we treated all genotypes as fixed effects

and location as random effects and calculated the best linear unbiased estimator (BLUE) values represent the genotypic in above traits.

The BLUE values for all genotypes were estimated via mixed linear model as:

$$y = G + E + GE + \varepsilon$$

Where y is the measured trait value, the G = Genotype, E = Environmental factors, and the interaction of genotype and environmental factors were treated as random effects. The calculated BLUE values for each trait were used in genomic prediction and genome wide association studies (GWAS). A genomic relationship matrix was used to represent the proportion of genetic similarity among individuals and was used to predict the BLUE for specific traits of individuals:

$$\hat{Y} = \mu + \hat{G} + \varepsilon$$

All genotypes in this study were uniform hybrids, and the genomic relationship matrices was estimated by additive effects (A), dominance effects (D), and Gaussian kernel (GK) by averaging the Euclidean distance of all SNPs. To explore the effects of LD and SNP number on across-population genomic prediction, unpruned markers were implemented for all three genomic relationship matrices for prediction. In prediction, we held fixed the yellow stripes in the training set and performed a rotation of three NAM subpopulations, treating two as the training set and the remaining one as the testing set. The correlation of the BLUE value and the predicted value was reported as prediction accuracy. The A and Dominant can be written as:

> $A_i \sim N(0, \sigma_A^2 A)$  $D_i \sim N(0, \sigma_D^2 D)$

The entries of the Gaussian kernel were computed as:

$$GK = \exp(-hd_{ii}^{\prime 2})$$

where dii' is the Euclidean distance between the individuals ith and i' th (i=1,2,....) given by the markers. The scaling factor is determined by the median of distances between markers, represented as h=1.The methodology mentioned is detailed in the research conducted by Crossa et al. (2010). The theoretical underpinnings of the Gaussian kernel, as applied in the context of kernel averaging for the Reproducing Kernel Hilbert Space (RKHS), were elucidated in the study by de los Campos et al. (2010) and Crossa et al., (2010). All model fitting and data analysis were platformed in R. All genomic prediction models were run in the BGLR package with 12000 iterations after 2000 burn-ins (Pérez and de Los Campos 2014).

## SNP Fst:

Fixation Index (Fst) quantifies the degree of genetic variation between subpopulations relative to the total genetic variation in the entire population. Higher values indicate greater genetic differentiation between subpopulations, lower values represent the opposite. The Fst values were calculated four times: first using the full combined set of three NAM subpopulations (PHN11vs MO44 vs MOG), and then by separating out test versus training set (leaving one population out) for each across-population's prediction.

#### **Results and Discussion**

The first two principal components (PC) explained 12.4% and 9.1% of the variance, respectively. In total, 29.4% of the variance could be explained by the first three PCs (Figure 3.1 & S Figure 1). The rest of the PCs explained relatively small variances. The three small NAM populations could be distinguished by the first principal component. In the biplot, Mo44 and MoG subpopulations showed some overlap, while PHN11 was relatively far away from these two families. The yellow stripe and three NAM subpopulations exhibited large differences in the second principal component with larger variance.



Figure 3.1: Numbers of genotypes in each sub-population and the distribution of subpopulations across the first two principal components.

After basic SNP quality control, the remaining SNP number was 263,273. We used LD levels of 0.99, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, and 0.1 as thresholds to prune the SNPs, retaining 92836, 44103, 28314, 20898, 15979, 12348 9634, 7624, 5681, and 3897 SNPs, respectively (Figure 3.2).





Michigan and Wisconsin showed environmental effects on all traits for each sub population. Overall, the plants in the Wisconsin location had higher ear height, height above ear, greater yield, and later anthesis and silking (Figure 3.3). The PHN11 subpopulation exhibited shorter height traits, but otherwise the phenotypic trait differences across families were not obvious. Long tails were observed in PHN11 families in Silking Days and yield (Figure 3.3).



Figure 3.3. Violin plots for Ear height, Plant height, Height Above Ear, Anthesis Day (D) after planting, Silking Day (D) after planting, plot yield for MO44, PHN11, MOG and Yellow Stripes (YS) in two locations.

#### **Random Sample Prediction vs Cross Subpopulation Prediction**

We utilized two of the three NAM families along with the yellow stripe hybrids to predict the traits of the remaining subpopulation. This approach is similar to what would occur in crosspopulation prediction in the germplasm introgression process (Figure 3.4). As a comparison, we randomly selected the same number of genotypes (90, 168, 185) as each of the three NAM subpopulations from the entire set and repeated this sampling 50 times to eliminate sampling bias. Compared with mean prediction accuracy for random sampling, using two NAM subpopulations along with the yellow stripe hybrids to predict the remaining subpopulation had significantly lower prediction accuracy across traits (Figure 3.5). Meanwhile, the randomly sampled test set showed very limited effects of test set size on mean prediction accuracy across 6 traits. In cross-subpopulation prediction, the prediction accuracy varied by trait and subpopulation. The prediction accuracy for ear height in the Mo44, MoG, and PHN11 subpopulations, respectively, were 0.17, 0.33, and 0.38. In terms of plant height, prediction of the Mo44 subpopulation exhibited the lowest accuracy, while the accuracies of MoG and PHN11 were comparable. The Diff HT prediction was a little bit different, as the accuracy of prediction was 0.17 for Mo44, 0.15 for MoG, and 0.07 for PHN11. The prediction accuracy of families showed a pattern for Anthesis Days and Silking Days (D). Anthesis D was slightly more predictable than silking D. PHN11 had the highest prediction accuracy, with Mo44 and MoG following behind. Yield is considered a highly complex trait which is hard to predict. In yield prediction, the Mo44 subpopulation achieved the highest accuracy, the PHN11 ranked second, and the MoG subpopulation was not predictable.

It has been reported that the test size may not affect prediction accuracy. The prediction accuracy of random samples of different numbers of test sizes also agrees with this claim. However, huge prediction accuracy gaps exist between random samples and cross-subpopulation prediction under the same training-test ratio. The phenotypic traits were biased for population structure, but not very significantly. Conserved allele frequency was observed in different population sets which could give incorrect marker effect estimates under cross-subpopulation prediction conditions. Overall, narrow genetic diversity and small size of the NAM population made this situation even more severe. In random sampling from a genetics pool, the population bias was largely eliminated. When performing cross-population genomic prediction, family/subpopulation factors cannot be easily interpreted or predicted by SNPs. In this study,

markers and genomic relationship largely explain the Mendelian sampling. Simply migrating the maker effects or reconstructing the genomic relationship matrix for a new family/population could result in a worse prediction accuracy than expected. Preliminary population and family analysis added into the genomic relationship or marker effects model might be a better way to avoid this issue.



Figure 3.5: The comparison of mean prediction accuracy using random sampling with 50 interactions, using the two NAM subpopulations and yellow stripes to predict the remaining NAM subpopulation (dash). The test sizes of 90, 168 and 185 equal the number of genotypes in the MO44, MOG and PHN11 subpopulations, respectively.

#### **Fst analysis**

Fst is commonly used in population genetics to measure genetic differentiation (Nei 1986). The mean of Overall Fst for the three NAM subpopulations was 0.0719 (Figure A3.1 & A3.2). This low mean Fst value of three NAM populations had quite low genetic differentiation, which is caused by shared ancestors. Fst values corresponding to the above three training and testing sets were calculated. When using Mo44, PHN11, and MoG as the predicted (test) set, the mean of SNP Fst was 0.024, 0.039 and 0. 041, respectively (Table 3.1). Fst results in this cross-subpopulation prediction revealed that Mo44 exhibited the least difference from the training set,

while PHN11 ranked second, and MOG displayed the largest differences. But the distribution of the SNP's Fst values for each of the three subpopulations as prediction and training set was quite different from mean Fst. When using Mo44 as the predict set, we observed the lowest average mean Fst compared to the training set. That subpopulation also had the largest count of extremely high-Fst SNPs, followed by MoG and PHN11 (Figure 3.6).



Table 3.1: The mean Fst with each subpopulation as the predicted group.Predict MO44Predict MOGPredict PHN11



In terms of predicting height above ear and yield, lower Fst values between the test and training sets were associated with higher prediction accuracy. However, these patterns were not observed for the remaining traits examined in this study. Several previous studies from simulated data and real data have shown higher reliability for genomic prediction when the prediction population (test set) was more closely related to the training population (Makowsky *et al.* 2011; Slavov *et al.* 2014; Wu *et al.* 2015). But in this study, compared with the mean SNP Fst value between predict and training set, only yield followed the above rules. MOG emerged as the most

predictable subpopulation for ear height with highest mean Fst value, compared to the other two subpopulations. In our study, prediction accuracy of traits exhibited varied patterns by mean Fst. Although quantitative traits are controlled by a large number of small effects from different loci, the overall mean Fst of all SNPs overrepresent the SNPs involved in phenotypic formation. Mo44 overall has the smallest mean Fst against the training set, but a small proportion of SNPs were distinct with the training set. The traits' own genetic characteristics and reverse observation of mean Fst and High-Fst SNP count in rotation across subpopulations could be a possible explanation for genomic prediction accuracy. When PHN11 was the predicted subpopulation, although it had the largest predicting training ratio and mid mean Fst, it also had low extreme Fst SNP counts which led to a relatively good prediction accuracy in almost every trait.

