

Project title: Methods in Marker Assisted Breeding in Octoploid Strawberry

Project number: CP163

Project leader: Richard Harrison, NIAB

Report: Final, April 2021

Previous report: Annual Report, September 2018

Key staff: Joe He

Location of project: NIAB EMR

Industry Representative: Tom Rogers, Soloberry Ltd.

Date project commenced: 01/10/2016

Date project completed 31/03/2021
(or expected completion date):

DISCLAIMER

While the Agriculture and Horticulture Development Board seeks to ensure that the information contained within this document is accurate at the time of printing, no warranty is given in respect thereof and, to the maximum extent permitted by law the Agriculture and Horticulture Development Board accepts no liability for loss, damage or injury howsoever caused (including that caused by negligence) or suffered directly or indirectly in relation to information and opinions contained in or omitted from this document.

© Agriculture and Horticulture Development Board 2021 No part of this publication may be reproduced in any material form (including by photocopy or storage in any medium by electronic mean) or any copy or adaptation stored, published or distributed (by physical, electronic or other means) without prior permission in writing of the Agriculture and Horticulture Development Board, other than by reproduction in an unmodified form for the sole purpose of use as an information resource when the Agriculture and Horticulture Development Board or AHDB Horticulture is clearly acknowledged as the source, or in accordance with the provisions of the Copyright, Designs and Patents Act 1988. All rights reserved.

All other trademarks, logos and brand names contained in this publication are the trademarks of their respective holders. No rights are granted without the prior written permission of the relevant owners.

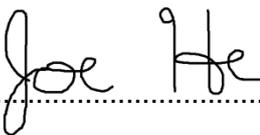
AUTHENTICATION

We declare that this work was done under our supervision according to the procedures described herein and that the report represents a true and accurate record of the results obtained.

Joe He

PhD Student

University of Reading

Signature ..  Date

[Name]

[Position]

[Organisation]

Signature Date

Report authorised by:

[Name]

[Position]

[Organisation]

Signature Date

[Name]

[Position]

[Organisation]

Signature Date

CONTENTS

GROWER SUMMARY	1
Headline.....	1
Background.....	1
Summary	2
Financial Benefits	3
Action Points.....	3
SCIENCE SECTION	4
Introduction	4
Materials and methods	8
Results.....	12
Discussion	19
Conclusions	19
Knowledge and Technology Transfer	21
References	22
Appendices	25

GROWER SUMMARY

Headline

Progress towards deployment of Genomic Selection (GS) as an advanced breeding technique in strawberries is ongoing.

Background

Strawberry breeders aim to generate novel genotypes that express traits suitable for the industry in their target region. Over the past 200 years, significant progress has been made in traits such as flavour, berry size, yield, disease resistance and cropping season duration. Current goals in strawberry breeding include improvements in maintenance of post-harvest fruit quality, yield, texture and flavour.

Traditionally, crossing is conducted based on identification of desirable traits in parental germplasm material. Offspring from a cross are assessed throughout the growing season and scored on a weighted index of favourable traits. The highest scoring individuals are selected to progress onto further larger scale trials, where additional information, such as yield and picking speed are gathered, and to confirm the presence of the favourable traits. Additionally, the selected genotypes are assessed for suitability across a range of environmental conditions, with particular focus on the target region. Overall, making crosses to release of a novel cultivar may take between 7 and 10 years.

Genetic markers are detectable features within the genome of a plant that may differ between individuals of the same species. Markers that are physically close to genetic variants controlling economically important traits tend to be co-inherited with the desirable genetic variant when the plant produces offspring, making some markers reliable proxies for these genes. Over the past 20 years, the number of known markers has dramatically increased and the cost of identifying them has greatly decreased. It is now possible to incorporate genomic information in the breeding process to aid breeders in selection of the optimal individuals.

