

# 1 ARTICLE TYPE

2 Original Research

## 3 CORE IDEAS

4 1. Identification, introgression, and frequency increase of large effect loci are important for cultivar  
5 development.

6 2. The *SstI* locus has a significant effect on cutting score in fields exposed to sawfly infestation.

7 3. Historical genetic information can be utilized to predict haplotypes for lines which have  
8 genome-wide genetic data.

9 4. An R package, HaploCatcher, has been developed to facilitate this analysis in other programs.

10

# **HAPLOCATCHER: AN R PACKAGE FOR PREDICTION OF HAPLOTYPES**

Zachary James Winn<sup>1,\*</sup>, Emily Hudson-Arns<sup>1</sup>, Mikayla Hammers<sup>1</sup>, Noah DeWitt<sup>2</sup>, Jeanette  
Lyerly<sup>3</sup>, Guihua Bai<sup>4</sup>, Paul St. Amand<sup>4</sup>, Punya Nachappa<sup>5</sup>, Scott Haley<sup>1</sup>, and Richard Esten  
Mason<sup>1</sup>

1. Department of Soil and Crop Sciences, Colorado State University, Fort Collins, CO, USA 80523

2. School of Plant, Environmental, and Soil Sciences, Louisiana State University, Baton Rouge,  
LA, USA 70803

3. Department of Crop and Soil Sciences, North Carolina State University, Raleigh, NC, USA  
27607

4. USDA Agricultural Research Service, Hard Winter Wheat Genetics Research Unit, Manhattan,  
KS, USA 66506

5. Department of Agricultural Biology, Colorado State University, Fort Collins, CO, USA 80523

## ABBREVIATIONS

AYN	Advanced Yield Nursery
BLUE	Best Linear Unbiased Estimate
IWGSC	International Wheat Genome Sequencing Consortium
KASP	Kompetative Allele Specific Polymerase Chain Reaction
KNN	K-Nearest Neighbors
LD	Linkage Disequilibrium
Mbp	Megabase pair
PCR	Polymerase Chain Reaction
RF	Random Forest
RGON	Regional Germplasm Observation Nursery
SNP	Single Nucleotide Polymorphism
SRPN	Southern Regional Performance Nursery
USDA	United States Department of Agriculture
WSS	Wheat Stem Sawfly Solid Stem Panel

26

# ABSTRACT

27 Wheat (*Triticum aestivum* L.) is crucial to global food security, but is often threatened by  
 28 diseases, pests, and environmental stresses. Wheat stem sawfly (*Cephus cinctus* Norton) poses a  
 29 major threat to food security in the United States, and solid-stem varieties, which carry the stem-  
 30 solidness locus (*SstI*), are the main source of genetic resistance against sawfly. Marker-assisted  
 31 selection uses molecular markers to identify lines possessing beneficial haplotypes, like that of the  
 32 *SstI* locus. In this study, an R package titled "HaploCatcher" was developed to predict specific  
 33 haplotypes of interest in genome-wide genotyped lines. A training population of 1,056 lines  
 34 genotyped for the *SstI* locus, known to confer stem solidness, and genome-wide markers was  
 35 curated to make predictions of the *SstI* haplotypes for 292 lines from the Colorado State University  
 36 wheat breeding program. Predicted *SstI* haplotypes were compared to marker derived haplotypes.  
 37 Our results indicated that the training set was substantially predictive, with kappa scores of 0.83  
 38 for k-nearest neighbors and 0.88 for random forest models. Forward validation on newly developed  
 39 breeding lines demonstrated that a random forest model, trained on the total available training data,  
 40 had comparable accuracy between forward and cross-validation. Estimated group means of lines  
 41 classified by haplotypes from PCR-derived markers and predictive modeling did not significantly  
 42 differ. The HaploCatcher package is freely available and may be utilized by breeding programs,  
 43 using their own training populations, to predict haplotypes for whole genome sequenced early  
 44 generation material.

## INTRODUCTION

Common bread wheat (*Triticum aestivum* L.) consumption represents nearly 20% of human caloric intake; however, current genetic gain of wheat grain yield is insufficient to meet the rise in demand as the global population increases (Poole et al., 2021; Ray et al., 2013; Shiferaw et al., 2013). Two major threats to grain yield stability in wheat are diseases and pests. One such pest, which presents a major risk to yield stability in the United States northern Great Plains and Mountain West regions, is wheat stem sawfly (*Cephus cinctus* Norton). In terms of domestic losses, the Colorado winter wheat growing region lost approximately 32.7 and 31.2 million dollars' worth of wheat in the years 2020 and 2021, respectfully (Erika et al., 2023). Yield losses for infested hollow-stem varieties can be anywhere from 90-120 kg ha<sup>-1</sup> for spring wheats, potentially resulting in multi-million dollar losses annually (Beres et al., 2007). Furthermore, Beres et al (2011) estimated that sawfly infestation may cost 350 million dollars annually to the United States northern Great Plains and Canadian provinces, making it a major concern of consumers and producers alike.

The wheat stem sawfly is an insect native to North America that infests wheat by ovipositing eggs within the stem of wheat plants from late May to early June (Weiss & Morrill, 1992). Once the egg has been deposited into the stem, the larva emerges and feeds upon the parenchyma and vascular tissue inside the stem (Weiss & Morrill, 1992). After receiving the correct combination of physical and photoperiodic signals (Holmes, 1975), the larvae will migrate downward from its hatching site to an area of the stem near the soil surface and create a notch, which is known as a hibernaculum, that it fills with the excrement of digested plant material (frass) (Weiss & Morrill, 1992). The wheat stem tends to break at the notch development site, causing the

substantial lodging of affected plants that gives the pest its common name. The larvae will enter a period of diapause to then pupate and emerge from the hibernaculum in the following year (Beres et al., 2011).

Although wheat stem sawfly are weak fliers which tend to oviposit near the area of emergence (Beres et al., 2011; Weiss & Morrill, 1992), their distribution is wide and integrated pest management is challenging. Removal of stubs was once recommended as a key cultural control of wheat stem sawfly (Fletcher, 1904), but contemporary research proved that that method was ineffective (Beres et al., 2011). Rotations of wheat followed by fallows also appeared to increase infestation rates, so it has been suggested to use a non-host in rotation as a “trap crop” (Beres et al., 2011; Seamens, 1929). More recently, trap crops have been suggested as a management tool where trap crops are planted as a border around hollow-stem varieties to act as a buffer-zone and prevent infestation of higher yielding hollow-stem varieties (Beres et al., 2009; Peirce, Cockrell, Ode, et al., 2022).

Solid-stem varieties of wheat have been available since the mid-twentieth century (Peirce, Cockrell, Mason, et al., 2022; Weiss & Morrill, 1992). In solid-stem varieties, undifferentiated parenchyma cells create a solid pith within the stem (Berzonsky et al., 2003) and this lessens the severity of yield losses in wheat plants (Beres et al., 2007). The genetic architecture of stem solidness appears oligogenic, with large effect loci being the main contributors to solidness (Peirce, Cockrell, Mason, et al., 2022). One gene found responsible for solidness is *Sst1* (Nilsen et al., 2020, 2017), which was first identified in a QTL study conducted by Cook et al (2004) on the long arm of wheat chromosome 3B (*Qss.msub-3BL*). Stem solidness is thus caused, in part, by tandem

repeats of the *TdDof* gene coding sequence which lead to the filling in of the pith within the wheat stem (Nilsen et al., 2020).

While visual rating of stem solidness can be a reliable method for selecting lines that express solid-stem phenotypes, many wheat breeders in the United States utilize molecular markers to haplotype the region containing *SstI* in a process termed marker-assisted selection. More recently, the United States Department of Agriculture (USDA) Central Small Grains Genotyping Lab located in Manhattan, Kansas has been producing haplotype information on many large effect loci, including *SstI*, for the lines entered into the Southern Regional Performance Nursery (SRPN) and Regional Germplasm Observation Nursery (RGON) [<https://www.ars.usda.gov/plains-area/lincoln-ne/wheat-sorghum-and-forage-research/docs/hard-winter-wheat-regional-nursery-program/research/>]. This service performed by the USDA lab is conducted to assist breeders in releasing lines with the solid-stem trait, and, as a result, it has also created a backlog of information on lines in the SRPN and RGON lines with known *SstI* haplotypes.

Moreover, the lines in the SRPN and RGON have been characterized for genome-wide single nucleotide polymorphisms (SNPs) by the Colorado State University (CSU) Wheat Breeding Program on an annual basis for more than a decade. Winn et al (2022) described a method where historical molecular and haplotype data are utilized to produce accurate haplotype predictions on lines which only have genome-wide molecular data by characterizing either homozygous resistant or homozygous susceptible varieties

In the current work we sought to (1) produce a deployable R statistical software compatible package to perform an analysis similar to the one performed in Winn et al (2022), (2) demonstrate

110 that the analysis preforms similarly in an unrelated germplasm pool for a different locus than those  
111 explored in Winn et al (2022), (3) predict the *SstI* haplotypes of genome-wide genotyped  
112 individuals, and (4) compare the effect of predicted *SstI* verses genotyped *SstI* on wheat stem  
113 sawfly related phenotypes in Colorado State University hard winter wheat germplasm.