### **GWAS results**

GWAS is used to search for associated SNPs and to find candidate genes. Associated SNPs have more potential to have large contributions in genomic prediction. We employed GWAS for the above traits and found the Fst for each significant SNP. In the GWAS for ear height, plant height, and height above ear, a total of 16, 10, and 14 SNPs were respectively identified (Figure 3.7 & Figure 3.8). In the context of Anthesis D and Silking D, the study identified 4 SNPs (Figure 3.8). In terms of yield, a total of 11 significant SNPs were discovered. The overall mean Fst for these were noted in Table 3.1. The Fst of SNP signals varied for traits. Several SNPs detected for plant heightm, ear height and height above ear exhibited exceptionally high Fst values. The Fst values of all significant SNPs associated with yield were relatively low; the genomic prediction accuracy followed by mean Fst of all when predicting subpopulations. High Fst SNPs trended towards population bias, even when unbiased in training and predict set. Relatively large numbers of high Fst SNP signals were found in ear height and plant height,

providing a possible additional support for the low prediction accuracy when Mo44 was designated as the test set yet high prediction accuracy when MOG was the test set. However, for part of the SNPs detected for Anthesis D and Silking D, the Fst values were relatively high.

NAM-type populations are designed to dissect the genetic basis of complex traits in plants, enabling analysis of a segregating population as well as high-resolution mapping of quantitative trait loci (QTLs) associated with complex traits. The double diploidy PHW65 NAM population has been used to study tassel morphology, where 155 significant SNPs were identified as associated with 15 tassel traits. In our study, though the PHW65 mini-NAM testcrossed with PHT69 showed considerable diversity in many agronomic traits. Although subpopulation differences were not very remarkable for the above measured traits, in scanning significant SNPs, the high Fst SNPs found to be the genetic basis of phenotypic variance were highly subpopulation biased. The close relationship among this small NAM population largely narrows down the genetic diversity which cannot represent the specific associated SNPs well in predicting subpopulations, which directly leads to low prediction accuracy. Diversifying possible genotypes in the training set can boost the accuracy of genomic prediction (Pszczola et al. 2012). In this case using a small breeding population based genomic prediction, the genetic diversity of a training set which cannot represent the genetic basis of the prediction set is possible. When designing a training population in a real breeding project, it's important to take into account not only the relationships among the training varieties but also the connections between the training set and potential predict set(s).



Figure 3.7: GWAS results for ear height, plant height and height above ear. The significant signals are labeled with their SNP Fst value.



Figure 3.8: GWAS results for Anthesis D, Silking D and Yield. The significant signals are labeled with their SNP Fst values.
### Genomic relationship matrices and LD pruning effects on prediction accuracy

To fill the gaps of prediction accuracy between random samples and across subpopulations, we attempted to use a D relationship matrix and GK-based relationship matrix to replace the additive relationship matrix in genomic prediction. The different prediction accuracy is a strong symptom of population bias existing. Differential prediction accuracy also indicated that population structure played a negative role in prediction. Overall, the genomic relationship matrix method had limited impact on prediction accuracy, regardless of whether comparisons were made across families or within individual traits.

Heterosis has been utilized to increase yield in hybrid maize, though heterotic effects vary for different traits. Previous studies have demonstrated that additive genetic variance predominantly influences the variation in flowering time and height, with minimal contribution from  $G \times E$  (gene-environment) interactions (Buckler *et al.* 2009; Romay *et al.* 2013; Peiffer *et al.* 2014). In contrast, maize yield is largely influenced by heterosis and, consequently, dominance effects, with significant contributions from  $G \times E$  interactions (Comstock and Robinson 1948; Hallauer *et al.* 2010). In this study, the plant materials were all hybrids. Most hybrids were test-crossed F1s of a NAM population, hence individual makers are either hybrid or inbred. The binary status and uniformity of the population (all hybrids) may cause small observed differences between additive and dominance methods in genomic prediction. The relatively narrow diversity may have made it hard to capture small effects, contributing to prediction accuracy.



Figure 3.9 The prediction accuracy of using the yellow stripes and two of three NAM subpopulations to predict the remaining NAM subpopulation across different LD threshold pruning and different genomic prediction methods, A (additive relationship matrix), D (dominant relationship matrix), GK (gaussian kernel relationship matrix).

#### Linkage disequilibrium

Because of the influence of population structure on the patterns of linkage disequilibrium (LD), we also explored how to use LD to prune SNPs for genomic prediction. When using LD as a threshold to prune the SNPs, mean Fst of all subpopulations decreased and reached the bottom at an LD threshold of 0.6, then increased back to 0.071 at a threshold of 0.1 (Figure A.31). Compared to the unpruned marker dataset, LD pruning did not have a noticeable pattern of prediction accuracy (Figure 3.9).

As the LD pruning threshold was varied, different levels of genomic relationship-based prediction accuracy displayed irregular fluctuations that intertwined with each other. However, there was one notable exception: when the combination of A, D, and GK was employed, it achieved superior prediction accuracy for MoG in height above ear compared to using A, D, and GK alone. The best LD pruning threshold and genomic relationship matrix varied for traits and training test set partition (Figure 3.9 & Figure A3.3). LD can reveal historical population information in natural evolution and human selection. The strong relationship of population substructure on patterns of LD is not ignorable. Previous studies pointed out BayesB can exploit LD information which leads it to outperform G-BLUP in genomic prediction (Habier et al. 2010). LD is also used to correct for population structure in association mapping studies and when constructing the genomic relationship matrix in genomic prediction (Mathew et al. 2018). LD pruning is most used to adjust the population bias. The benefits of LD pruning also can reduce the SNP number to improve computational efficiency and is considered to yield better prediction accuracy than pruning random makers (Vilhjálmsson et al. 2015). However, the threshold for LD pruning is not well illustrated. In this study, we used different SNP sets after LD pruning to implement genomic prediction. The results show the effect of LD pruning has

various patterns for different traits. This may be due to the different genetic architecture of these traits. But overall, LD pruning didn't provide a benefit in cross-subpopulation prediction accuracy. Other studies also demonstrated that reducing SNP density based on LD has a limited effect on genomic prediction accuracy in a full-sib family in aquaculture species (Song and Hu 2022). Similar results also have been reported in pigs and birds, where LD pruning played little role in prediction (Song et al. 2019; Song and Hu 2022). One potential explanation is that the LD pruning encompasses both SNPs and QTLs under imperfect LD conditions, which can decrease the reliability of predictions. Compared with complicated LD-weighting of the SNPs and LDcorrected genomic relationship matrix generation, LD pruning lost the informative markers (Vilhjálmsson et al. 2015) and complicated the situation of LD between QTL and SNPs in heterogeneous regions; therefore, using a threshold for all SNPs may cause bias in genomic prediction (Gusev et al. 2013; Yang et al. 2015, 2017). In this study, the benefits of LD pruning in population structure may be a trade off with the loss of informative markers. Using LD pruning to reduce the marker size didn't lead to a decline in the prediction accuracy, so overall it could be considered a good method to save computing resources.

# Conclusion

The accuracy of genomic prediction can be significantly influenced by the population structure, which cannot be adequately captured using a summary of marker effects or genomic relationship. The conserved component of a population continues to play a crucial role in shaping phenotypes. LD is a good way to prune SNP numbers and reduce the computational complexity in genomic prediction but has little effect to improve the prediction accuracy or capture the population structure in this study. The different genomic prediction methods didn't change the accuracy much in these hybrids in cross-subpopulation prediction. When the test and

training sets have less genetic bias, increasing the testing and training ratio for a specific number of populations has a minimal impact on prediction accuracy. Based on the above observations and analyses, we propose some ideas that may be helpful in practical genomic prediction in the future. First, increase the observation size and diversity of the training set, as it allows for a better representation of the population structure and genetic variability. Perhaps a universal dataset built from different experiments could be used for building a large model. Alternatively (or additionally), deep learning has potential applications in genomic prediction to capture complex patterns and interactions within genomic data.