GS offers a range of benefits relative to conventional breeding approaches. Firstly, it allows for greater selection accuracy as the confounding environmental effects on a trait can be eliminated. Secondly, it allows for strong selection on traits that are expensive or difficult to assess or selection on traits that are apparent only under rare environmental conditions. Thirdly, as multiple traits can be assessed, GS potentially allows selection at the juvenile stage, reducing the duration of the breeding cycle. Moreover, GS is particularly suitable for identification of traits that are controlled by many genes (polygenic traits) as its simultaneous regression of all markers on all traits reduces the likelihood of over/underestimation of effect

size. GS also potentially allows control of inbreeding and elimination of certain field experiments.

Summary

Strawberry is an economically important crop with global and UK production on an upwards trend. Strawberry breeding efforts attempt to generate novel varieties that have increased yields, resistance to pathogens, good eating quality and high nutritional content. Genomic prediction (GP) is an advanced marker assisted prediction (MAP) technique that makes predictions about agronomically important traits in crops. Three areas of improvement were identified to assist commercial deployment of GP for strawberry.

Breeding efforts currently rely primarily on visual and mechanical measurement of plant phenotypes, which is slow, imprecise and liable to human biases. A strawberry phenotyping platform was developed that captured images from 360 degrees around the strawberry fruit to generate 3D representations. Seven fruit quality traits were calculated from the representations, which showed good concordance with manual measurements. Deployment of the system could lower phenotyping costs, increase throughput, increase precision and thus improve GP accuracy.

Current genotyping approaches for dense marker panels in strawberry are too expensive for commercial deployment. A rational design process was implemented to generate amplicon sets that would genotype a panel of markers to maximise information for GP, with scalability to accommodate resources available to different breeding programmes. The design process failed to generate marker information due to unexpected interaction in the multiplexed PCR reactions.

The relative effectiveness of phenotypic prediction and MAP in strawberries is unclear. Moreover, existing models of GP in strawberry does not represent all the traits of interest to breeders. Between years predictions of 15 fruit quality traits were implemented using phenotype only, traditional MAP (tMAP) and GP models. GP had similar selection accuracy compared to phenotype only prediction, but tMAP performed significantly worse than the other models. It was concluded that GP would likely yield benefits to strawberry breeding in the context of speed breeding. Models for GP for 15 strawberry fruit quality traits are available for breeders to deploy in their breeding populations.

Financial Benefits

Due to lack of evidence of the efficacy of the proposed genotyping approach at £5 - £10 per sample, it is assumed that genotyping costs will be using the 35K SNP array Axiom IStraw35 384HT array. Under these conditions, it is unlikely that GP can be cost-effectively integrated into UK strawberry breeding programmes. Further research is needed to decrease the cost of genotyping strawberries whilst maintaining efficacy of GP.

Action Points

Phenotyping of seven external strawberry fruit quality traits can be sped up five-fold, reducing labour costs, using a novel automated 3D image capture and analysis platform. Where genotypic data exists, 15 strawberry fruit quality traits can be predicted with higher accuracy than tMAS using the described novel models.

These techniques need further investigation before they can be effectively integrated into current breeding programmes.

SCIENCE SECTION

Introduction

Genomic selection (GS) is an advanced breeding technique that utilises a densely genotyped and phenotyped training population, from which associations are made relating the magnitude and direction of quantitative trait loci (QTLs) associated with agronomically important traits (Meuwissen, Hayes, Goddard, et al. 2001). GS has been successfully deployed in a range of crops, including grape (Viana et al. 2016), wheat (Heffner, Jannink, and Sorrells 2011; Thavamanikumar, Dolferus, and Thumma 2015), maize (Shikha et al. 2017) and strawberry (Gezan et al. 2017). Deployment requires a training population, which is densely genotyped and phenotyped for the agronomically relevant traits. A statistical model is developed, which associates the genotype and phenotype. Solely on the basis of the genotype and statistical model, breeding values for breeding material is estimated and selections are made (Heffner, Sorrells, and Jannink 2009; Meuwissen, Hayes, and Goddard 2001).