114



115

## MATERIALS AND METHODS

### 116 Germplasm

117 Two separate sets of germplasm were utilized in this study. The first population used in  
118 this study was a historical panel of lines submitted to the SRPN and RGON. This panel of lines  
119 consisted of 1,056 distinct genotypes, and all lines in the panel were genotyped genome-wide for  
120 SNPs and haplotyped via a diverse panel of markers for the *SstI/Qss.msub-3BL* locus. The second  
121 population utilized in this study represented contemporary lines in the CSU Wheat Breeding  
122 Program from the 2022 advanced yield nursery (AYN) and the 2022 wheat stem sawfly solid stem  
123 panel (WSS), which were phenotyped for sawfly reaction traits, genotyped for SNPs across the  
124 genome, and screened with kompetitive allele specific polymerase chain reaction (KASP) assays  
125 for the *SstI* locus. The AYN consisted of 107 distinct genotypes and the WSS consisted of 185  
126 distinct genotypes (292 total genotypes). Individuals in the WSS had not gone through any  
127 phenotypic or marker-assisted selection for solid stem or wheat stem sawfly resistance, while  
128 individuals in the AYN had already undergone one generation of field selection for resistance.

### 129 Phenotyping

130 In the 2021-2022 wheat growing season, the AYN and WSS were planted in Akron,  
131 Colorado and a second location of the AYN was planted in New Raymer, Colorado. These sites  
132 were selected for sawfly trials due to the historical presence of sawfly within these regions and the  
133 consistent infestation that they receive (Cockrell et al., 2021; Irell & Peairs, 2014). Furthermore,  
134 in areas surrounding the field sites, 100 sweeps were taken along the field edge bordering an  
135 adjacent wheat stubble field. Sampling began during mid-jointing and continued weekly until no

adult sawfly were found in sweep samples. This data further confirms if infestation pressure was adequate for data collection (Nachappa, 2023; Nachappa & Peirce, 2022).

In each site, all plots were sown in mid-September using a 1.5m wide no-till drill seeder that was guided by a cable and had 4.9m spacing between centers. Following spring green up, centers were pruned using glyphosate (Bayer, St Louis, Missouri, USA) applied by a 1.2m wide hooded sprayer. After end trimming, this resulted in a measurable area of 1.5m by 3.7m. The AYN and WSS were planted in partially replicated designs arranged in rows and columns, with repeated checks included at both row and column levels.

After physiological maturity, when lodging due to sawfly cutting was apparent, a visual cutting score was assigned to each plot in each location. The visual cutting score was assigned as an index of percent plot affected by cutting, which is the physical process by which insect injury detaches most of the wheat stem from the base of the plant. Visual scores were assigned via an ordinal scale ranging from 1-9, where one is fully resistant and erect despite wheat stem sawfly pressure, and nine indicates the whole plot is affected, cut, and prostrate.

## **Genome-Wide Genotyping**

Ten seeds were planted for each line and a 2-3 cm of leaf tissue sample was taken from each plant and bulked for DNA extraction. Genomic DNA was extracted from the samples using MagMax (ThermoFisher Scientific; Waltham, Massachusetts, USA) plant DNA kits following the manufacturer's instructions and quantified using PicoGreen (ThermoFisher Scientific; Waltham, Massachusetts, USA) kits. Extracted DNA was normalized to a concentration of 20 ng  $\mu\text{L}^{-1}$  and sequencing libraries were prepared following the protocol established by Poland et al (2012). The multiplexed libraries were sequenced on a NovaSeq 6000 (Illumina, San Diego, California, USA)

sequencer at 384-plex density per lane. The resulting reads were aligned to the International Wheat Genome Sequencing Consortium (IWGSC) wheat reference sequence RefSeq v2.0 (Appels et al., 2018) using the burrow-wheeler aligner (Li & Durbin, 2009).

The TASSEL 2.0 standalone pipeline (Glaubitz et al., 2014) was used to process the reads obtained from alignment, and markers were organized into compressed variant calling format files (Danecek et al., 2011). Initial variant calling format files were filtered using the following parameters: monomorphic SNPs, insertions, and deletions were removed, SNPs with 85% or less missing data were retained, SNPs with a read depth of more than one or less than 100 were retained, SNPs with a minimum allele frequency of less than 5% were removed, SNPs with more than 10% heterozygosity were removed, and all unaligned SNPs were removed. After filtration, missing data were imputed using the Beagle algorithm V5.4 (Browning et al., 2018), and a synthetic wheat biparental cross between 'W7984' and 'Opata' was used to derive a recombination distance-based map for imputation (Gutierrez-Gonzalez et al., 2019).

## Historical Haplotype Information Curation

Information on the *SstI* locus was curated from historical marker calling files generated by the USDA Central Small Grains Genotyping Lab [<https://www.ars.usda.gov/plains-area/lincolne/wheat-sorghum-and-forage-research/docs/hard-winter-wheat-regional-nursery-program/research/>]. The information for lines in both the RGON and SRPN was standardized to a biallelic haplotype of homozygous *SstI*, heterozygous, and homozygous wildtype represented as “+/+”, “+/-”, and “-/-”, respectively.

Haplotype calls within the CSU wheat breeding program were made using a single marker identified as diagnostic for the *SstI* locus. Extracted and purified DNA, ranging between 50n and

150 ng  $\mu\text{L}^{-1}$  were plated in 96-well plates in 4  $\mu\text{L}$  volumes. Plates contained both test genotype DNA as well as positive, heterozygous, negative, and non-template controls. Each well 4  $\mu\text{L}$  of 2X KASP (LGC Genomics, Middlesex, UK) reaction mixture and 0.11  $\mu\text{L}$  of KASP primer assay mixture. The assay mixture contained an equal mixture of 100  $\mu\text{M}$  of FAM and HEX fluorescence labeled forward primers, as well as 2.5 concentration of 100  $\mu\text{M}$  reverse primer, suspended in molecular grade sterile water (Table 1). Assays were run on a Bio-Rad (Bio-Rad; Hercules, California, USA) CFX96 RT PCR machine and results were read using a single endpoint measurement of florescence. Haplotype calls were made by visual discrimination of florescence groupings. The frequency of allelic states for *SstI* in the training and testing set are also provided (Table 2).

Table 1. List of primers and sequences for the *Usw275* marker.

Marker	Chromosome	Position (Mbp <sup>a</sup> )	Primer	Sequence
<i>Usw275</i>	3B	843.6	HEX Forward <sup>b</sup>	GAAGGTCGGAGTCAACGGAT TAAAGAAAACAAAACCTGTC AAAAAC
			FAM Forward <sup>c</sup>	GAAGGTGACCAAGTTCATGC TAAAGAAAACAAAACCTGTC AAAAAT
			Common Reverse	GAATTTTCGGAGTTACAGAT TGC

Note: <sup>a</sup> Megabase pair position, <sup>b</sup> HEX labeled primer is diagnostic for the solid allele of the *SstI* locus, <sup>c</sup> FAM labeled primer is diagnostic for the non-solid allele of the *SstI* locus

196 Table 2. Number of observations and frequency of *SstI* haplotypes in the training and test  
197 population.

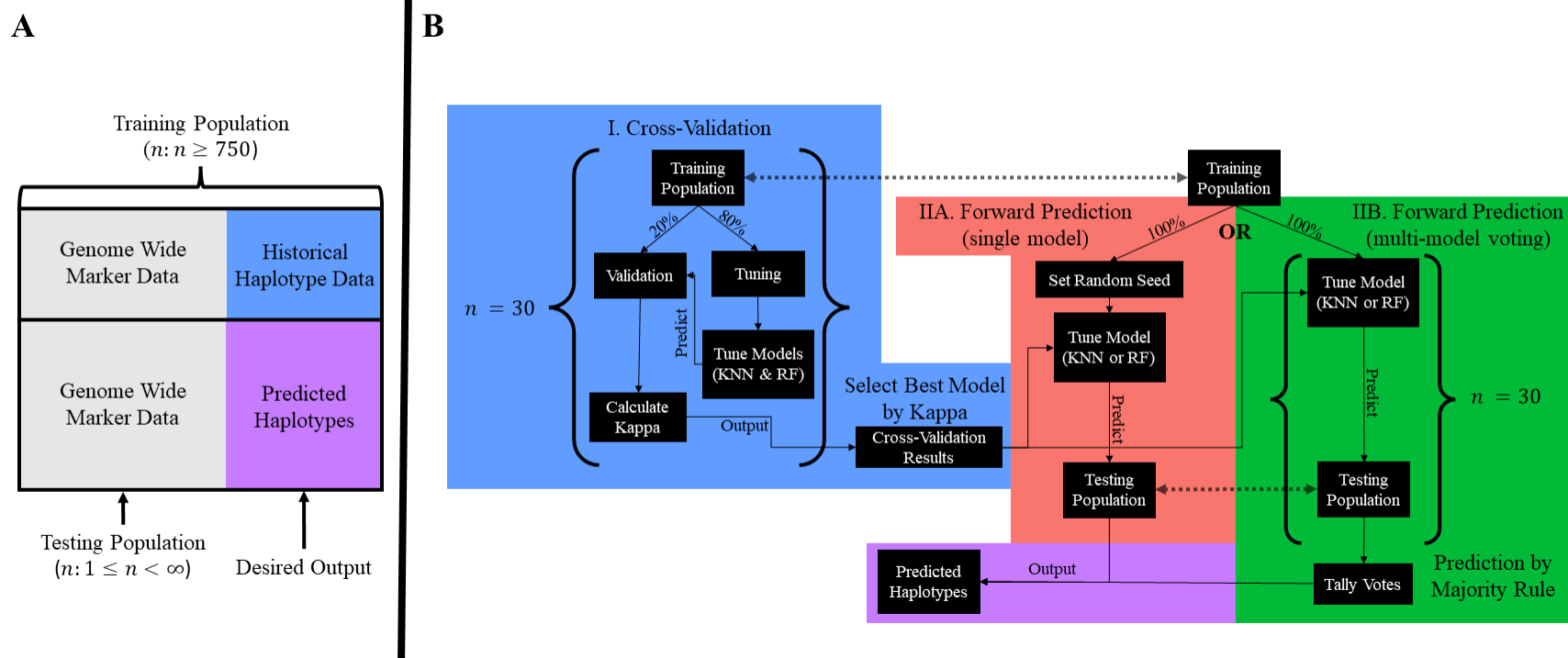
Haplotype	Training		Testing	
	n <sup>a</sup>	Frequency	n	Frequency
+/ <sup>+</sup> <sup>b</sup>	303	0.29	156	0.53
+/- <sup>c</sup>	104	0.10	25	0.09
-/- <sup>d</sup>	649	0.61	111	0.38