#### REFERENCES

- Buckler, E. S., J. B. Holland, P. J. Bradbury, C. B. Acharya, P. J. Brown *et al.*, 2009 The genetic architecture of maize flowering time. Science 325: 714–718.
- Comstock, R. E., and H. F. Robinson, 1948 The Components of Genetic Variance in Populations of Biparental Progenies and Their Use in Estimating the Average Degree of Dominance. Biometrics 4: 254–266.
- Habier, D., R. L. Fernando, and J. Dekkers, 2007 The impact of genetic relationship information on genome-assisted breeding values. Genetics 177: 2389–2397.
- Hallauer, A. R., M. J. Carena, and J. de Miranda Filho, 2010 *Quantitative genetics in maize breeding*. Springer Science & Business Media.
- Jannink, J.-L., A. J. Lorenz, and H. Iwata, 2010 Genomic selection in plant breeding: from theory to practice. Briefings in functional genomics 9: 166–177.
- Lorenz, A. J., 2013 Resource Allocation for Maximizing Prediction Accuracy and Genetic Gain of Genomic Selection in Plant Breeding: A Simulation Experiment. G3 Genes|Genomes|Genetics 3: 481–491.
- Maccaferri, M., M. C. Sanguineti, E. Noli, and R. Tuberosa, 2005 Population structure and longrange linkage disequilibrium in a durum wheat elite collection. Molecular Breeding 15: 271–290.
- Makowsky, R., N. M. Pajewski, Y. C. Klimentidis, A. I. Vazquez, C. W. Duarte *et al.*, 2011 Beyond missing heritability: prediction of complex traits. PLoS genetics 7: e1002051.
- Mathew, B., J. Léon, and M. J. Sillanpää, 2018 A novel linkage-disequilibrium corrected genomic relationship matrix for SNP-heritability estimation and genomic prediction. Heredity 120: 356–368.
- Meuwissen, T. H., B. J. Hayes, and M. Goddard, 2001 Prediction of total genetic value using genome-wide dense marker maps. genetics 157: 1819–1829.
- Nei, M., 1986 Definition and estimation of fixation indices. Evolution 40: 643–645.
- Peiffer, J. A., M. C. Romay, M. A. Gore, S. A. Flint-Garcia, Z. Zhang *et al.*, 2014 The genetic architecture of maize height. Genetics 196: 1337–1356.
- Pérez, P., and G. de Los Campos, 2014 Genome-wide regression and prediction with the BGLR statistical package. Genetics 198: 483–495.
- Pszczola, M., T. Strabel, H. A. Mulder, and M. P. L. Calus, 2012 Reliability of direct genomic values for animals with different relationships within and to the reference population. Journal of Dairy Science 95: 389–400.

- Romay, M. C., M. J. Millard, J. C. Glaubitz, J. A. Peiffer, K. L. Swarts *et al.*, 2013 Comprehensive genotyping of the USA national maize inbred seed bank. Genome biology 14: 1–18.
- Schlenker, W., 2020 Environmental drivers of agricultural productivity growth and socioeconomic spillovers.
- Slavov, G. T., R. Nipper, P. Robson, K. Farrar, G. G. Allison *et al.*, 2014 Genome-wide association studies and prediction of 17 traits related to phenology, biomass and cell wall composition in the energy grass Miscanthus sinensis. New phytologist 201: 1227–1239.
- Song, H., and H. Hu, 2022 Strategies to improve the accuracy and reduce costs of genomic prediction in aquaculture species. Evolutionary Applications 15: 578–590.
- Song, B., A. E. Woerner, and J. Planz, 2021 mixIndependR: a R package for statistical independence testing of loci in database of multi-locus genotypes. BMC bioinformatics 22: 1–21.
- VanRaden, P. M., 2008 Efficient methods to compute genomic predictions. Journal of dairy science 91: 4414–4423.
- Vilhjálmsson, B. J., J. Yang, H. K. Finucane, A. Gusev, S. Lindström *et al.*, 2015 Modeling linkage disequilibrium increases accuracy of polygenic risk scores. The american journal of human genetics 97: 576–592.
- Wu, X., M. S. Lund, D. Sun, Q. Zhang, and G. Su, 2015 Impact of relationships between test and training animals and among training animals on reliability of genomic prediction. Journal of Animal Breeding and Genetics 132: 366–375.
- Ye, S., N. Gao, R. Zheng, Z. Chen, J. Teng *et al.*, 2019 Strategies for obtaining and pruning imputed whole-genome sequence data for genomic prediction. Frontiers in Genetics 10: 673.

# APPENDIX



Figure A3.1 The mean Fst of all SNPs for three NAM subpopulations across different LD pruning thresholds.



Figure A3.2 Manhattan plot of Fst values for all SNPs across three NAM subpopulations.



Figure A3.3. The prediction accuracy of using the yellow stripes and two of three NAM subpopulations to predict the remaining NAM subpopulation across different LD pruning thresholds and different genomic prediction methods, A (additive relationship matrix), D (dominant relationship matrix), GK (gaussian kernel relationship matrix).

# CHAPTER 4: TIME-SERIES PHENOTYPING REFLECTS THE MAIZE GROWTH AND APPLICATION OF YIELD PREDICTION

### Abstract

Plant phenotyping is widely applied in agricultural research and breeding, utilizing multispectral cameras' vegetation indices to assess plant status at scale, potentially replacing labor-intensive fieldwork. This study employs time series normalized difference vegetation index (NDVI) data from unoccupied aerial systems (UAS) to characterize plant growth patterns, demonstrating a high correlation between NDVI values and manually measured traits at various time points. A Genome-wide Association Study (GWAS) using static NDVI values and handmeasured phenotypic traits detected consistent and significant signals, confirming the physiological and genetic relevance of NDVI. A novel method was developed, regressing discrete NDVI measurements against growing degree days (GDD), aligning data across multiple years and locations. Analyzing the NDVI-GDD curve revealed NDVI change as a crucial indicator of growth dynamics. Using silking time and NDVI change, differential GDD (DiffGDD) emerged as strongly correlated with yield. This new metric was incorporated into phenomic and genomic prediction models, yielding promising predictive outcomes.

# Introduction

In the quest for food security and sustainable agricultural practices, enhancing crop yield prediction and understanding the underlying factors that control canopy dynamics are paramount. Maize, as one of the world's most crucial staple crops, plays a pivotal role in global food production and food security. Maximizing maize yield requires a comprehensive understanding of the factors that govern canopy growth and development throughout its growth stages. Traditionally, yield prediction models have relied on a multitude of environmental variables and agronomic practices to estimate potential crop productivity (Moriondo et al., 2007). However,

these models often lack precision because they fail to capture the complex interactions between genetic factors and environmental conditions that govern plant growth (Paine et al., 2012). As a result, there remains considerable room for improvement in yield prediction accuracy, particularly under diverse and dynamic agricultural environments.

Recent advances in high throughput phenotyping technologies have opened new avenues for monitoring and assessing crop growth at high spatial and temporal resolutions (Gill et al., 2022). Among these techniques, the use of Normalized Difference Vegetation Index (NDVI) stands out as the most widely used and stable performance index in the characterization of crop canopies. NDVI is a remote sensing-derived index that quantifies the photosynthetic activity and greenness of vegetation, making it an invaluable indicator of plant health and growth status (Baret & Guyot, 1991; Rouse et al., 1974). Previous research has developed a variety of applications for NDVI, such as disease monitoring (Kumar et al., 2016; Moshou et al., 2005), leaf prediction (Lai et al., 2018; Steltzer & Welker, 2006), drought monitoring (Gu et al., 2008; Peters et al., 2002), and yield prediction (Danilevicz et al., 2021; Moriondo et al., 2007). Time series analysis of canopy NDVI has emerged as a powerful approach to monitor crop development over the growing season. By capturing the temporal patterns of NDVI throughout various growth stages, researchers can gain deeper insights into the phenotypic variation and dynamics exhibited by different genotypes in response to changing environmental conditions (Christopher et al., 2014).

Genomic prediction traditionally uses markers or genomic relationships to predict the breeding value of new individuals (Meuwissen et al., 2001). Recent studies have claimed that the inclusion of measurements of "secondary" traits, which are easier or quicker to measure, can lead to improved predictions (Lado et al., 2018; Pszczola et al., 2013). This is particularly beneficial

when the focal trait exhibits low heritability. By leveraging "secondary" traits with high heritability and substantial genetic and non-genetic correlations, multi-trait prediction can effectively boost the precision of selection for the main trait (Runcie & Cheng, 2019). NDVI acts as a comprehensive trait that reflects plant status and is highly correlated with yield.

However, it is hard to utilize NDVI surveys from different environments or locations together in a uniform way. The time point of NDVI measurements could be a large problem, as the weather can change the plant growth and development processes. The variance in comparison of NDVI data from different environments may be caused by different growth stages and could change based on how plants response to the environment.

The first aim of this study was to explore the potential of time series canopy NDVI data to establish its relationship with critical growth and development traits. Using this approach, we then aimed to solve problems with cross-environment NDVI alignment. Consequently, we interpreted the resulting data for the suitability of its use in physiological growth models for improved yield prediction in maize. We then applied the dynamic time series NDVI data to unravel the genetic basis of canopy traits, which could contribute to improved yield prediction in maize. Finally, the use of NDVI in a secondary-trait-assisted genomic prediction model was tested to capture the environmental factors and GXE factor for agricultural and breeding research.

### Method

This study was based on the field experiment of the Michigan location of the Genomes to Fields (G2F) initiative (www.genomes2fields.org) in 2020 and 2021. In these two years, the G2F hybrid materials were developed from a W10004 Recombinant Inbred Line (RIL) family tested crossed with PHP02. The field site for this study was at the Michigan State University agronomy

farm in East Lansing, Michigan. In each year, the previous crop was soybean. The trial was planted in a randomized complete block design in 2 blocks with 375 hybrids in each block. Weather data were collected by a Watchdog weather station installed in the field on planting day. Dates recorded during the study were converted to growing degree days (GDD), a meaningful measure for explaining physiological changes beyond simple day counts after planting. Daily Growing Degree Days (GDD) were computed using the equation: GDD = (Tmax + Tmin) /2 – Tbase. Here, Tmax represents the daily maximum temperature, Tmin stands for the daily minimum temperature, and Tbase is a constant set to 50 °F, serving as the reference base temperature for maize. In instances where temperatures surpassed 86 °F, Tmax was constrained to 86 °F to prevent overestimation. Conversely, to establish a lower limit, Tmin temperatures were fixed at a minimum of 50 °F.