There are a range of benefits associated with GS. Firstly, assuming that there are sufficient markers available, GS has been demonstrated to generate greater prediction accuracy than conventional selection. This is largely due to the approach ignoring the variable and non-hereditary environment (De Los Campos et al. 2009). Secondly, GS allows the regression of a genotype onto multiple individuals, allowing increase in power of detection of small effects. Moreover, this allows selection on rare, expensive, or otherwise difficult to phenotype traits. For example, consider a rare event, such as a harsh drought in the summer that is expected to occur only every 10 years. It is possible that this does not occur during the entire breeding cycle of a novel cultivar, and thus no phenotypic information about crop performance can be gathered. Under conventional selection, no predictions can be made regarding plant performance under harsh drought, but under genomic selection, predictions can be made based on the genotype of the novel cultivar and data gathered from genotypes from all instances of harsh drought (Heffner et al. 2010; Peace et al. 2017). Perhaps most importantly, GS allows for the reduction of the duration of the breeding cycle as plants can be genotyped and selected for in the seedling phase. Additionally, GS allows for elimination of some field experiments and better planning of crosses by providing information on relatedness (Gezan et al. 2017). GS performs better than marker assisted breeding (MAS) because it simultaneously estimates the effect size for all markers, and thus is less likely to overestimate the "significant" and underestimate the "insignificant" MAS markers.

Various models have been deployed to predict plant performance based on previous trait data, with linear models being commonly used. There are a large number of classes of equations to model the effects of genetics on phenotype, but often the solutions to such

systems are mathematically challenging and requires knowledge of constants which may be difficult to estimate. There is no particular reason to assume that the effect of the variation in genetic data on phenotype is linear, but restricting the class of equations to linear additive models simplifies the modelling and allows utilisation of a range of tools already developed for linear analysis. Typically, the only constants that need to be known are the first and second moments of the variables to be estimated (Henderson 1975). Moreover, empirically, linear models have shown success in prediction of plant and animal performance in breeding populations (Bates et al. 2015).

Linear mixed model (LMM) is often deployed to predict plant traits based on data gathered from previous years or other populations. The LMM assumes the observed phenotype is a linear combination of a set of fixed and random effects.

$$y = Xb + Zu + e$$

where y is a vector of observed phenotypes, X is a design matrix for the fixed effects, b is a vector of fixed effects, Z is a design matrix for the random effects and e is the error. The random effects are assumed to be drawn from a normal distribution with mean 0 and a known variance covariance matrix $u \sim MVN(0, G)$; $e \sim MVN(0, R)$. In the context of plant breeding, G represents the genetic covariance between the individuals and R represents their environmental covariance.

In cases where the number of effects to be estimated differs from the number of unique measurements, then no unique solution can be obtained (more precisely, in the system $Ax = b$ with augmented matrix $[A|b]$, a unique solution exists if and only if $rank[A] = rank[A|b]$).

Under this model, it is common to treat plant performance as a random variable, making predictions with the best linear unbiased prediction (BLUP) (Molenaar, Boehm, and Piepho 2018). Originally developed for animal breeding, BLUP models the genotype effect on the observed phenotype as a random variable, with other effects including location, block and experimental year classed as fixed variables. BLUPs have properties that are desirable in prediction models: of all the linear models where predictions are unbiased, it has the minimum variance. Additionally, it incorporates shrinkage towards the mean, which is a desirable statistical property of an estimator, as it increases accuracy, leading to a smaller mean squared error (Piepho et al. 2008). BLUP and best unbiased linear estimation (BLUE) solutions to the linear mixed model can be computed using Henderson's mixed model equations (Henderson 1984).