198  
199 Note: <sup>a</sup> number of observations, <sup>b</sup> homozygous *SstI* haplotype, <sup>c</sup> heterozygous *SstI* haplotype, <sup>d</sup>  
200 homozygous wildtype haplotype

201

## Package Development and Analysis Pipeline

The “HaploCatcher” package was developed using the “devtools” package (Wickham et al., 2022) in R statistical software (R Core Team, 2022) via the RStudio (Posit; Boston, Massachusetts, USA) development environment on a computer with a Microsoft® Windows operating system. Data inputs required of the package are a marker matrix containing both individuals in the training population and those in the testing population, a historical haplotype classification file for individuals in the training population, and a marker information file which denotes the name, chromosome, and position of each marker in the genotypic matrix (Figure 1A). The package is comprised of several core functions which are then streamlined into the function “auto\_locus”. The “auto\_locus” function conducts a similar analysis pipeline to Winn et al (2022) through the “caret” package (Kuhn, 2008, 2022), while requiring minimal intervention from users (Figure 1B).



214

215 Figure 1. A diagram of [A] input data structure and [B] the “auto\_locus” function pipeline. Panel [A] shows a total data set that is  
 216 partitioned into a training and test population. The training population in panel [A] shows a population of individuals, that is suggested  
 217 to be comprised of more than 750 individuals, which have both genome-wide marker and historical haplotype data. The testing  
 218 population in panel [A] shows a testing population, which may be any size greater than zero, which only has genome-wide marker data.  
 219 Panel [B] shows the workflow of the “auto\_locus” function. In the cross-validation step [I], the total training population is split in a user  
 220 defined way (default is 80:20 split) and the 80% tuning population is used to train and select optimal hyper-parameters for a k-nearest  
 221 neighbors (KNN) and random forest (RF) model. The trained models are then used to predict the haplotype of the validation population.  
 222 The predicted haplotype is then compared to the ‘true’ haplotype and kappa (and accuracy) are calculated. This is repeated a user set  
 223 number of times (default is 30). The best performing model based on accuracy or kappa (default is kappa) is then taken as the model to  
 224 be used in forward prediction. There are two options post cross-validation: [IIA] a single model with a set seed for repeatability or [IIB]  
 225 a user set number of random models (default is 30) used to create a consensus haplotype prediction.



The “auto\_locus” function is comprised of two major phases: cross validation by random partitioning into a user specified ratio (default is 80:20 training-testing; Figure1B - Step I) over a set number of permutations (default is 30) and forward prediction of training population candidates by the best model in cross validation (Figure1B- Step IIA and Figure1B- Step IIB). Prediction of training population haplotypes can be performed by either setting a random seed for reproducibility and performing the model once (Figure 1B - Step IIA) or by running the optimal model with no set seed over a user specified number of permutations (default is 30; Figure 1B - Step IIB) and producing a haplotype prediction by majority rule.

Cross-validation results were visualized using functions from the packages “ggplot2” and “patchwork” (Pedersen, 2022; Wickham et al., 2016). Both the cross-validation and forward prediction by voting steps in the “auto\_locus” function can be run either sequentially or in parallel using the R packages “parallel”, “doParallel”, and “foreach” (Microsoft & Weston, 2022a, 2022b; R Core Team, 2022). Users can specify if the analysis is to be done in parallel (default argument is FALSE) or sequentially. Users may also define the number of processing cores desired for analysis or use the default setting which uses the function “detectCores” from the parallel package to determine the number of system cores and subtract that value by one.

The computer used for development of the package had a hexacore 2.6GHz Intel® (Intel; Santa Clara, California, USA) i7-10750H processor with 12 logical processors, 32 gigabytes of DDR4 RAM and a dedicated NVIDIA (NVIDIA; Santa Clara, California, USA) GeForce® RTX 2070 graphics card. Using the example datasets available in the HaploCatcher package, the “auto\_locus” function performed in parallel with 100 permutations of cross-validations and 100 votes for majority rule resulted in a total runtime of eight minutes and 36 seconds.

## Statistical Analysis

All statistical analysis was conducted in R statistical software version 4.2.2 (R Core Team, 2022). Cutting visual score data was checked for normality by visualization of the distribution of observations using histograms and QQ-plots. Upon evaluation, all data exhibited near-normality or somewhat skewed normal distributions. Mixed linear models were run using the function “mmer” in the package “sommer” (Covarrubias-Pazaran, 2016, 2018). Across locations the following model was utilized to estimate the effect of the *SstI* locus:

$$y_{ijklm} = \mu + H_i + g_j + e_k + r: c_{lm} + \varepsilon_{ijklm}$$

Where  $y_{ijklm}$  is the response,  $\mu$  is the population mean,  $H_i$  is the haplotype fixed effect of the  $i^{\text{th}}$  haplotype,  $g_j$  is the genotypic random effect of the  $j^{\text{th}}$  genotype effect whose variance is defined by the additive relationship matrix among individuals derived by markers (VanRaden, 2008),  $e_k$  is the random environment effect of the  $k^{\text{th}}$  environment that is identically and independently distributed across levels,  $r: c_{lm}$  is the random row by column interaction effect of the  $l^{\text{th}}$  row and the  $j^{\text{th}}$  column whose variance is defined by the two-dimensional penalized tensor-product of spline relationship between row and column effects as described by Lee et al (2013), and  $\varepsilon_{ijklm}$  is the residual error that is identically and independently distributed across all levels.

To compare KASP genotyped haplotype vs machine learning predicted haplotype effects, the same mixed linear model was run twice to estimate an  $H_i$  haplotype fixed effect first using the “true” haplotype calls derived by KASP genotyping and then using the machine learning predicted haplotype information. Fixed effect group mean estimates for both the observed and predicted haplotype effects were estimated via the “predict.mmer” function in the package “sommer”. Visual comparison of effect estimates was summarized using functions from the “ggplot2” package.

Narrow-sense, per-plot, genomic heritability ( $h_g^2$ ) of cutting visual score ratings within environment were estimated using the following mixed linear model:

$$y_i = \mu + g_i + \varepsilon_i$$

Where  $y_i$  is the observation,  $\mu$  is the population mean,  $g_i$  is the random genotype effect of the  $i^{\text{th}}$  genotype whose variance is defined by the marker-derived additive relationship matrix calculated by the “A.mat” function from the “sommer” package, and  $\varepsilon_i$  is the residual error whose variance is identically and independently distributed. Variance components were used in the function “vpredict” in the “sommer” package to estimate  $h_g^2$  using the following formula:

$$h_g^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_\varepsilon^2}$$

Where  $h_g^2$  is the narrow-sense, per-plot, genomic heritability,  $\sigma_g^2$  is the genotypic variance, and  $\sigma_\varepsilon^2$  is the residual error variance.

Importance of defining variables (genome-wide SNP markers) in KNN and RF algorithms was calculated for each iteration of the 100 permutations of cross-validation by using the function “varImp” in the “caret” package. Variable importance, or more aptly put the importance of genome-wide SNP markers in defining haplotypes, was scaled between zero and 100 for comparability across models and the average importance of markers across all permutations was reported in images generated by the “ggplot2” package. Linkage disequilibrium (LD) was calculated for all markers identified as important using the function “LD” in the package “gaston” and results were reported in images derived by functions in the “ggplot2” package (Perdry & Dandine-Roulland, 2018).

Confusion matrices were calculated by comparing the predicted haplotype to the observed haplotype state in the WSS and AYN combined via the function “confusionMatrix” in the “caret” package. Model performance parameters were calculated across iterations of the 100 permutations of cross-validation and the forward prediction of the WSS and AYN. The reported measures of model performance were accuracy, sensitivity, specificity, and unadjusted Cohen’s kappa (McHugh, 2012). Accuracy was calculated as:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Where TP is the number of true positive cases, TN is the number of true negative cases, FP is the number of false positive cases, and FN is the number of false negative cases. Sensitivity and specificity were calculated as such:

$$Sensitivity = \frac{TP}{TP + FN}$$

$$Specificity = \frac{TN}{TN + FP}$$

Cohen’s kappa was calculated as:

$$\kappa = \frac{\Pr(a) - \Pr(e)}{1 - \Pr(e)}$$

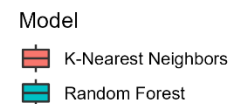
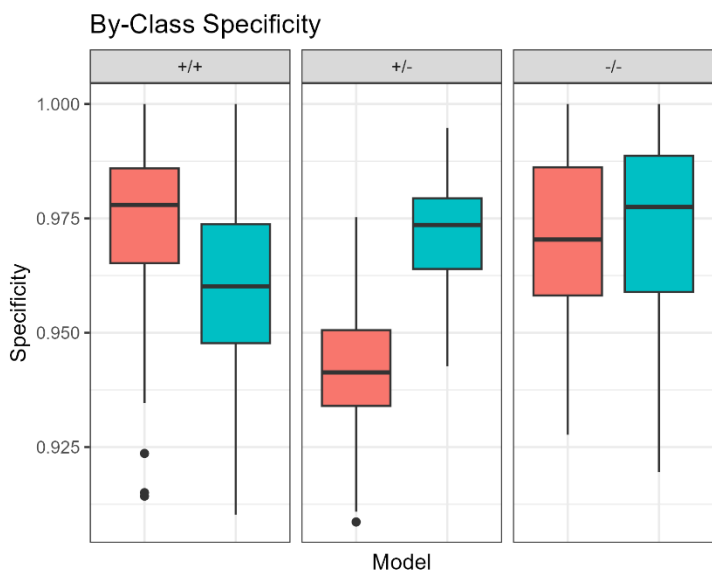
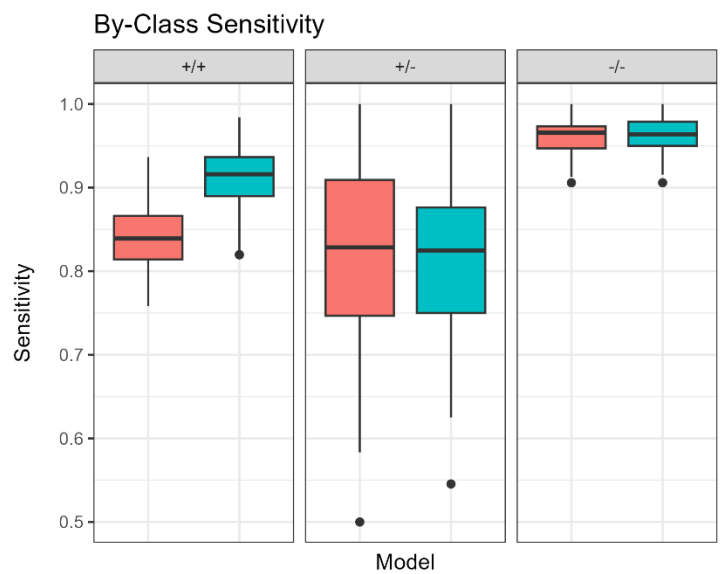
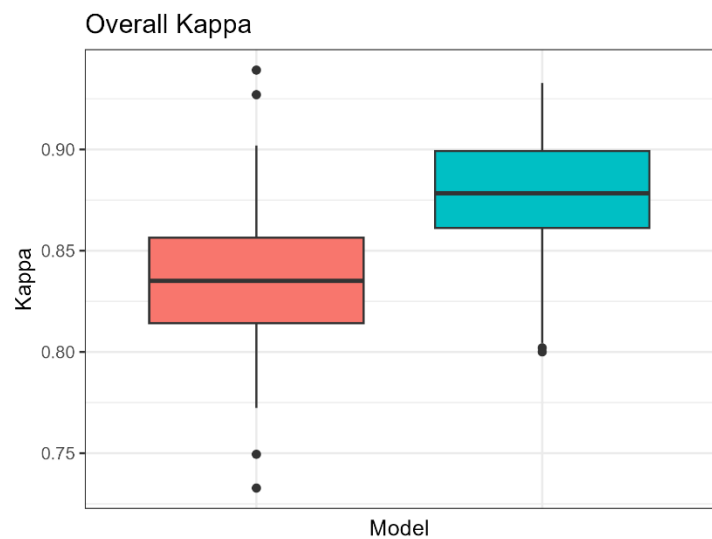
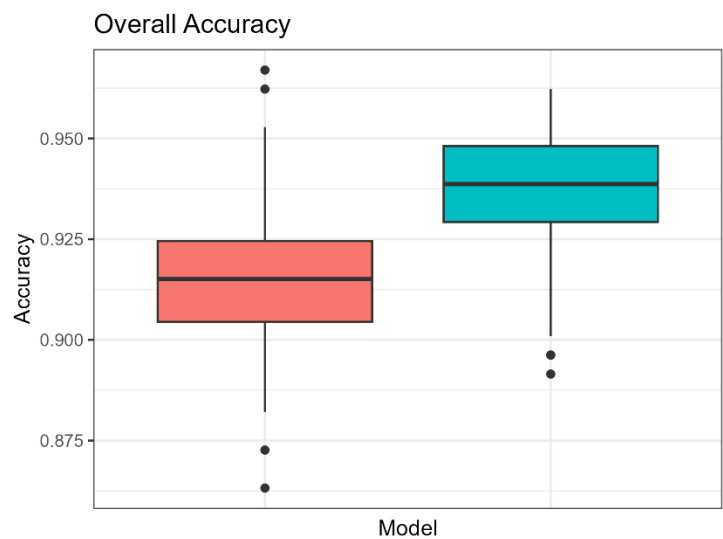
Where  $\Pr(a)$  is the probability of observed agreement, and  $\Pr(e)$  represents the expected rate of chance agreement. Kappa is often considered more robust than accuracy as a measurement parameter of reliability in categorization models because it is not easily biased by sample size (McHugh, 2012).

Kappa may be understood as a measurement which is bound between -1 and 1 where a value of 1 represents a perfectly categorizing model, 0 is the same as chance agreement, and a value of -1 is categorization that is worse than chance agreement (Viera et al., 2005). Historically, a kappa value of 0.8 to 1 is considered to be either “substantial” to “almost perfect” in its predictive ability (Landis & Koch, 1977). All model parameters were either reported in tables or visualized using functions from the “ggplot2” package.

## RESULTS

### *SstI* Prediction Cross-Validation

Cross-validation indicated that the training data was well suited for analysis and substantially predictive based on reported kappa values (Figure 2). Over the 100 permutations of cross-validation, the average kappa value for the KNN model was  $\kappa = 0.83$  and  $\kappa = 0.88$  for RF. The average accuracy for the KNN model was 0.91 and 0.94 for RF. By-class sensitivity varied by haplotype. For homozygous *SstI* calls, KNN had a mean sensitivity of 0.84 and RF had a mean sensitivity of 0.91. For heterozygous *SstI* calls, KNN had a mean sensitivity of 0.82 and RF had a mean sensitivity of 0.81. For homozygous wildtype calls, both KNN and RF had a sensitivity of 0.96.

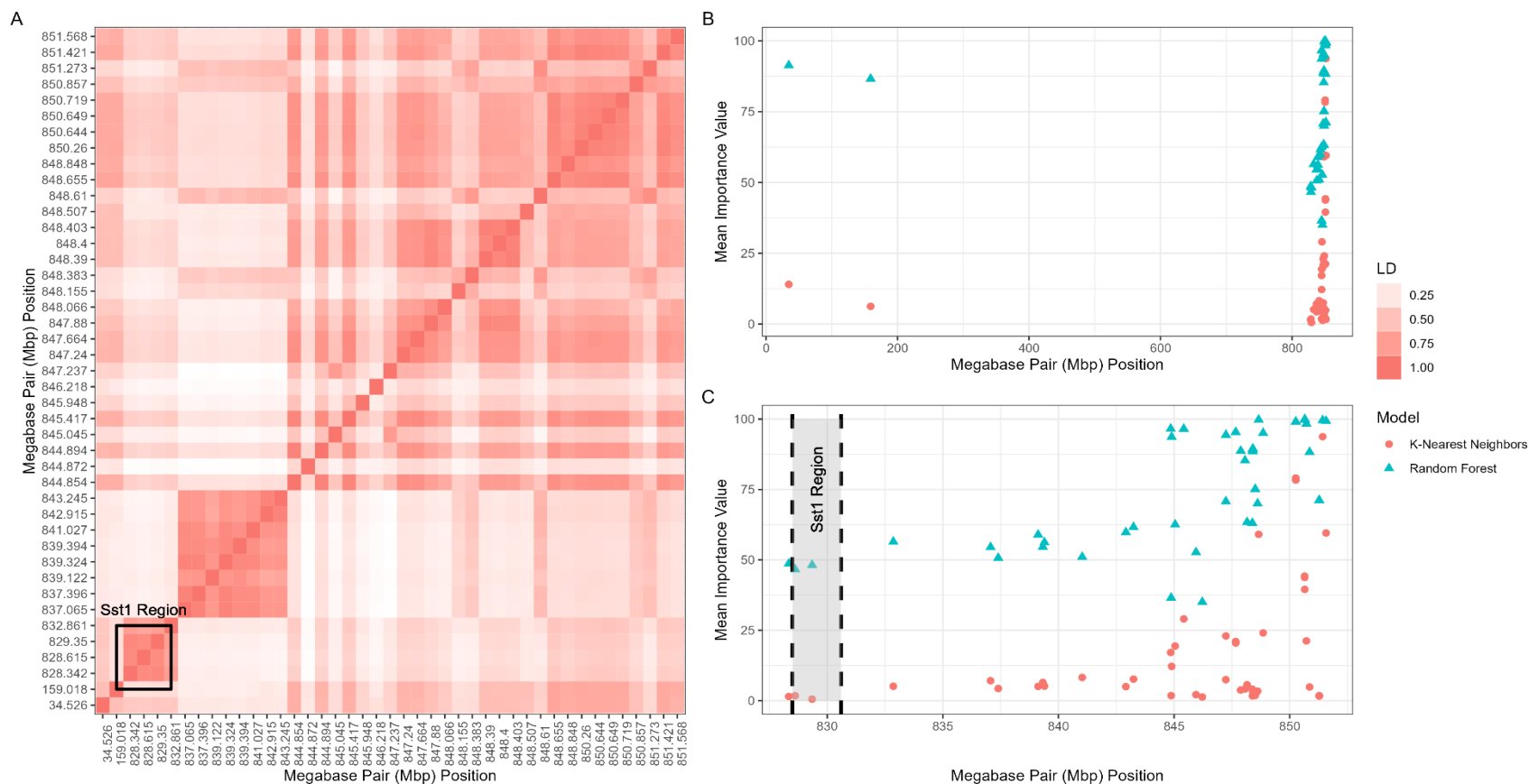


327 Figure 2. A visualization output by the “auto\_locus” function of overall accuracy (**A**), kappa, and by-class sensitivity and specificity  
328 value distributions over 100 permutations of cross validation. The figure legend on the right of the total figure displays the color that  
329 corresponds to which model. The top left panel displays the overall accuracy of each model in boxplots. The top right panel displays the  
330 overall kappa of each model in boxplots. The bottom left panel displays the by-class sensitivity values in boxplots for homozygous *SstI*  
331 individuals (+/+), heterozygous individuals (+/-) and homozygous wildtype individuals (-/-). The bottom right panel displays the  
332 specificity of each model for each classification in boxplots. The x-axis is the model in each figure. The y-axis corresponds to the value  
333 of interest displayed within the graph.