All field data were collected following the Genomes to Fields standard operating procedure (https://www.genomes2fields.org/resources/). One month after planting, stand counts were recorded for each plot. We recorded the anthesis date and silking dates for appearance of anthers on half of the main tassel spike on half the plants, and silks emerging on half the plants, respectively. Before flowering, we counted green snap occurrences. At harvest, we assessed stand count (SDC), stalk lodging (STL), and root lodging (RTL). Two representative plants per plot were measured for plant height (PHT) and ear height (EHT). Additionally, we marked specific leaf numbers weekly during the vegetative stage to determine ear leaf number (ELN) and total leaf number (TLN) accurately. The largest leaf for these representative plants were identified (LLN) and measured for leaf length (LLL) and width (LLW). The largest leaf area (LLA) was calculated by LLL\*LLW\*0.75.

## UAV data

During the growing season, unoccupied aerial systems (UAs) were flown as a platform equipped with natural color (RGB) and multispectral cameras for image data collection. The multispectral sensor employed in this study captured five bands: blue (475 nm with 32 nm FWHM), green (560 nm with 27 nm FWHM), red (668 nm with 14 nm FWHM), red edge (717 nm with 12 nm FWHM), and near-infrared (842 nm with 57 nm FWHM) (MicaSense Inc., Seattle, WA, USA; http://www.micasense.com/). The UAS data collections were conducted by Remote Sensing and Geographic Information Science (RS&GIS) at Michigan State and were guided by GPS using a pre-programed path under clear and windless conditions, with priority between 10am and 12pm. To maintain a clean field and ensure high-quality aerial image data, manual weeding was carried out. The flight survey frequency and time point are listed in Table 1.

RGB and single-band reflectance images captured by UAV were processed and integrated using Pix4Dmapper software, thereby generating accurate geo-orthophotos. Subsequently, a shapefile containing plot information was meticulously drawn utilizing the UAStools R package. In ArcGIS, these shapefile polygons underwent further refinement via high-resolution RGB orthophotos. By leveraging NIR and red band orthophotos from each flight, we computed the Normalized Difference Vegetation Index (NDVI). This process involved utilizing the polygons to extract the pixel mean NDVI values for each plot across all flight surveys.

#### NDVI curves

Each plot in the NDVI temporal curve of our maize field was created through nonparametric loess regression, which involved fitting NDVI values to corresponding Growing

Degree Days (GDD) across various time points. In these loess regressions, a span of 0.6 was chosen, and all other hyperparameters were kept at their default values (Figure 4.1A). By utilizing these fitted curves, we were able to interpolate specific NDVI values for given GDD values. The differential GDD between start point and end point was called DiffGDD (Figure 4.1B) . This approach enhanced our ability to understand the relationship between NDVI dynamics and the maize growth process.



Figure 4.1: A: Changes in plot NDVI by GDD across the growing season in 2020 and 2021. B: Two examples of plot-level DiffGDD, starting from silking time and ending when NDVI dropped below 0.8 in 2020 and 2021.

The Single Nucleotide Polymorphism (SNP) genotype data set was filtered in the following standard manner: Minor Allele Frequency (MAF)  $\geq$ 5%, Max missing <=0.5. Missing SNPs were imputed in BEAGLE v5.0 (Browning et al. 2018) with 10 iterations for initial burn-

in, 15 sampling interactions, and an effective population size of 50,000. The dataset after MAF and Max missing cleaning and imputation resulted in approximately 326,000 genotyping by sequencing (GBS) SNP markers.

A Genome-Wide Association Study (GWAS) was performed for 13 hand-measured traits and 9 NDVI surveys in 2020 and 17 NDVI surveys in 2021 using 375 WI 1004 x PHP02 hybrids. GWAS was conducted using mixed linear models (FARMCPU) implemented by rMVP. The mixed linear models used for fitting incorporated principal components (PCs) and a genomic relationship matrix (kinship) to account for population structure and relatedness. The number of PC was customized for each trait. A Bonferroni-corrected P-value of 0.05 was used as the threshold for statistical significance of SNP signals. Subsequently, we focused on annotated maize genes within 100-kb regions surrounding these signals. Plausible candidate genes were selected based on their involvement in plant growth and development, photosynthesis pathways, or regulation of leaf morphology.

#### **Genomic prediction**

Genomic prediction ensembled with fixed effects was implemented in the BGLR package. Predictions were applied for three practical scenarios: tested genotypes in untested environments, untested genotypes in tested environments, and untested genotypes in untested environments. For scenarios involving prediction of untested genotypes, one-fold cross validation was used. For scenarios involving prediction of untested environments, three of the four total environments were used to predict the remaining one environment. The combination of the above two scenarios were used in predicting untested genotypes in untested environments. **Results** 

Major differences in phenotypic traits were found between 2020 and 2021. In 2020,

plants had higher flag leaf height and ear height (Figure 4.2 and Figure A4.1). However, the yield in 2020 was much lower than in 2021. The flowering time (GDD based) in 2021 was later than 2020. The difference in leaf number traits between the two years was quite small. On average, the ear leaf number, largest leaf number, and total leaf number in 2021 exhibited a marginal decrease of only 0.6, 0.2, and 0.6, respectively, when compared to the corresponding values observed in 2021.



Figure 4.2: Violin plots for yield (A) and plant height (PHT) (B) in the years 2020 and 2021. **NDVI** 

Throughout the growing season, NDVI exhibited a distinct pattern: it rapidly increased to reach its peak, then maintained this elevated level before gradually declining. However, the NDVI curve also displayed differing patterns across the two years (Figure 4.1A). In 2020, the NDVI peaked earlier and remained at the peak for a shorter duration compared to 2021. Notably, the disparity in the duration of the high NDVI plateau between 2020 and 2021 corresponded to a significant gap in harvest yields during these two years. The NDVI exhibited noticeable variation across accessions, particularly starting at the fourth time point.

NDVI relationship with other traits

The NDVI correlation patterns in 2020 and 2021 exhibited distinct differences (Figure

4.3). In 2020, PHT demonstrated a strong positive correlation with early NDVI surveys. EHT along with ELN and TLN exhibited their highest correlation at the 7/21 time point. In contrast, in 2021, both PHT and EHT showed correlations with several NDVI surveys prior to flowering. ELN, TLN, LLN in 2020 displayed weak correlations with NDVI surveys throughout the season. For both 2020 and 2021, STL and RTL exhibited weak correlations with all NDVI surveys across the growing season. The flowering times in GDD (AGDD and SGDD) displayed an opposite pattern in the time-series NDVI correlation between the two years. In 2020, AGDD and SGDD showed strong correlation with the 7/21 NDVI survey, which closely aligned with the flowering time. Due to the variance in flowering time, NDVI surveys near flowering time displayed weak correlation with late NDVI surveys. In 2021, the region of positive correlation for LLL, LLW Leaf LLA, and SDC coincided with the dates 6/28, 7/06, and 7/14. In 2020, positive correlations of LLW and LLA were observed with late-season surveys, whereas LLL and SDC displayed positive correlations with earlier time points.

Previous studies have already underscored the significance of NDVI, highlighting its high correlation with other traits. In the present study, the correlation patterns of yield also exhibited variation across the two years. In 2020, the highest correlation between yield and NDVI was evident. Conversely, in 2021, the correlation between yield and NDVI was initially low, gradually increasing until 7/14, and then exhibiting a gradual decline towards the end of the season.



Figure 4.3: Heat plot of correlation between different timepoints of NDVI and other hand measured traits in 2020 (left) and 2021(right).

# NDVI vs fully expanded leaf number

During the early stages of plant growth in both years, leaf counts were conducted. Using the count time points and their corresponding growing degree days (GDD) and leaf numbers, we calculated the leaf initiation rate. The leaf initiation rates displayed variation among plots in both years. While the overall mean leaf initiation rate was relatively consistent between the two years, significant disparities were observed in the intercepts (Figure 4.4 A & B).

By utilizing the leaf count time points and NDVI loess curve, we extrapolated GDD values for each plot. The relationship between leaf numbers and their corresponding interpolated NDVI values exhibited a robust positive correlation across the two years (Figure 4.4C). Additionally, the correlation between NDVI and fully expanded leaf numbers displayed certain trends. The rate of NDVI increase starts to decline when the leaf count exceeds ten. This trend could potentially be attributed to variations in leaf area that are associated with different leaf numbers, as well as the potential saturation of NDVI values. Upon comparing the leaf numbers and NDVI values between the two years, it becomes evident that the effect of leaf number

remained consistent.



Figure 4.4: The regression of fully expanded leaf number and GDD for year 2020 (A) and 2021 (B); the interpolation of NDVI from NDVI/GDD curves in 2020 and 2021 (C).