Mean and variances of the plant traits are to be estimated using the available data, with maximum likelihood (ML) being the most popular method. Assuming that the phenotypic measurements are drawn from a normal distribution, we can calculate the likelihood of μ and

σ^2 given the data. Maximising the likelihood (or log-likelihood in practice as this function is monotonic and easier to perform calculus on) gives us the maximum likelihood (ML) estimate of σ^2 . It is known that the ML estimation of variance is biased downwards in cases where a finite population is sampled. The bias in ML variance estimation arises due to the loss of one degree of freedom required in estimation of each of the fixed effect parameters (i.e. the mean) (Foulley 1993). The restricted maximum likelihood (REML) maximises a modified likelihood that has no mean component and thus avoids bias.

In organisms that cannot be clonally propagated and exhibit sexual dimorphism, such as dairy cattle, best linear unbiased prediction (BLUP) has been used for decades to select sires with the highest estimated breeding values (based on measurements of the offspring from previous matings with each sire) to breed superior families (Henderson 1984). For plants, BLUP has been deployed to breed for various traits in ornamental carnation and geraniums, finding selection based on family indices worked at least as well as individual selection (Molenaar et al. 2018). In potato, BLUP was deployed to breed for resistance to late blight resistance (Sood et al. 2020). BLUP was also deployed to predict for expansion volume and yield and select families in maize (Viana et al. 2011).

When marker information is available, MAP can be implemented. MAP estimates the magnitude of effect of some marker(s) on the trait and estimates the plant trait (in the case of a linear models) as the sum of the effect of each marker that an individual has. The Kruskal Wallis (KW) test, a non-parametric implementation of ANOVA, tests if samples originate from the same distribution and can be applied to each marker to determine if it is associated with differences in traits (Broman 2003). The generated H value is to be compared with the KW distribution (with an appropriate number of degrees of freedom) to determine a p -value, but when the number of groups or individuals are large, computation of the KW distribution becomes computationally infeasible. It remains unknown the best method to approximate the KW distribution when groups are large (as in the case of field experiments where hundreds of genotypes are assessed), but the χ^2 distribution is often used as an approximation (Meyer and Seaman 2008).

When marker data are not densely or evenly spread throughout the genome, power to detect QTLs in sparse regions falls (Akond et al. 2019). Consideration of QTLs in small intervals, using nearby markers increases the power to detect such QTLs and is termed interval mapping (IM) (Lande and Thompson 1990). A logarithm of odds (LOD) score measures the likelihood that a particular interval is associated with a QTL. A LOD score of 3 is often considered evidence of a true marker, but the probability of a false positive is dependent on the number of markers, size of the intervals considered and the heritability of the trait.

Frequently, an empirical test for significance is performed by bootstrapping with a permutation test (Churchill and Doerge 1994).

A range of GP models have been deployed on both simulated and experimental datasets. The first model for GS utilises BLUPs to predict effects of the markers on traits (Meuwissen, T. H. E. , Hayes, B. J., & Goddard 2001). Genomic BLUP (GBLUP) utilises the linear model:

$$y = \mu + \sum_{i=1}^n Wqi + e$$

Where y is the observed phenotype, μ is the mean, W is the genotypic design matrix, qi is the effect of each SNP and n is the total number of markers. The variances of each SNP is assumed to be equal. As the number of markers (effects) is typically greater than the number of phenotypic records (measurements), the system is typically underdetermined and some form of regularisation is deployed to solve the system. Although stepwise algorithms exist to select a subset of markers (Habier, Fernando, and Dekkers 2007), this method remains highly biased when strongly correlated markers are present. Ridge regression adds λ to the least squares estimator as a penalty, shrinking the effect size of each marker equally towards 0 to overcome the underdetermined system, whilst still using all markers (Piepho et al. 2008). The ridge regression parameter controls the magnitude of the penalty and parameterises the relative importance of the data-dependent empirical error (Ogutu, Schulz-Streeck, and Piepho 2012). If variances of markers are assumed to be equal, λ is the ratio of the residual and marker variances, usually estimated through maximum likelihood methods(Endelman 2011). Under this model, when residual errors are large, more shrinkage allows for control of bias, and when marker effects are large, shrinkage is reduced to allow for estimation of true positives.