Specificities had a narrow range among haplotype classifications and models. Average specificities for homozygous *SstI* individuals were 0.97 for KNN and 0.96 for RF. Specificities for heterozygous *SstI* individuals were 0.94 for KNN and 0.97 for RF. Specificities for homozygous wildtype individuals were 0.97 for both KNN and RF. These results indicate that KNN tended to under-identify true negatives in heterozygous cases, meaning that it tended to overclassify non-heterozygous individuals as heterozygous. Furthermore, the lower sensitivity scores of both the RF and KNN models (as compared to the higher sensitivity in the homozygous cases) indicates that both models were not as well suited for classifying heterozygous individuals as they were for homozygous individuals. Based on highest achieved average kappa value, the random forest model was selected for use in forward prediction.

Models in cross-validation mainly selected markers in or near the known region of *SstI*, however there were two outliers on the distal short arm of 3B at approximately 34 megabase pairs (Mbp) and 159 Mbp (Figure 3B). When looking at LD among the markers selected for use in the models, it appears that the outlier markers and markers in the 828-852 Mbp region share minor-to-substantial LD ( $r^2$ :  $0.20 < r^2 < 0.70$ ; Figure 3A). More specifically, the LD appears to be very strong between these two outliers and markers in the 848-850 Mbp region ( $\overline{r^2} \approx 0.66$ ) which indicates that the markers may not be inherited independently. This may be the result of misalignment of markers to the wrong arm of the 3B chromosome. Alternatively, this may be a signature of true linkage disequilibrium, indicating that some region on the short arm of 3B is being inherited frequently with the *SstI* locus.



355

356 Figure 3. A visualization of (A) linkage disequilibrium (LD) among the most important markers identified between the k-nearest  
 357 neighbors and random forest model, (B) the importance values of markers across the genome and (C) the importance values of markers  
 358 proximal to the known position of *Sst1*. Panel (A) displays the linkage disequilibrium of each marker identified as important by the  
 359 models. The x and y axes display the marker megabase pair (Mbp) position of each marker. The color within the plot on panel (A) that  
 360 corresponds with the figure legend located on the right indicates the magnitude of LD between those markers. The known location of  
 361 *Sst1* falls within the black box. The graph in panel (B) shows the importance of markers averaged over the 100 iterations. The y-axis

362 displays the importance value of the marker which is represented by the colored dot. The x-axis represents the Mbp position of the  
363 marker. The point color corresponds with which model the point belongs to, which is denoted by the figure legend to the right. The  
364 graph in panel (C) is a zoomed in version of panel (B) where the known location of *SstI* is labeled with a gray shaded box flanked by  
365 dashed lines.

366

When looking at derived importance values within the region, it appears that those outlier markers on the distal short arm of 3B are highly important ( $x > 0.75$ ) for the RF model and less so for the KNN model ( $x < 0.25$ , Figure 3B). Taking a closer look at the known location of *SstI*, it appears that markers within the region are moderately important ( $x \approx 0.50$ ) in RF models while they were non-important for KNN models ( $x < 0.10$ ) (Figure 3B). Interestingly, the most important markers ( $x \geq 0.95$ ) identified by KNN and RF models were in the 845-853 Mbp region.

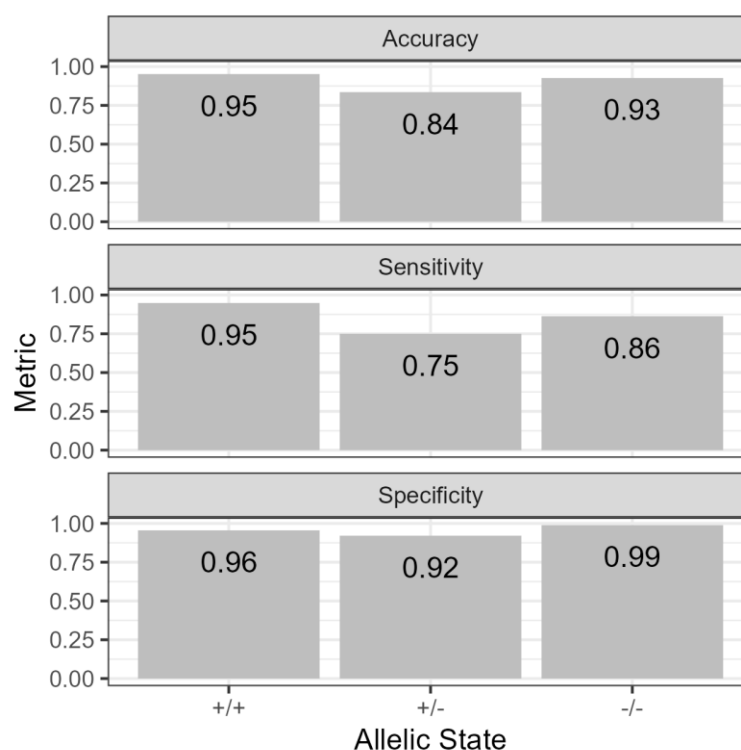
This region of highly important markers is located nearly 15-20 Mbps away from the known location of *SstI*. However, historical markers used to haplotype the *SstI* in the RGON and SRPN are not in perfect linkage with the causal polymorphism. Furthermore, the marker used by the CSU wheat breeding program, which is diagnostic of the *SstI* locus, is found at approximately 843 Mbp, which is directly adjacent to the most important markers for classification. These results may be due to the use of haplotype designations derived from markers which do not lie within or in direct proximity to the *SstI* locus. Regardless, model performance parameters, namely kappa, indicate that both models are capable of “substantial predictions” using historical scales for kappa interpretation (Landis & Koch, 1977).

### ***SstI* Prediction Forward Validation**

Forward validation on the WSS and AYN using a RF model trained on the total available training data produced similar results to that of cross-validation (Figure 4). Accuracies for homozygous wildtype and *SstI* individuals were 0.95 and 0.93, respectively. As observed in the cross-validation results, the accuracy for identification of heterozygous individuals was lower at 0.84. Specificities for homozygous *SstI*, heterozygous *SstI*, and homozygous wildtype were 0.96, 0.92, and 0.99. Sensitivities followed the same trend as cross-validation, where the true positive

389 identification rate for homozygous *SstI* and wildtype individuals was higher (0.95 and 0.85,  
390 respectively) than identification of heterozygous individuals (0.75).

391



392

393 Figure 4. Visualization of performance parameters of predictions by a random forest model trained  
 394 on all available training data. Each subgraph represents a separate measurement of model  
 395 performance. The y-axis displays the magnitude of the metric displayed in each subgraph. The  
 396 allelic state on the x-axis denotes individuals who are homozygous *SstI* (+/+), heterozygous (+/-),  
 397 and homozygous wildtype (-/-). The value of the metric for each allelic state is displayed within  
 398 each bar.

399

Based on the confusion matrix (Table 3) of predicted vs observed haplotypes, the RF algorithm misidentified heterozygous individuals as wildtype frequently. Homozygous wildtype individuals were most often correctly identified (two cases misidentified), followed by homozygous *SstI* individuals (seven cases misidentified). These results may indicate that, while not completely uninformative, these methods may be best suited for identifying homozygous individuals, like in Winn et al (2022), rather than trying to identify heterozygous individuals as well.

Table 3. Confusion matrix of predicted *SstI* haplotypes calls vs haplotypes calls made by kompetative allele specific polymerase chain reaction (KASP).