# **GWAS**

To delve into the dynamic regulation of NDVI, a comprehensive Genome-Wide Association Study (GWAS) was conducted using NDVI values from 9 and 17 time points in the years 2020 and 2021, alongside hand-measured traits recorded in both years. Through this study, a total of 102 SNPs were identified, each associated with distinct NDVI time points across the two growing seasons. Among the NDVI datasets from the two years, the greatest number of significant associations was detected in the NDVI values recorded on 7/09/2020, followed by the NDVI data captured on 6/28/2021.

However, in the early and late seasons of both years, the detection of GWAS signals was rare. Meanwhile, a total of 109 SNPs associated with hand-measured traits across these two years were identified (Table A4.1). There were several cases where SNPs identified by multiple NDVI surveys and other hand measured traits had no corresponding candidate genes within a 100kb window. Moreover, nine promising candidate genes were highlighted, contributing to plant growth development, the photosynthesis pathway, or the regulation of leaf morphology.

Out of all the potential genes, Phytosulfokine receptor 1 (PSKR1) emerged as the most strongly associated one in both time series and hand-measured traits (Figure 4.5). Notably, it exhibited associations with 8 flight NDVI surveys in 2021 and SDC for 2020 and 2021. These findings align with a previous study conducted in Arabidopsis, indicating that this gene can indeed regulate plant growth, which corresponds with the observed patterns in SDC and NDVI (Ladwig et al., 2015). Moreover, two genes directly related to photosynthesis, namely Photosystem I Assembly 3 (PSA3) and Photosystem I reaction center subunit VI (PSAH), were identified in correlation with ELN in 2021 and NDVI on 7/09/2020. PSA3's role in aiding photosynthesis assembly has been documented (Shen et al., 2017). Similarly, PSAH1 plays a crucial role in photosynthesis and additionally impacts plant growth.



Figure 4.5: Genome-wide Association Study results for eight flights' NDVI in 2021 and stand count (SDC) in 2020 and 2021. The shared significant signals were found for candidate genes of interest.

In 2021, the Growth Regulation Factor (GRF) gene was identified through its association with Total Leaf Number (TLN). GRF is known to govern the growth of stems and leaves. In

maize, GRF1 mutants exhibit reduced plant height and narrower leaves, yet leaf number remains unaffected (Zhang et al., 2018).

Furthermore, the Floral Homeotic Protein (APETALA 2) was identified through its correlation with the NDVI survey conducted on 8/05/2020. This gene has been reported to regulate flowering time and floral organ development (Aukerman & Sakai, 2003). Interestingly, two members of the AP2 gene family were screened as candidate genes regulating ear height and the ratio of ear to plant height (Fu et al., 2023; Vanous et al., 2018)

Additionally, several chloroplastic-related genes were identified, including CCDA1, DHQS1, RHL9, Putative cyclic nucleotide-gated ion channel 20 chloroplastic, and Putative WEB family protein chloroplastic, as well as WEL2. These associations were linked to various traits: EHT in 2020, ELN in 2021, SDC in 2021, NDVI on 7/15/2020, stalk lodging in 2021, TLN in 2021, and PHT and EHT in 2020.

The Maize barren stalk1 (BA1) gene was identified through its correlation with the NDVI on 7/09/2020 and root lodging in the year 2020 (Figure 4.6). BA1 represents a gene involved in the patterning of maize inflorescences, and, in conjunction with the teosinte branched1 gene, it regulates the development of vegetative lateral meristems (Schmitz & Theres, 2005; Woods et al., 2011). Additionally, the WAB1 gene was detected in relation to NDVI on 6/28 and NDVI on 7/06, as well as stand count in 2021 (Figure 4.5). WAB1 is classified as a TEOSINTE BRANCHED/CYCLOIDEA/PCF (TCP) family transcription factor. Interestingly, it has also been associated with cold resistance in maize during the seedling stage (Yan et al., 2017).



Figure 4.6: Genome-wide Association Study results for NDVI of 7-9-2020 and root lodging (RTL) in 2020. The shared significant signals were found for candidate genes of interest.

# **Application of NDVI-GDD curve**

A significant correlation was observed between single time points and yield, indicating the strong influence of NDVI on yield determination. However, it is important to note that NDVI is highly susceptible to environmental conditions, making it susceptible to the variations caused by different weather conditions that can either accelerate or delay plant growth. Most previous studies treated vegetation indices as a decretal data in model fitting to predict traits of interests for single year one location experiment. Consequently, using NDVI for cross-environmental comparisons can be challenging.

In this study, loess regression was applied for time series NDVI for each plot. We observed that the high NDVI plateau closely followed the yield gap between the two years. To mitigate the environmental variability, we utilized silking Growing Degree Days (SGDD) as the starting point and the point at which NDVI fell back to 0.8 as the end point (Figure 4.1B). By interpolating the GDD and NDVI values for each plot using the NDVI loess function, we calculated the Diff GDD for this specific time period.

To confirm these results, we added the 2019 Michigan Genome to Fields data, which utilized completely different varieties. Remarkably, the diff GDD for this period exhibited a

highly significant correlation (r=0.84) with yield across the three years. When considered individually, a positive correlation of 0.54, 0.64 and 0.35 was observed for the yields of 2019, 2020 and 2021, respectively (Figure 4.7 A).

Furthermore, in our physiological yield prediction model, we incorporated the variables of stand count, stalk lodging and root lodging. Two of the three years of data were used as the training data set to predict the remaining year's data, and promising yield prediction outcomes were achieved. The correlation with 2019, 2020, and 2021 as the test set was 0.61, 0.75, and 0.50, respectively, with mean square of errors (MSE) of 11.06, 11.57, and 13.72 versus observed yield (Figure 4.7 B).



Figure 4.7: (A): Correlation of DiffGDD and yield for three years in Michigan Location. (B): Rotated using two years of data to predict the remaining year. The correlation and MSE of observed and predicted yield are indicated.

It is noteworthy that the DiffGDD and yield were positively correlated with stalk lodging. However, it is worth mentioning that in 2020, two lodging counts were recorded on 10/19/2020 and 10/26/2020. The first counting recorded very few instances of lodging. Interestingly, when we interpolated the NDVI to fall to 0.8, the corresponding GDD date was earlier than the date of the first lodging count. This suggests that the NDVI variability in the late season did not accurately reflect the current lodging status.

# Genomic prediction model using DiffGDD as fixed index

The DiffGDD demonstrated a promising correlation with yield and was integrated into the genomic prediction model. Two datasets from the 2021 Wisconsin G2F project were incorporated to validate this model. Due to differences in sensors and preprocessing methods between Michigan and Wisconsin, the spatial distribution of NDVI varied significantly. Consequently, the previous threshold of NDVI equating to 0.8 was no longer applicable in Wisconsin's data. Upon analysis, regressing NDVI to 0.8 indicated an approximate 15% decrease from the NDVI observed at silking time in the Michigan location (Figure 4.1B). To standardize, we adapted an NDVI of 0.85% at silking time for both Michigan and Wisconsin datasets, truncating the DiffGDD values. These adjusted DiffGDD values exhibited an overall correlation across four distinct environments. However, in the Wisconsin location, the local correlation between DiffGDD and yield was notably weak due to the aforementioned disparities in data and sensors, leading to a completely different spatial distribution compared to Michigan.

To mitigate the effects of differing DiffGDD in various locations, we calculated the mean DiffGDD as fixed effects, representing environmental effects in the genomic model for yield prediction (Figure 4.8A). Predicting genotypes in untested environments yielded promising results with an average correlation of 0.32 and relatively small Mean Squared Error (MSE) of 17.91 (Figure 4.8 B). High correlations (average 0.78) and minor MSE (average 13.05) were observed when predicting yield for untested genotypes in tested environments (Figure 4.8 C). The most challenging scenario, predicting untested genotypes in untested environments, also yielded a good correlation (average 0.28) and well-controlled MSE (average 20.36) (Figure 4.8D).



Figure 4.8: DiffGDD assisted genomic prediction of yield. (A) Correlation of yield and Diff GDD for each individual plot in four different environments. (B) Predicting tested genotypes in untested environments. (C) Predicting untested genotypes in tested environments. (D) Predicting untested genotypes in untested environments. In B, all genotypes in each set of three environments were used to predict all genotypes in the remaining untested environment. In C, five-fold cross validation (CV) was used to predict untested genotypes in all environments. One of the five-fold cross validation results are shown on the plot, while the r and MSE are the mean of fivefold CV. In D, one of the five folds of genotypes in each set of four untested environments were plotted. The r and MSE are the mean of fivefold CV.

# Discussion

In the initial phases of the growing season, huge differences were observed in weather conditions, leading to variations in ear height and plant height. Notably, in 2021, the flowering time was delayed in comparison to 2020, possibly attributed to the impact of cold and rainy weather on seed germination. This effect was clearly reflected in the intercept of the leaf initiation fitting. As the growing season progressed, maize crops in 2020 encountered weather-related stress after flowering, leading to a rapid decrease in NDVI values and ultimately resulting

in a diminished harvest yield.