Although most experimental implementations of GS deploy multiple models and assess the prediction accuracy of the different models, in general, different models of GS perform similarly. Comparison of 11 GP models found that most models generated similar accuracies, with slightly better performance when deploying RKHS (Heslot et al. 2012). Comparison of GBLUP, RKHS and BayesCII for wheat yield showed little difference in selection accuracy(He et al. 2016). Comparing three different Bayesian models, no significant improvement in selection accuracy was observed (Habier et al. 2011). In strawberry, selection accuracies were slightly higher using RKHS than BayesB or GBLUP, but other factors had greater effects (Gezan et al. 2017).

When dealing with polyploids, such as strawberry, markers may be in LD with QTLs on only one subgenome. When the resolution to a homeologous subgenome is unclear, the linkage of a detected marker to the QTL may be unclear. Although effort has been deployed to

generate phasing techniques for polyploids (He et al. 2018), GP has also been deployed in polyploids with little modification of models. To evaluate the utility of models, comparisons of selection accuracy between models can be made, as well as comparisons to estimations of heritabilities of traits. Under a randomised block design, there is assumed to be no genotype by environment correlations (Kruijer et al. 2014).

In order for genomic selection to be viable for commercial implementation, it must be more cost effective than the currently employed method. Tools have been developed to perform cost-benefit analysis and to optimise resource allocation for implementation of genomic selection. Analysis with DeltaGen on a forage breeding population, for example, suggests that genotyping for genomic selection approximately doubles the cost, but also nearly doubles the increase in genetic gain per cycle when compared to selection without genotypic information (Jahufer and Luo 2018). Genomic selection has been experimentally implemented in strawberry utilising the IStraw90 Axiom SNP array (Bassil et al. 2015) to generate genotypic information. High prediction accuracies were observed for a range of agronomically important traits, but it was acknowledged that the cost of the SNP array was likely too high for commercial deployment (Gezan et al. 2017).

There were two main aims in this section. Firstly, based on a biparental strawberry mapping population, three between years prediction approaches (phenotype only, tMAS and GP) of 15 strawberry fruit quality traits relevant to breeders were to be assessed. Secondly, biological correlations and efficacy of selection was to be computed. Together, these datasets offer models for strawberry breeders on the methods and traits suitable for selection.

Materials and methods

The biparental mapping population was used previously for genetic mapping (Antanaviciute 2016). Briefly, 188 seedlings were raised from a cross between two *F. ananassa* cultivars 'Redgauntlet' and 'Hapil', of which 120 were randomly selected. These individuals were clonally propagated with six replicates in the Autumn of 2015. Additionally, the parental genotypes and two check varieties, 'Sonata' and 'Elsanta', were included in the experiment, making a total of 744 individuals. The experiment took place at East Malling Research, at 51° 17'15"N 0° 27'12"E.

Seedlings were distributed in a randomised block design within three tunnels, with three beds per tunnel and two rows per bed. Each block was one-and-a-half rows. Seedlings were planted in a double row zig-zag 35cm high, 50cm wide with 40cm between plants. Plants were allowed to establish over winter and dead material were removed in early 2016 and again in May 2016. Irrigation and fertigation was installed and performed according to conventional practice. Plants were also sprayed against common pests and diseases

according to common practice. Harvesting took place three times a week (Monday, Wednesday and Friday) from when the first fruits developed until all fruits were harvested (17/06/2016 - 21/07/2016). In each harvest, all ripe fruits were collected from all plants for phenotypic analysis, except during the peak season, where only one or two tunnels were harvested for logistical reasons; in any week, all plants were harvested at least once. Harvesting was initiated in early morning at approximately 05:00 and classed *in situ* as marketable or unmarketable before delivery to a centralised location, where phenotyping of other traits took place.