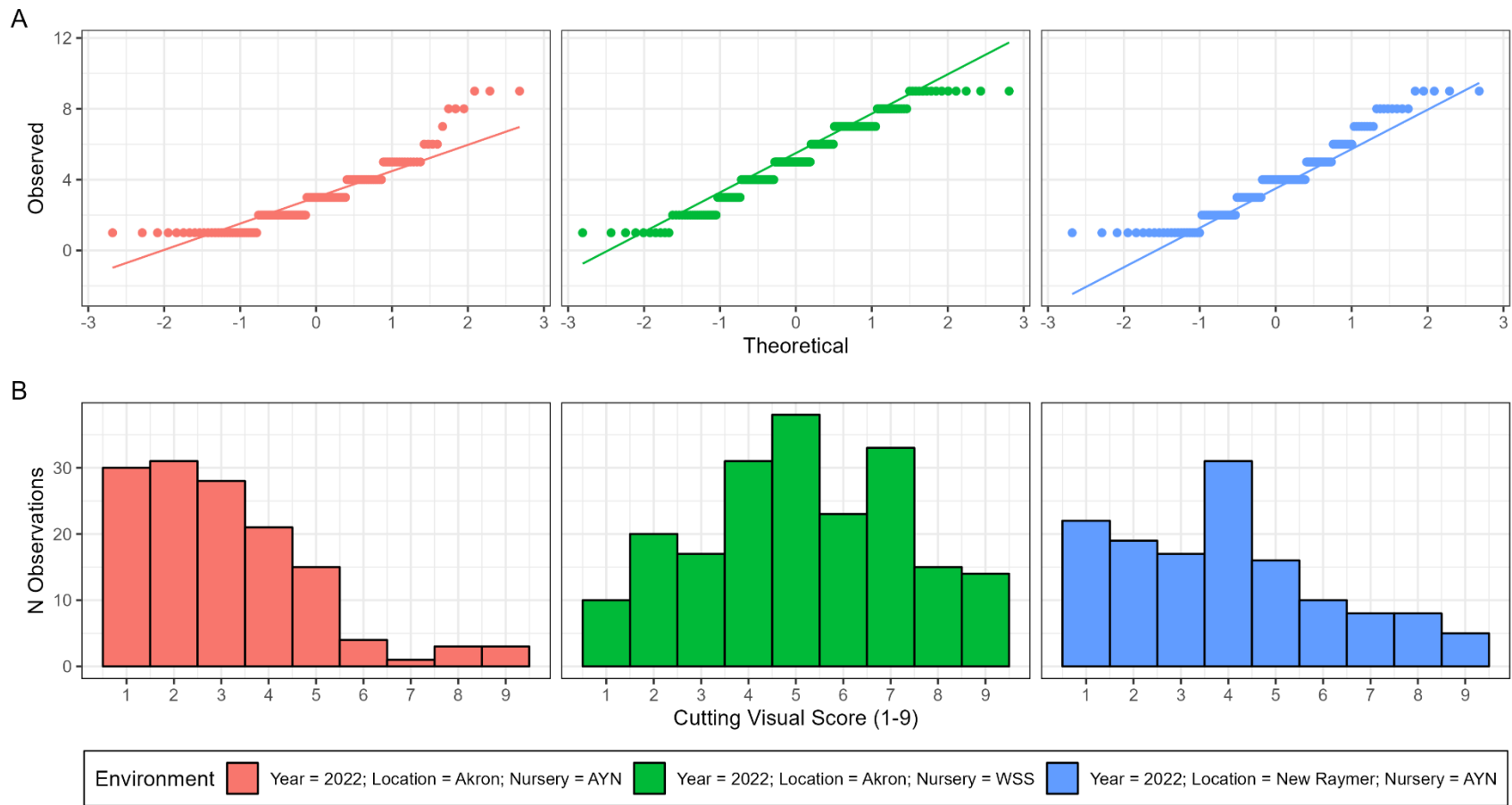
	Predicted	KASP (Observed)		
		+/+	+/-	-/-
	+/+ <sup>a</sup>	147	7	0
	+/- <sup>b</sup>	6	19	15
	-/- <sup>c</sup>	2	0	96

Note: <sup>a</sup> homozygous *SstI* calls, <sup>b</sup> heterozygous *SstI* calls, <sup>c</sup> homozygous wildtype calls



## **Effect of *SstI*: Predicted vs Genotyped Haplotype**

Visual examination of distributions for both locations of the AYN revealed somewhat skewed data distributions while observations from the single location of the WSS followed an approximately normal distribution (Figure 5). Notably, the AYN exhibited a distribution skewed towards lower values of cutting visual scores at both Akron, CO and New Raymer, CO. This is most likely because the AYN is one generation later than the WSS in the breeding process and has already gone through one cycle of selection for wheat stem sawfly resistance. Summary statistics of the location mean, minimum, maximum, standard deviation, heritability and standard error of heritability are also provided (Table 4).



423

424 Figure 5. A visualization of (A) qqplots for each locations data and (B) histogram of the cutting visual score within each environment.  
 425 In panel (A) the y axis represents the observed cutting visual score and the x axis represents the theoretical quantiles. The line going  
 426 across observation points shows the pattern of expected vs observed visual scores for a normal distribution. Panel (B) displays  
 427 histograms of each location where the y axis is the count of observations within the bin and the x axis is the cutting visual score. The  
 428 legend at the bottom of the image displays each environment which corresponds to the color of each subgraph.

429 Table 4. Table of descriptive statistics for each environment.

Year	Location	Nursery	N Observations	N Genotypes	Min	Mean	Max	SD <sup>a</sup>	$h_g^2$ <sup>b</sup>	SE <sup>c</sup>
2022	Akron	AYN	136	107	1.00	3.07	9.00	1.86	0.62	0.10
2022	New Raymer	AYN	136	107	1.00	3.97	9.00	2.23	0.70	0.08
2022	Akron	WSS	201	185	1.00	5.12	9.00	2.18	0.49	0.11

430

431 Note: <sup>a</sup> standard deviation of cutting visual score, <sup>b</sup> narrow-sense, per-plot, genomic heritability, <sup>c</sup> standard error of heritability  
432 measurements.

433

Narrow-sense, per-plot, genomic heritabilities varied among locations. The lowest heritability ( $h_g^2 = 0.49 \pm 0.11$ ) was observed in Akron, CO for the WSS and the highest ( $h_g^2 = 0.70 \pm 0.08$ ) was observed in New Raymer, CO for the AYN. The mean cutting score in both Akron and New Raymer was lower for the AYN ( $\bar{x} = 3.07$  and  $\bar{x} = 3.97$ , respectively) than Akron for the WSS ( $\bar{x} = 5.12$ ), however both nurseries across locations contained visual scores between one and nine. This implies that the generation of selection prior to the AYN did shift the population mean towards resistance, yet it did not cull out all susceptible genotypes, which is expected.

Both predicted and KASP-genotyped *SstI* haplotype calls had significant effects on cutting score ( $P(F) < 0.05$ ). Homozygous *SstI* and heterozygous *SstI* individuals did not have substantially different cutting scores when classifying based on either KASP-genotyped or predicted *SstI* haplotypes. Estimates of *SstI* effects made by prediction were not significantly different from KASP-genotyped *SstI* effects within each haplotype (Figure 6). In the case of KASP-genotyped *SstI* effects, the homozygous wildtype individuals significantly differentiated themselves from both the homozygous *SstI* and heterozygous individuals; however, predicted haplotypes for homozygous wildtype individuals did not significantly differentiate from heterozygous individuals. This is because the prediction method tended to incorrectly classify homozygous wildtype individuals as heterozygous (Table 3).

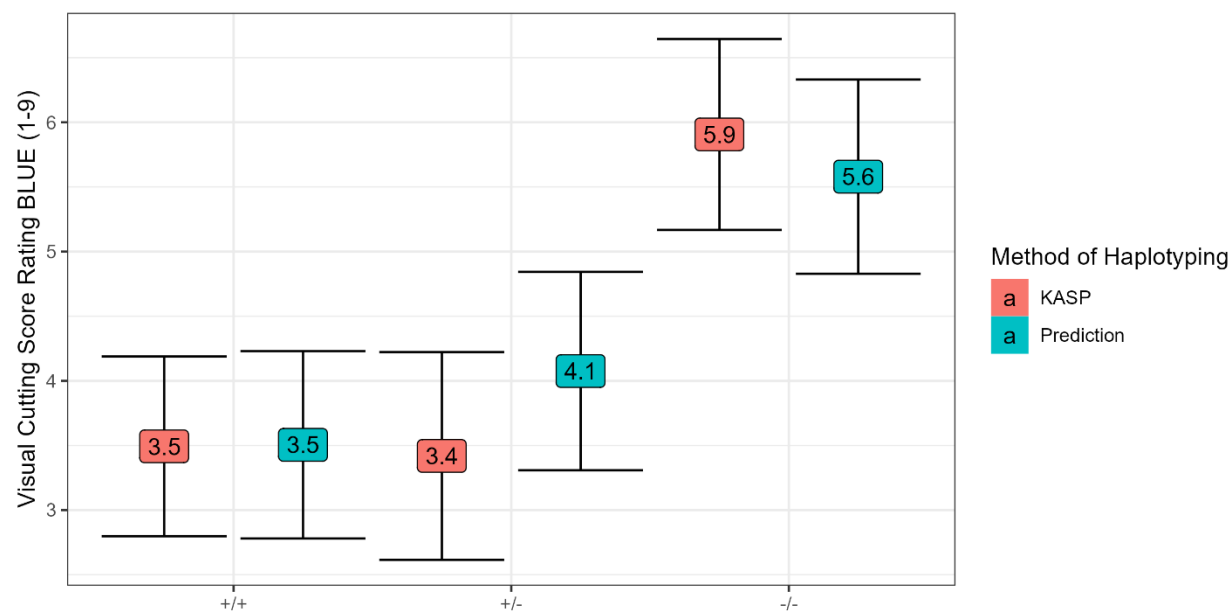


Figure 6. KASP derived *SstI* haplotype call effects vs predicted *SstI* haplotype effects. On the x-axis is the allelic state of the *SstI* locus where +/+ represents individuals homozygous for *SstI*, +/- represents individuals as heterozygous for the *SstI* locus, and -/- represents individuals homozygous for the wildtype allele at the *SstI* locus. The y-axis displays the visual cutting score rating best linear unbiased estimate (BLUE) for the estimated effect of *SstI*. The legend on the right indicates what color-coded box corresponds to which type of *SstI* haplotype assignment. The black bars around each point estimate represents a 95% confidence interval about the estimate.

## DISCUSSION

In the current work, we developed an R user accessible R package by the name “HaploCatcher” which can predict haplotypes using historical information derived from molecular marker assays on genome-wide genotyped lines. The function, “auto\_locus”, allows users to produce predictions for the many lines submitted for sequencing which are not KASP genotyped on an annual basis. Just as in the work performed by Winn et al (2022), we suggest that this method may be deployed in generations where genome-wide sequencing is performed on a very large number of lines which would otherwise not be screened via PCR based assays for these loci. While these predictions were not perfect in their predictive accuracies ( $k = 1$ ), they were substantial in their predictive ability ( $k \geq 0.80$ ) and similar in respect to the results of Winn et al (2022) (Landis & Koch, 1977).

Furthermore, this method is directly accessible to breeding programs, researchers, and students due to its development and deployment through R, a free and accessible statistical computing language. Here, we demonstrated that this method is successful in predicting the *SstI* locus which has a direct impact on improving sawfly resistance in areas threatened by this emerging pest. Moreover, applying this method to resistance loci beyond *SstI* could lead to further progress in the pyramiding and maintenance of wheat stem sawfly resistance.