Traditionally, Growing Degree Days (GDD) hold a crucial role in determining the timing of various phenological events. GDD is closely tied to the accumulation of heat energy during plant growth. Precisely quantifying the timing of a crop's growth stage and predicting the date when it will reach a predefined developmental phase are of utmost importance.

The year-to-year difference in correlation of NDVI and yield reflected two aspects of yield components. NDVI has been widely employed for assessing vegetation health and drought conditions. In 2020, the drought and resulting decay of NDVI starting at 8/18/2020 survey could serve as an indicator that the plants ceased to accumulate yield. Towards the late season of 2020, plots with higher NDVI values tended to maintain a healthier canopy, facilitating yield accumulation. In 2021, due to weather conditions, germination was not uniformly distributed and experienced an overall delay compared to 2020. This was further verified by a common germination related candidate gene found in GWAS for stand count and NDVI in 2021. Additionally, the flowering time in 2021 exhibited a larger variation than in 2020.

Conversely, in 2021, the strong correlation between NDVI and high yield was evident on 7/14, several days before the peak NDVI and flowering time. Notably, greater variability in flowering time was observed in 2021 compared to 2020. The NDVI reading on 7/14 could potentially indicate the flowering time and the response to adverse environments. The diminished correlation in the late season can be explained by two possibilities: 1) the NDVI variability in the late season might be influenced by uncertain germination at the early stage and the wide range of flowering times, rendering it unable to accurately reflect yield accumulation, and/or 2) after an extended period of high NDVI plateau, the accumulation of pure photosynthesis approaches a saturation point for the plants in the field, suggesting that other

factors could assume a dominant role in this scenario.

NDVI is a direct observation obtained from plot-based plant canopies, providing insight into their status at a specific time point. This status is influenced by genetics, environmental conditions, and their interactions, making it promising for predicting other traits. However, leveraging multiple NDVI survey time points across different environments presents challenges. Due to the impact of environmental factors, the growth stage corresponding to NDVI varies significantly, making direct comparisons difficult. To harness the potential of NDVI across diverse environments, the alignment of NDVI data becomes essential. To address this, we regressed a continuous NDVI-GDD curve. This approach seeks to standardize and align the NDVI data across different environments, allowing for more meaningful and accurate comparisons. Leaf numbers were measured in two years, and the interpolation of NDVI using corresponding GDD seems to be consistent, confirming the reliability of this method.

In 2021, the highest yield correlation was observed in the 7/14 NDVI survey, occurring approximately two weeks before flowering time. The varying flowering time in 2021 meant that the 7/14 NDVI reading reflected distinct growth stages. Flowering is considered as the starting point to yield accumulation. Given that, we aligned the flowering time in GDD as the start point. NDVI = 0.8 or derived 85% of silking time was selected as a cutoff in the plant physiological model to approximate the R5-R6 growth stage when the maize had already finished most of its yield accumulation and approached physiological maturity. But NDVI = 0.8 is an approximate value is its potential to be applied more universally for other locations and environments. However, it could be further fine-tuned. For example, it has been claimed that the stage immediately prior to flowering also plays a crucial role in yield formation (Tollenaar & Daynard, 1978; Zhai et al.,

2022). Also, because of its ease to fit and interpolate, we used the loess function to fit the curve. It could be that other sophisticated growth curves and senescence curves could identify the transition point better and provide more verified biological meanings. Hence, future work still has potential to further polish this method to achieve better prediction.

Consequently, we utilized the NDVI-GDD curve as a developmental accumulation trait, providing a more accurate method to quantify energy accumulation for yield. The robust correlation between yield and DiffGDD across three years with two different genetic pools supported our hypothesis. Notably, stand count and lodging emerged as significant factors influencing final harvestable plants. Therefore, our physiological model integrates DiffGDD, stand count, stalk lodging, and root lodging to make predictions. This model yielded promising results for three years. However, in 2021, the correlation between DiffGDD and yield was not as strong as in 2019 and 2020, and the model's performance was slightly inferior compared to the 2019 and 2020 data. This discrepancy could be attributed to the inconsistent flowering time within a plot, leading to potential inaccuracies in recording flowering events. Additionally, during the high NDVI stage, the DiffGDD might exhibit diminishing marginal effects, resulting in diminishing returns for adding more stag green ability compared to earlier stages.

In genomic prediction model, for a compromised of different source and spatial distribution of Michigan and Wisconsin data, mean of diffGDD between silking time and cut-off point in each environment as fixed effects in genomic prediction model. In this situation, genetics and environment were considered in yield prediction, but the GXE effects are not. Overall, the prediction results were encouraging with relatively high correlation and low MSE. Among datasets in this study, environmental effects were a more important factor in yield. Prediction of untested genotypes in tested environments that takes full advantage of

environmental information achieved the highest correlation. Using DiffGDD to quantify the environment effect in genomic prediction can help control the MSE very well. Previous studies have shown that the GXE effects are notable in yield prediction. Our observation of low correlation and high MSE in prediction of yield in untested MI\_2020 environment confirmed that GXE explained even more yield variation under an adverse environment. Using DiffGDD for each plot from the same sensor in a prediction model to handle the GXE effects could achieve a better prediction result.

The success of phenomic and genomic prediction showed that the stay green ability is a main factor affecting yield. Previous pure weather data-based models have not performed well in quantifying the environmental index for yield prediction. Our NDVI-GDD curve and cut-off point can determine the environmental effect on plants leading to the decline of NDVI. The NDVI-GDD also can help lead to the development of a canopy-based crop growth model, or replace the complex leaf-based parameters in current comprehensive crop growth models.

#### **GWAS for NDVI**

This GWAS was conducted in a relatively small hybrid maize population, resulting in slower LD decay. This slower LD decay could potentially impact the accuracy of candidate gene detection. Compared with previous studies (Adak et al., 2023; Wang et al., 2021), we identified more genetic loci associated with NDVI. NDVI reflects the comprehensive status of the plant canopy. Several genes were simultaneously identified for NDVI and other hand measured traits. Several genes for NDVI were identified in both years which can be considered less affected by the environment. Prior to flowering time, stand count was highly correlated with NDVI, and overlapping SNP signals were detected in both traits. Meanwhile, the NDVI, like other canopy traits, dynamically changed across the growing season, and many SNPs were identified by

multiple close flight dates. In the two years of the study (2020, 2021), the weather conditions over the growing season were different. In 2021, cold temperatures and rain occurred after planting which resulted in inconsistent germination and slow growth. WAB1 was detected by NDVI 6/28 and NDVI 7/06 in 2021. This gene is a TCP family transcription factor also associated with seedling cold resistance in maize and rice, consistent with the pattern in 2021. The GRF gene was associated with total leaf number, but previous research reported that GRF changed leaf size (leaf length and width) but did not affect the leaf number in a mutant study. This conflicting result may be caused by distinct alleles or mutations of the gene, only some of which change the leaf number.

## Conclusion

In this study we explored time series canopy NDVI, which reflected the complex nested relationship of canopy leaf traits and other agronomic traits. We also explored the genetic basis of canopy NDVI. A series of growth, development, and morphological candidate genes reveal the genetic variance of different genotypes and further verified the potential application of NDVI to represent multiple physiological responses to the environment. Harnessing the power of time series canopy NDVI data and GDD holds the promise of transforming a maize growth curve which could be useful to maize research and crop management. DiffGDD combined with lodging gave a promising result in yield prediction. This leads to improved understanding of canopy dynamics, paving the way for improved yield prediction models and contributing to global food security. Using DiffGDD as a secondary trait can help bridge the gap between genetics and phenotypic expression. This study aspires to drive innovation in maize breeding and agronomy, offering a more sustainable and efficient approach to address the challenges of feeding a growing world population.