Phenotypic assessment of the plants took place on the same day as harvest, except during peak periods, where assessment took place over the day of harvest and the day after. Where assessment took place the next day, fruits were stored at 4 °C overnight. Assessment was conducted using a modified RosBREED protocol for strawberry, with their standards defining the extremes and midpoints of the fruit quality traits where applicable (Mathey et al. 2013). Assessment was primarily conducted by J. He and A. Karlstrom with occasional assistance from others. For all individuals, examples of phenotypes corresponding with measurements on the appropriate scales were demonstrated and agreed before phenotyping took place.

A total of 15 traits were assessed. Marketable and unmarketable fruits were collected and weighed separately for each plant, and summed for all harvests throughout the season. For each other phenotype, a single value was generated from assessment of ten randomly selected fruits (where available) from the marketable portion of each plant at each harvest, except pH, soluble solids, and firmness, where three, twenty and ten fruits were randomly selected for analysis over the season respectively.

pH was measured by releasing a drop of strawberry juice onto a pH meter; firmness was assessed by gentle depression of the fruit by a robotic arm and measurement of the deformation (Firmtech Umweltanalytische Produkte GmbH). Soluble solids content was measured by releasing a drop of juice from a randomly selected fruit onto an interferometer; cap size was a visual assessment of the width the cap relative to the neck of the fruit; appearance was a visual assessment of the fruit ranging from very malformed to symmetrical and attractive; external colour was a visual assessment of the fruit colour; glossiness was a visual assessment of the shine of the fruit; achene position was a visual assessment of how protruding the achenes were; seediness was a relative measure of the density of visible achenes; fruit shape was a visual assessment of the ratio of fruit height to width; neck line was visual assessment of the shape of the neck; skin strength was the number of fruits with broken skin after ten fruits were rubbed gently with a thumb; and internal colour was the relative colours of the inside of the fruit after bisection (Mathey et al. 2013). In addition to data collected in 2016, phenotypic data from a previous study on the population from the 2013 to

2015 were included in the analysis. The 2013 - 2015 dataset was defined as the prediction dataset and the 2016 dataset was defined as the validation dataset.

DNA extraction for genotyping was performed on single young leaves using the Qiagen DNeasy kit according to the manufacturer's instructions. Genotyping was performed using the 90K array with genotypes being calculated in accordance to the manufacturer's instructions (Cockerton, unpublished). Filtering was performed on the dataset to remove non-segregating markers and remove redundant information by maintaining only one instance of markers that segregated identically. After filtration, 3436 segregating markers were identified and included for genotypic analysis. Genetic rogues, defined as individuals with non-parental genotypes or individuals that were genetically identical to apparently other genotypes were excluded from analysis. After data filtration, 66297 phenotypic records were used in the training dataset and 24908 phenotypic in the validation set, amounting to a 77% and 23% data split respectively.

116 individuals including the parental and check varieties were included for phenotypic analysis. To explore the data, the phenotypic values were plotted against the genotypes, with the parental and check cultivars highlighted. For all traits, the mean across all blocks in the years was computed, except marketable yield and unmarketable yield, where the sum of all records were computed, and pH, where the mean of the concentration of hydronium ions was calculated, and the result converted to the logarithmic pH scale. The differences in rank of the parental strains were also computed.