The cross-validation results in the current work were similar to those in Winn et al (2022); however, unlike Winn et al (2022) we included the option “include\_hets” in the “auto\_locus” function, which allows for the prediction of biallelic loci with a heterozygous state. In our results,

we observed that both the KNN and RF models were not as capable of identifying heterozygous cases as they were homozygous cases. There may be several reasons for this phenomenon.

Firstly, lines in the CSU program are often initially sequenced in the  $F_{3:5}$  generation with recurrent sequencing in each subsequent year. Leaf tissue from ten seeds of each line is bulked and used to prep libraries for sequencing. Therefore, if a heterogenous line for the *SstI* locus was selected in the  $F_3$  generation, the DNA extracted may be a small 1:2:1 mixture of *SstI* haplotypes, and because of this, the sequence of the region may not be truly representative of a heterozygous *SstI* haplotype, leading to misidentification by this method.

Secondly, we curated KASP data produced by the USDA Central Small Grains genotyping lab over years for training. This data, while highly informative, showed some inconsistency across years. Marker platforms and locus region sizes change over years, and this can lead to unexpected association of markers with the locus. More specifically, we observed that markers 15-20 Mbps away from the known region of the locus were identified as “highly important”. This may be because markers which were used to haplotype the region were not in direct linkage with the causal polymorphism, and this led to the detection of markers on the distal long arm of 3B as important.

Furthermore, we aligned our genetic data to the IWGSC wheat reference genome version 2.0 (Appels et al., 2018). This genome is genetically distant from the wheat germplasm located in the Great Plains area of the United States and may have led to misalignments of sequencing reads, like those potentially observed in the marker importance figure (Figure 3A). Moreover, this genetic data was imputed using Beagle (Browning et al., 2018), which is also not perfectly predictive. Therefore, the summed errors of genomic sequencing method, historical data curation, misalignment, and imputation may have contributed to the lower predictability of heterozygous

classes. We therefore suggest that if users can do so, that they produce their own training populations within their own programs and genotype them with a consistent set of markers for best results. Regardless, our proposed models are still substantially informative.

When comparing KASP-based and predicted *SstI* haplotype call group mean estimates, we observed that the predicted and KASP-based haplotype group means were not significantly different from each other within haplotype. We did observe that homozygous wildtype and heterozygous *SstI* haplotype group means did not significantly differentiate in the prediction as the RF model used to make this prediction often grouped heterozygous individuals with wildtype haplotypes (Table 3).

Irrespective of these shortcomings, this method provides a way to assess haplotypes of interest, with a measurable margin of error, in generations that would otherwise not be screened for these. More importantly, this package now provides an easily accessible method of pipeline implementation for breeding programs. While targeted sequencing platforms (Lundberg et al., 2013) may reduce the need for a method like this, it will remain useful for programs without access to targeted sequencing platforms or programs missing probes for specific loci of interest. Furthermore, this method can accommodate many different sequencing platforms (diversity arrays, genotyping-by-sequencing, amplicon, etc.), does not require physical position information, and can potentially be widely applied across species. Lastly, we demonstrated that this method can be applied within breeding programs and produce comparable results to PCR based marker calls; more specifically, we showed that this method could be a viable way of screening early development germplasm for the *SstI* locus, and thus increase the frequency of this locus earlier in the development pipeline.



524

## CONCLUSIONS

525       The utility of marker-assisted selection has been vetted through the vast literature available  
526 for the method. However, with whole genome sequencing technologies being applied in early  
527 generations for use in genomic prediction, there lies an opportunity to acquire data on haplotypes  
528 of important loci. The method proposed in Winn et al (2022) allows breeding programs to organize  
529 their historical marker-assisted selection data to produce predictive haplotype calls for lines in  
530 generations where PCR-based assays for loci of interest are not run due to increased time, labor,  
531 and genotyping cost. This can allow breeders to observe locus profiles of potential varieties much  
532 earlier in the breeding process than before. Here, we chose wheat stem sawfly – an emerging pest  
533 that threatens grower profitability and the dryland cropping agroecosystem– as a test case to  
534 demonstrate the effectiveness of this method. We used existing genotypic data sets to deliver  
535 breeders precise predictions of the presence of a major resistance gene, *Sst1*, allowing for improved  
536 selection for stem sawfly resistance at an earlier generation. With the development of the  
537 HaploCatcher package, there is now a freely accessible software for easier implementation of this  
538 method in other breeding pipelines.

539

540

## ACKNOWLEDGMENTS

541           This research was made possible by funds derived from the competitive grant 2022-68013-  
542 36439 (WheatCAP) from the USDA National Institute of Food and Agriculture. Mention of trade  
543 names or commercial products in this publication is solely for the purpose of providing specific  
544 information and does not imply recommendation or endorsement by the US Department of  
545 Agriculture. The USDA is an equal opportunity provider and employer.

546

## CONFLICT OF INTEREST

547   The authors declare no conflict of interest.

548

## ORCID

549   Zachary James Winn - 0000-0003-1543-1527  
550   Mikayla Hammers - 0000-0003-2661-8389  
551   Jeanette Lyerly - 0000-0003-3853-9581  
552   Noah DeWitt - 0000-0001-9055-993X  
553   Scott Haley - 0000-0002-1996-9570  
554   Guihua Bai - 0000-0002-1194-319X  
555   Paul St. Amand - 0000-0002-5432-5033  
556   Punya Nachappa- 0000-0002-9673-9939

557

## DATA AVAILABILITY

558           Code and data utilized in this study may be found at  
559 <<https://github.com/zjwinn/HAPLOCATCHER-A-PACKAGE-FOR-PREDICTION-OF->

560 HAPLOTYPES>. The package developed for this project, “HaploCatcher”, can be directly  
561 downloaded to an R installation using `devtools::install_github(“zjwinn/HaploCatcher”)` to directly  
562 install from GitHub or `install.packages(“HaploCatcher”)` to install from the CRAN database.

563

## REFERENCES

- 564 Appels, R., Eversole, K., Feuillet, C., Keller, B., Rogers, J., Stein, N., Pozniak, C. J., Choulet, F.,  
565 Distelfeld, A., & Poland, J. (2018). Shifting the limits in wheat research and breeding using  
566 a fully annotated reference genome. *Science*, 361(6403), eaar7191.
- 567 Beres, B. L., Cárcamo, H. A., & Bremer, E. (2009). Evaluation of Alternative Planting Strategies  
568 to Reduce Wheat Stem Sawfly (Hymenoptera: Cephidae) Damage to Spring Wheat in the  
569 Northern Great Plains. *Journal of Economic Entomology*, 102(6), 2137–2145.  
570 <https://doi.org/10.1603/029.102.0617>
- 571 Beres, B. L., Cárcamo, H. A., & Byers, J. R. (2007). Effect of Wheat Stem Sawfly Damage on  
572 Yield and Quality of Selected Canadian Spring Wheat. *Journal of Economic Entomology*,  
573 100(1), 79–87. <https://doi.org/10.1093/jee/100.1.79>
- 574 Beres, B. L., Dosdall, L. M., Weaver, D. K., Cárcamo, H. A., & Spaner, D. M. (2011). Biology  
575 and integrated management of wheat stem sawfly and the need for continuing research.  
576 *The Canadian Entomologist*, 143(2), 105–125. <https://doi.org/10.4039/n10-056>
- 577 Berzonsky, W. A., Ding, H., Haley, S. D., Harris, M. O., Lamb, R. J., McKenzie, R., Ohm, H. W.,  
578 Patterson, F. L., Peairs, F., & Porter, D. R. (2003). Breeding wheat for resistance to insects.  
579 *Plant Breeding Reviews*, 22, 221–296.
- 580 Browning, B. L., Zhou, Y., & Browning, S. R. (2018). A one-penny imputed genome from next-  
581 generation reference panels. *The American Journal of Human Genetics*, 103(3), 338–348.  
582 <https://doi.org/10.1016/j.ajhg.2018.07.015>

583 Cockrell, D. M., Randolph, T., Peirce, E., & Peairs, F. B. (2021). Survey of Wheat Stem Sawfly  
584 (Hymenoptera: Cephidae) Infesting Wheat in Eastern Colorado. *Journal of Economic*  
585 *Entomology*, 114(2), 998–1004. <https://doi.org/10.1093/jee/toab015>

586 Cook, J. P., Wichman, D. M., Martin, J. M., Bruckner, P. L., & Talbert, L. E. (2004). Identification  
587 of Microsatellite Markers Associated with a Stem Solidness Locus in Wheat. *Crop Science*,  
588 44(4), 1397–1402. <https://doi.org/10.2135/cropsci2004.1397>

589 Covarrubias-Pazaran, G. (2016). Genome assisted prediction of quantitative traits using the R  
590 package sommer. *PLoS ONE*, 11, 1–15.

591 Covarrubias-Pazaran, G. (2018). Software update: Moving the R package sommer to multivariate  
592 mixed models for genome-assisted prediction. *Biorxiv*.

593 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Handsaker, R. E.,  
594 Lunter, G., Marth, G. T., Sherry, S. T., McVean, G., Durbin, R., & Genomes Project  
595 Analysis, G. (2011). The variant call format and VCFtools. *Bioinformatics*, 27(15), 2156–  
596 2158. <https://doi.org/10.1093/bioinformatics/btr330>

597 Erika, P., Nachappa, P., Hill, R., Mason, E., Erker, B., & Denninghoven, T. (2023). *Wheat Stem*  
598 *Sawfly Economic Impact Study*. Colorado Association of Wheat Growers.  
599 [https://coloradowheat.org/wp-content/uploads/2022/07/WSS-Economic-Impact-](https://coloradowheat.org/wp-content/uploads/2022/07/WSS-Economic-Impact-Study_06212022.pdf)  
600 [Study\\_06212022.pdf](https://coloradowheat.org/wp-content/uploads/2022/07/WSS-Economic-Impact-Study_06212022.pdf)

601 Fletcher, J. (1904). Experimental farms reports–Report of the Entomologist and Botanist.  
602 *Appendix to the Report of the Minister of Agriculture, Sessional Paper, 16*, 172–173.