#### REFERENCES

- Adak, A., Murray, S. C., & Anderson, S. L. (2023). Temporal phenomic predictions from unoccupied aerial systems can outperform genomic predictions. *G3*, *13*(1), jkac294.
- Aukerman, M. J., & Sakai, H. (2003). Regulation of flowering time and floral organ identity by a microRNA and its APETALA2-like target genes. *The Plant Cell*, 15(11), 2730–2741.
- Baret, F., & Guyot, G. (1991). Potentials and limits of vegetation indices for LAI and APAR assessment. *Remote Sensing of Environment*, 35(2), 161–173. https://doi.org/10.1016/0034-4257(91)90009-U
- Christopher, J. T., Veyradier, M., Borrell, A. K., Harvey, G., Fletcher, S., & Chenu, K. (2014). Phenotyping novel stay-green traits to capture genetic variation in senescence dynamics. *Functional Plant Biology*, *41*(11), 1035–1048.
- Danilevicz, M. F., Bayer, P. E., Boussaid, F., Bennamoun, M., & Edwards, D. (2021). Maize yield prediction at an early developmental stage using multispectral images and genotype data for preliminary hybrid selection. *Remote Sensing*, 13(19), 3976.
- Fu, J., Wang, L., Pei, W., Yan, J., He, L., Ma, B., Wang, C., Zhu, C., Chen, G., Shen, Q., & Wang, Q. (2023). ZmEREB92 interacts with ZmMYC2 to activate maize terpenoid phytoalexin biosynthesis upon Fusarium graminearum infection through jasmonic acid/ethylene signaling. *New Phytologist*, 237(4), 1302–1319. https://doi.org/10.1111/nph.18590
- Gill, T., Gill, S. K., Saini, D. K., Chopra, Y., de Koff, J. P., & Sandhu, K. S. (2022). A comprehensive review of high throughput phenotyping and machine learning for plant stress phenotyping. *Phenomics*, *2*(3), 156–183.
- Gu, Y., Hunt, E., Wardlow, B., Basara, J. B., Brown, J. F., & Verdin, J. P. (2008). Evaluation of MODIS NDVI and NDWI for vegetation drought monitoring using Oklahoma Mesonet soil moisture data. *Geophysical Research Letters*, 35(22).
- Kumar, S., Röder, M. S., Singh, R. P., Kumar, S., Chand, R., Joshi, A. K., & Kumar, U. (2016). Mapping of spot blotch disease resistance using NDVI as a substitute to visual observation in wheat (Triticum aestivum L.). *Molecular Breeding*, 36, 1–11.
- Lado, B., Vázquez, D., Quincke, M., Silva, P., Aguilar, I., & Gutiérrez, L. (2018). Resource allocation optimization with multi-trait genomic prediction for bread wheat (Triticum aestivum L.) baking quality. *Theoretical and Applied Genetics*, 131, 2719–2731.
- Ladwig, F., Dahlke, R. I., Stührwohldt, N., Hartmann, J., Harter, K., & Sauter, M. (2015). Phytosulfokine Regulates Growth in Arabidopsis through a Response Module at the Plasma Membrane That Includes CYCLIC NUCLEOTIDE-GATED CHANNEL17, H+-ATPase, and BAK1. *The Plant Cell*, 27(6), 1718–1729. https://doi.org/10.1105/tpc.15.00306

- Lai, Y., Pringle, M., Kopittke, P. M., Menzies, N. W., Orton, T. G., & Dang, Y. P. (2018). An empirical model for prediction of wheat yield, using time-integrated Landsat NDVI. *International Journal of Applied Earth Observation and Geoinformation*, 72, 99–108.
- Meuwissen, T. H., Hayes, B. J., & Goddard, M. (2001). Prediction of total genetic value using genome-wide dense marker maps. *Genetics*, 157(4), 1819–1829.
- Moriondo, M., Maselli, F., & Bindi, M. (2007). A simple model of regional wheat yield based on NDVI data. *European Journal of Agronomy*, 26(3), 266–274.
- Moshou, D., Bravo, C., Oberti, R., West, J., Bodria, L., McCartney, A., & Ramon, H. (2005). Plant disease detection based on data fusion of hyper-spectral and multi-spectral fluorescence imaging using Kohonen maps. *Real-Time Imaging*, *11*(2), 75–83.
- Paine, C. T., Marthews, T. R., Vogt, D. R., Purves, D., Rees, M., Hector, A., & Turnbull, L. A. (2012). How to fit nonlinear plant growth models and calculate growth rates: An update for ecologists. *Methods in Ecology and Evolution*, 3(2), 245–256.
- Peters, A. J., Walter-Shea, E. A., Ji, L., Vina, A., Hayes, M., & Svoboda, M. D. (2002). Drought monitoring with NDVI-based standardized vegetation index. *Photogrammetric Engineering and Remote Sensing*, 68(1), 71–75.
- Pszczola, M., Veerkamp, R., De Haas, Y., Wall, E., Strabel, T., & Calus, M. (2013). Effect of predictor traits on accuracy of genomic breeding values for feed intake based on a limited cow reference population. *Animal*, 7(11), 1759–1768.
- Rouse, J. W., Haas, R. H., Schell, J. A., & Deering, D. W. (1974). Monitoring vegetation systems in the Great Plains with ERTS. *NASA Spec. Publ*, 351(1), 309.
- Runcie, D., & Cheng, H. (2019). Pitfalls and Remedies for Cross Validation with Multi-trait Genomic Prediction Methods. *G3 Genes*|*Genomes*|*Genetics*, 9(11), 3727–3741. https://doi.org/10.1534/g3.119.400598
- Schmitz, G., & Theres, K. (2005). Shoot and inflorescence branching. *Cell Signalling and Gene Regulation*, 8(5), 506–511. https://doi.org/10.1016/j.pbi.2005.07.010
- Shen, J., Williams-Carrier, R., & Barkan, A. (2017). PSA3, a Protein on the Stromal Face of the Thylakoid Membrane, Promotes Photosystem I Accumulation in Cooperation with the Assembly Factor PYG7. *Plant Physiology*, 174(3), 1850–1862. https://doi.org/10.1104/pp.17.00524
- Steltzer, H., & Welker, J. M. (2006). Modeling the effect of photosynthetic vegetation properties on the NDVI–LAI relationship. *Ecology*, 87(11), 2765–2772.
- Tollenaar, M., & Daynard, T. (1978). Kernel growth and development at two positions on the ear of maize (Zea mays). *Canadian Journal of Plant Science*, *58*(1), 189–197.

- Vanous, A., Gardner, C., Blanco, M., Martin-Schwarze, A., Lipka, A. E., Flint-Garcia, S., Bohn, M., Edwards, J., & Lübberstedt, T. (2018). Association mapping of flowering and height traits in germplasm enhancement of maize doubled haploid (GEM-DH) lines. *The Plant Genome*, 11(2), 170083.
- Wang, J., Li, X., Guo, T., Dzievit, M. J., Yu, X., Liu, P., Price, K. P., & Yu, J. (2021). Genetic dissection of seasonal vegetation index dynamics in maize through aerial based highthroughput phenotyping. *The Plant Genome*, 14(3), e20155.
- Woods, D. P., Hope, C. L., & Malcomber, S. T. (2011). Phylogenomic analyses of the BARREN STALK1/LAX PANICLE1 (BA1/LAX1) genes and evidence for their roles during axillary meristem development. *Molecular Biology and Evolution*, 28(7), 2147–2159.
- Yan, J., Wu, Y., Li, W., Qin, X., Wang, Y., & Yue, B. (2017). Genetic mapping with testcrossing associations and F2: 3 populations reveals the importance of heterosis in chilling tolerance at maize seedling stage. *Scientific Reports*, 7(1), 3232.
- Zhai, J., Zhang, Y., Zhang, G., Xu, W., Xie, R., Ming, B., Hou, P., Wang, K., Xue, J., & Li, S. (2022). Nitrogen application and dense planting to obtain high yields from maize. *Agronomy*, 12(6), 1308.
- Zhang, D., Sun, W., Singh, R., Zheng, Y., Cao, Z., Li, M., Lunde, C., Hake, S., & Zhang, Z. (2018). GRF-interacting factor1 regulates shoot architecture and meristem determinacy in maize. *The Plant Cell*, 30(2), 360–374.

APPENDIX



Figure 4.1A: The violin plots for ear leaf height (EHT), ear leaf number (ELN), total leaf number (TLN), Anthesis GDD, Silking GDD, stand count, stalk lodging, and root lodging.

Traits	SNP	P-value	Traits	SNP	P-value
2020ELN	S2_103532815	2.77E-08	2020709NDVI	S1_77737455	6.96E-20
2020ELN	S3_171119596	1.39E-07	2020709NDVI	S1_154230727	5.78E-09
2020ELN	S5_213181752	6.30E-08	2020709NDVI	S1_202911306	1.06E-11
2020ELN	S10_82247162	8.79E-08	2020709NDVI	S1_239373253	1.06E-11
2020PHT	S2_219942158	2.02E-09	2020709NDVI	S1_266324915	3.43E-15
2020PHT	S5_21639760	9.07E-08	2020709NDVI	S1_286390929	5.57E-22
2020PHT	S5_50610597	5.00E-08	2020709NDVI	S3_8560498	3.29E-22
2020PHT	S8_120910486	1.99E-08	2020709NDVI	S3_188233224	1.14E-33
2020PHT	S9_154214284	1.00E-09	2020709NDVI	S5_23526521	1.41E-27
2020PHT	S10_15308971	8.35E-09	2020709NDVI	S5_50610597	3.00E-20
2020PHT	S10_142145822	2.14E-14	2020709NDVI	S5_198656059	2.18E-17
2020LLL	S1_76061031	7.50E-08	2020709NDVI	S6_175525345	1.68E-16
2020EHT	S1_23429765	1.17E-25	2020709NDVI	S7_167201885	2.72E-21
2020EHT	S1_69017979	3.64E-19	2020709NDVI	S9_141927374	3.28E-23
2020EHT	S1_302246542	4.08E-26	2020715NDVI	S1_202634381	4.13E-11
2020EHT	S2_149328356	3.31E-17	2020715NDVI	S1_239373253	6.81E-20
2020EHT	S3_10043303	5.66E-18	2020715NDVI	S3_8560498	1.37E-10
2020EHT	S3_10045725	1.91E-16	2020715NDVI	S3_206965633	7.69E-09
2020EHT	S3_143329661	1.18E-15	2020715NDVI	S4_234149651	8.57E-15
2020EHT	S3_170762933	4.17E-29	2020715NDVI	S5_23526521	2.54E-22
2020EHT	S3_216055853	7.16E-43	2020715NDVI	S5_182469085	2.62E-13
2020EHT	S5_32061239	3.63E-28	2020715NDVI	S5_197863733	2.60E-12
2020EHT	S5_208607683	1.65E-35	2020715NDVI	S6_169702527	9.42E-11
2020EHT	S6_127927120	1.98E-15	2020715NDVI	S7_127369376	1.84E-10
2020EHT	S7_131775079	8.91E-33	2020715NDVI	S8_93890164	2.20E-10
2020EHT	S7_136709248	4.65E-37	2020715NDVI	S8_175980979	3.75E-09
2020EHT	S10_138944107	9.66E-40	2020721NDVI	S1_81744169	5.46E-09
2020EHT	S10_140903982	1.33E-26	2020721NDVI	S3_14333328	1.97E-08
2020EHT	S10_142145822	3.59E-09	2020721NDVI	S3_174502643	2.67E-12
2020SDC	S1_77652891	2.28E-22	2020721NDVI	S4_54910368	4.47E-08
2020SDC	S1_239373253	6.74E-22	2020805NDVI	S1_124386467	1.62E-08
2020SDC	S2_2501058	3.07E-26	2020805NDVI	S2_5740979	1.22E-09
2020SDC	S2_230603944	7.55E-23	2020805NDVI	S5_71539516	9.68E-09
2020SDC	S4_192451668	2.78E-08	2020805NDVI	S7_21933725	1.35E-07
2020SDC	S5_23526521	2.97E-23	2020805NDVI	S7_131773926	6.60E-13
2020SDC	S7_173778758	6.95E-12	2020818NDVI	S7_139331533	4.37E-08
2020RTL	S1_99502064	1.39E-20	2021628NDVI	S1_30555408	6.97E-18
2020RTL	S1_195610213	6.07E-149	2021628NDVI	S1_209345164	1.66E-13
2020RTL	S2_189879296	3.80E-59	2021628NDVI	S2_2501058	3.18E-32
2020RTL	S2_226967594	5.20E-190	2021628NDVI	S2_24151614	2.13E-16
2020RTL	S3_13847319	1.30E-89	2021628NDVI	S2_184953337	1.80E-38

Table A4.1: The GWAS results of all hand measurements and time-series NDVI.
7	Table A4.1	(cont'd)				
	2020RTL	S3_188233224	4.13E-286	2021628NDVI	S3_177243465	3.70E-13
	2020RTL	S4_225887951	8.36E-117	2021628NDVI	S5_122957855	6.82E-26
	2020RTL	\$6_12327231	1.18E-162	2021628NDVI	S7_2071029	8.81E-14
	2020RTL	S6_108410579	4.11E-256	2021628NDVI	S7_171625869	1.09E-22
	2020RTL	S7_37518734	1.28E-160	2021628NDVI	S8_110345126	1.92E-10
	2020RTL	S7_127042941	1.84E-55	2021628NDVI	S8_152927808	7.11E-18
	2020RTL	S10_113918976	3.89E-66	2021628NDVI	S9_103983568	3.40E-21
	2021ELN	S1_59902886	3.33E-18	2021628NDVI	S10_140613851	6.70E-13
	2021ELN	S1_60583097	5.92E-23	2021706NDVI	S1_113202879	2.70E-13
	2021ELN	S1_95182112	9.58E-16	2021706NDVI	S2_2501058	3.49E-16
	2021ELN	S1_302972127	6.67E-32	2021706NDVI	S2_43873827	2.04E-15
	2021ELN	S1_306408585	9.95E-13	2021706NDVI	S2_184955958	1.02E-08
	2021ELN	S2_235341086	5.48E-25	2021706NDVI	S5_122957855	1.93E-15
	2021ELN	S2_237243782	8.65E-10	2021706NDVI	S5_210343357	4.10E-12
	2021ELN	S3_70300821	3.97E-22	2021706NDVI	S7_171625869	3.13E-19
	2021ELN	S4_246759605	1.74E-18	2021706NDVI	S9_5027916	1.46E-10
	2021ELN	S5_8158964	5.78E-18	2021706NDVI	S10_115417781	5.61E-10
	2021ELN	S5_76690653	2.66E-30	2021714NDVI	S2_2501058	1.48E-21
	2021ELN	S5_219087318	7.37E-21	2021714NDVI	S2_32536347	2.22E-08
	2021ELN	S7_71158997	8.30E-29	2021714NDVI	S5_212859097	1.03E-07
	2021ELN	S7_115256826	5.79E-13	2021726NDVI	S1_188128301	6.42E-08
	2021ELN	S8_125046889	4.14E-24	2021726NDVI	S1_293140231	2.39E-11
	2021ELN	S9_102709272	5.70E-26	2021726NDVI	S2_2501058	7.88E-22
	2021ELN	S9_143487529	1.31E-24	2021726NDVI	S4_248390503	9.55E-09
	2021TLN	S4_249561710	2.02E-08	2021726NDVI	S6_12327231	1.00E-09
	2021TLN	S5_153556609	8.26E-10	2021726NDVI	S9_161868429	1.20E-07
	2021TLN	\$5_219674523	1.42E-08	2021802NDVI	S1_293140231	8.49E-10
	2021TLN	S6_158989650	5.78E-08	2021802NDVI	S2_2501058	5.11E-19
	2021TLN	S8_125046889	6.03E-10	2021802NDVI	S3_180784697	4.59E-08
	2021TLN	S10_140137777	1.52E-07	2021802NDVI	S6_12327231	1.20E-09
	2021PHT	\$1_221762525	8.43E-13	2021802NDVI	S7_131773926	3.71E-08
	2021PHT	S1_251378724	1.12E-07	2021810NDVI	S1_296260794	1.11E-07
	2021PHT	S2_152998397	3.26E-10	2021810NDVI	S2_2501058	1.11E-08
	2021PHT	S3_19831158	1.04E-08	2021810NDVI	S4_6398794	7.11E-08
	2021PHT	S3_76697237	6.66E-08	2021810NDVI	S6_146035799	2.59E-17
	2021PHT	S3_179051632	1.45E-08	2021810NDVI	S9_150170383	7.79E-10
	2021PHT	S4_31087800	1.85E-13	2021816NDVI	S2_2501058	2.01E-11
	2021PHT	S5_212860411	4.85E-19	2021816NDVI	S2_188999500	1.22E-07
	2021EHT	S1_154444331	1.33E-09	2021816NDVI	S4_6398794	2.20E-09
	2021EHT	S2_28102654	3.93E-09	2021816NDVI	S7_131773926	1.23E-07
	2021EHT	S2_176205195	1.70E-08	2021816NDVI	S8_91249953	3.32E-08
	2021EHT	S3_10507362	9.23E-15	2021816NDVI	S8_136905052	3.50E-08

Table A4.1 (cont'd)								
	2021EHT	S3_135082304	1.00E-08	2021816NDVI	S9_150170383	3.11E-11		
	2021EHT	S3_170815446	5.95E-09	2021826NDVI	S4_39434672	5.98E-09		
	2021EHT	S3_175078566	8.03E-16	2021826NDVI	S9_150170383	1.35E-07		
	2021EHT	S4_38936348	1.13E-09	2021826NDVI	S10_140726257	8.86E-08		
	2021EHT	S4_190582593	4.28E-10	2021830NDVI	S1_272265543	1.31E-11		
	2021EHT	S4_192732980	4.78E-09	2021830NDVI	S2_19115614	4.12E-10		
	2021EHT	S5_153556874	1.03E-08	2021830NDVI	S2_186946859	2.22E-08		
	2021EHT	S5_219604287	3.47E-08	2021830NDVI	S2_232795789	1.14E-08		
	2021EHT	S7_184728968	2.31E-10	2021830NDVI	S4_6816706	2.74E-10		
	2021EHT	S10_77658898	1.40E-12	2021830NDVI	S5_189357381	9.11E-11		
	2021STL	S1_288573479	1.00E-07	2021830NDVI	S6_115443061	1.00E-07		
	2021STL	S5_93919517	9.42E-10	2021830NDVI	S8_154498762	1.47E-07		
	2021STL	S8_5844094	6.87E-14	2021830NDVI	S10_140724592	1.06E-07		
	2021STL	S8_174824245	3.72E-11	2021910NDVI	S3_236260639	6.52E-08		
	2021SDC	S1_266324915	3.79E-19					
	2021SDC	S2_2501058	2.56E-15					
	2021SDC	S2_10741793	2.80E-13					
	2021SDC	S2_184955958	8.51E-17					
	2021SDC	S2_202823606	2.82E-18					
	2021SDC	S3_1541624	2.74E-19					
	2021SDC	S3_29594978	1.04E-20					
	2021SDC	S4_29501021	2.08E-10					
	2021SDC	S5_21894535	1.41E-45					
	2021SDC	\$5_153557930	1.69E-09					
	2021SDC	S7_138758150	2.47E-13					
	2021SDC	S9 130847656	2.85E-12					