603    Glaubitz, J. C., Casstevens, T. M., Lu, F., Harriman, J., Elshire, R. J., Sun, Q., & Buckler, E. S.  
604           (2014). TASSEL-GBS: a high capacity genotyping by sequencing analysis pipeline. *PLoS*  
605           *One*, 9(2), e90346.

606    Gutierrez-Gonzalez, J. J., Mascher, M., Poland, J., & Muehlbauer, G. J. (2019). Dense genotyping-  
607           by-sequencing linkage maps of two Synthetic W7984×Opata reference populations  
608           provide insights into wheat structural diversity. *Scientific Reports*, 9(1), 1793.  
609           <https://doi.org/10.1038/s41598-018-38111-3>

610    Holmes, N. (1975). Effects of moisture, gravity, and light on the behavior of larvae of the wheat  
611           stem sawfly, *Cephus cinctus* (Hymenoptera: Cephidae). *The Canadian Entomologist*,  
612           107(4), 391–401.

613    Irell, B., & Pears, F. (2014). *Wheat Stem Sawfly: A New Pest of Colorado Wheat*.  
614           [https://extension.colostate.edu/topic-areas/insects/wheat-stem-sawfly-a-new-pest-of-](https://extension.colostate.edu/topic-areas/insects/wheat-stem-sawfly-a-new-pest-of-colorado-wheat-5-612/)  
615           colorado-wheat-5-612/

616    Kuhn, M. (2008). Building predictive models in R using the caret package. *Journal of Statistical*  
617           *Software*, 28(1), 1–26.

618    Kuhn, M. (2022). *caret: Classification and Regression Training*. [https://CRAN.R-](https://CRAN.R-project.org/package=caret)  
619           project.org/package=caret

620    Landis, J. R., & Koch, G. G. (1977). The Measurement of Observer Agreement for Categorical  
621           Data. *Biometrics*, 33(1), 159–174. JSTOR. <https://doi.org/10.2307/2529310>

622 Lee, D.-J., Durbán, M., & Eilers, P. (2013). Efficient two-dimensional smoothing with P-spline  
623 ANOVA mixed models and nested bases. *Computational Statistics & Data Analysis*, 61,  
624 22–37. <https://doi.org/10.1016/j.csda.2012.11.013>

625 Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler  
626 transform. *Bioinformatics*, 25(14), 1754–1760.  
627 <https://doi.org/10.1093/bioinformatics/btp324>

628 Lundberg, D. S., Yourstone, S., Mieczkowski, P., Jones, C. D., & Dangl, J. L. (2013). Practical  
629 innovations for high-throughput amplicon sequencing. *Nature Methods*, 10(10), 999–1002.

630 McHugh, M. L. (2012). Interrater reliability: The kappa statistic. *Biochemia Medica*, 22(3), 276–  
631 282.

632 Microsoft, & Weston, S. (2022a). *doParallel: Foreach Parallel Adaptor for the “parallel”*  
633 *Package*. <https://CRAN.R-project.org/package=doParallel>

634 Microsoft, & Weston, S. (2022b). *foreach: Provides Foreach Looping Construct*.  
635 <https://CRAN.R-project.org/package=foreach>

636 Nachappa, P. (2023). *Personal Communication of Unpublished Results*.

637 Nachappa, P., & Peirce, E. (2022). *Wheat Stem Sawfly in Colorado – Frequently Asked Questions*.  
638 Colorado State University.  
639 [https://webdoc.agsci.colostate.edu/csucrops/reports/winterwheat/2022/Sawfly\\_2022.pdf](https://webdoc.agsci.colostate.edu/csucrops/reports/winterwheat/2022/Sawfly_2022.pdf)

640 Nilsen, K. T., N’Diaye, A., MacLachlan, P. R., Clarke, J. M., Ruan, Y., Cuthbert, R. D., Knox, R.  
641 E., Wiebe, K., Cory, A. T., Walkowiak, S., Beres, B. L., Graf, R. J., Clarke, F. R., Sharpe,

642 A. G., Distelfeld, A., & Pozniak, C. J. (2017). High density mapping and haplotype analysis  
643 of the major stem-solidness locus SSt1 in durum and common wheat. *PLOS ONE*, 12(4),  
644 1–19. <https://doi.org/10.1371/journal.pone.0175285>

645 Nilsen, K. T., Walkowiak, S., Xiang, D., Gao, P., Quilichini, T. D., Willick, I. R., Byrns, B.,  
646 N'Diaye, A., Ens, J., Wiebe, K., Ruan, Y., Cuthbert, R. D., Craze, M., Wallington, E. J.,  
647 Simmonds, J., Uauy, C., Datla, R., & Pozniak, C. J. (2020). Copy number variation of  
648 *TdDof* controls solid-stemmed architecture in wheat. *Proceedings of the National Academy*  
649 *of Sciences*, 117(46), 28708–28718. <https://doi.org/10.1073/pnas.2009418117>

650 Pedersen, T. L. (2022). *patchwork: The Composer of Plots*. [https://CRAN.R-](https://CRAN.R-project.org/package=patchwork)  
651 [project.org/package=patchwork](https://CRAN.R-project.org/package=patchwork)

652 Peirce, E. S., Cockrell, D. M., Mason, E., Haley, S., Peairs, F., & Nachappa, P. (2022). Solid Stems  
653 and Beyond: Challenges and Future Directions of Resistance to Wheat Stem Sawfly  
654 (Hymenoptera: Cephidae). *Journal of Integrated Pest Management*, 13(1).  
655 <https://doi.org/10.1093/jipm/pmac023>

656 Peirce, E. S., Cockrell, D. M., Ode, P. J., Peairs, F. B., & Nachappa, P. (2022). Triticale as a  
657 Potential Trap Crop for the Wheat Stem Sawfly (Hymenoptera: Cephidae) in Winter  
658 Wheat. *Frontiers in Agronomy*, 4. <https://doi.org/10.3389/fagro.2022.779013>

659 Perdry, H., & Dandine-Roulland, L. (2018). Gaston—Genetic Data Handling (QC, GRM, LD,  
660 PCA) & Linear Mixed Models. *R Package*, 83, 1–29.

661 Poland, J., Endelman, J., Dawson, J., Rutkoski, J., Wu, S., Manes, Y., Dreisigacker, S., Crossa, J.,  
662 Sánchez-Villeda, H., Sorrells, M., & Jannink, J.-L. (2012). Genomic Selection in Wheat



663 Breeding using Genotyping-by-Sequencing. *The Plant Genome*, 5(3).  
664 <https://doi.org/10.3835/plantgenome2012.06.0006>

665 Poole, N., Donovan, J., & Erenstein, O. (2021). Viewpoint: Agri-nutrition research: Revisiting the  
666 contribution of maize and wheat to human nutrition and health. *Food Policy*, 100, 101976.  
667 <https://doi.org/10.1016/j.foodpol.2020.101976>

668 R Core Team. (2022). *R: A Language and Environment for Statistical Computing*. R Foundation  
669 for Statistical Computing. <https://www.R-project.org/>

670 Ray, D. K., Mueller, N. D., West, P. C., & Foley, J. A. (2013). Yield trends are insufficient to  
671 double global crop production by 2050. *PloS One*, 8(6), e66428.

672 Seamens, H. (1929). The Value of Trap Crops in the Control of the Wheat Stem Sawfly in Alberta.  
673 In *59th Annual Report Entomological Society of Ontario 1928*.

674 Shiferaw, B., Smale, M., Braun, H.-J., Duveiller, E., Reynolds, M., & Muricho, G. (2013). Crops  
675 that feed the world 10. Past successes and future challenges to the role played by wheat in  
676 global food security. *Food Security*, 5, 291–317.

677 VanRaden, P. M. (2008). Efficient methods to compute genomic predictions. *Journal of Dairy*  
678 *Science*, 91(11), 4414–4423.

679 Viera, A. J., Garrett, J. M., & others. (2005). Understanding interobserver agreement: The kappa  
680 statistic. *Fam Med*, 37(5), 360–363.

681 Weiss, M. J., & Morrill, W. L. (1992). Wheat Stem Sawfly (Hymenoptera: Cephidae) Revisited.  
682 *American Entomologist*, 38(4), 241–245. <https://doi.org/10.1093/ae/38.4.241>

683 Wickham, H., Chang, W., & Wickham, M. H. (2016). Package ‘ggplot2.’ *Create Elegant Data*  
684 *Visualisations Using the Grammar of Graphics. Version*, 2(1), 1–189.

685 Wickham, H., Hester, J., Chang, W., & Bryan, J. (2022). *devtools: Tools to Make Developing R*  
686 *Packages Easier*. <https://CRAN.R-project.org/package=devtools>

687 Winn, Z. J., Lyerly, J., Ward, B., Brown-Guedira, G., Boyles, R. E., Mergoum, M., Johnson, J.,  
688 Harrison, S., Babar, A., Mason, R. E., Sutton, R., & Murphy, J. P. (2022). Profiling of  
689 Fusarium head blight resistance QTL haplotypes through molecular markers, genotyping-  
690 by-sequencing, and machine learning. *Theoretical and Applied Genetics*, 135(9), 3177–  
691 3194. <https://doi.org/10.1007/s00122-022-04178-w>

692