In the absence of genotypic data, given the unbalanced data, it is conventional to deploy BLUP to predict plant performance. Calculation of the BLUP was performed using the 'lmer' command from the 'lme4' package in R (Bates et al. 2015). The prediction model included all data from 2013 - 2015, treating the year and blocks as fixed effects and the genotypes as random effects. It was assumed that there was no differential interaction effect between genotype and block or year. In order to compute a comparable figure for the validation dataset, a similar linear model was fitted for the 2016 data, treating blocks as fixed effects and the genotypes as random effects. The variance estimation method for both models was 'REML'. Random effects were extracted from the model using the 'ranef' command and their Pearson's correlation coefficients were computed using the 'cor' function as a measure of prediction accuracy. The concordance correlation coefficient (CCC) of the predictions were also computed using the 'CCC' function from the 'DescTools' package (Signorell et al. 2021). To estimate correlations between traits, the Pearson's correlation coefficients of BLUPs for every pair of traits from the prediction, validation and total datasets were also computed. *p*-values under the null hypothesis that the correlations were not different from 0 were calculated

and a Bonferroni correction was performed using the number of pairwise tests performed ($p < 0.05$, $n = 315$).

103 progeny were included in genotypic analyses. In order to maximise power to detect markers associated with QTL, all phenotypic data from all years were included in the marker discovery phase. Two methods of marker discovery were implemented: the Kruskal Wallis (KW) test and interval mapping (IM). Marker discovery for both methods were conducted using MapQTL5 in accordance to the user manual (van Ooijen 2009). The mean of all traits were calculated as input for the 'qua' file, except marketable yield and unmarketable yield, where the sum of all records were computed, and pH, where the concentration of aqueous hydronium ions were calculated.

For the KW analysis, the resulting H statistic (and their associated degrees of freedom) was extracted from MapQTL5 for p value estimation and multiple testing correction. As an approximation to the KW distribution, the H statistic was compared to the χ^2 distribution with the appropriate number of degrees of freedom to yield a p value for each marker being associated with a QTL. Computation of the χ^2 distribution was performed using the 'pchisq' function from the 'stats' package in R. The Benjamani-Hochburg (BH) correction was applied to adjust for multiple testing of markers to control false discovery rate. The critical value for false discovery rate was set to an exploratory rate of 0.2. Computation of the BH correction was performed using the 'p.adjust' function, also from the 'stats' package in R.

For the IM test, significance thresholds were first generated. A permutation test was conducted using the 'permutation test' function of MapQTL5. 100 permutations were simulated for all traits and the 95th percentile of the genome-wide significance levels were taken as the threshold for a statistically significant marker. As markers physically close together are likely to be in LD with each other as well as QTLs, a simple clustering algorithm was implemented to determine if a set of markers with significant LOD values described the same QTL. When a LOD peak was identified at a locus, all other markers and intervals were scanned from both directions until a marker was identified with a LOD score 2 units less than peak. All markers scanned were clustered as describing the same QTL, with the marker with the highest LOD score selected as representative of that peak.

Markers from KW and IM were pooled and used for prediction of plant performance. Prediction was calculated as the mean of the trait values estimated by MapQTL5 for a given allele for each marker. Prediction accuracy was defined as the Pearson's correlation coefficient between the tMAP estimations of the prediction data and the BLUPs of the validation data as described previously. Additionally, the CCC between the values was calculated.

GP was implemented through a two step process. In the first step, BLUPs were calculated for the prediction and validation data as described previously. GBLUP was conducted using the 'mixed.solve' function from the 'rrBLUP' package (Endelman and Jannink 2012). The Y and Z matrices were defined as the phenotypic BLUPs and the design matrix of the genotypic markers respectively. The predictor of an individual is the GEBV and was defined as the sum of all corresponding marker effects of the individual. To evaluate accuracy, two methods of prediction accuracies were calculated. The BV BLUP cor method correlates the GEBV of the training population with the BLUPs of the validation population and the GEBV cor method correlates the GEBVs of the training population and the GEBVs of the validation population.

Results

To explore the range of data, the mean phenotypic traits of each genotype along with the parental and check varieties were plotted. Comparison with traits described in literature may be unreliable because the scores of individuals are dependent on their environment. For example, the average colour for 'Sonata' and 'Elsanta' were found to be statistically different, but almost indistinguishable to the human eye. Total sugar content in 'Sonata' is higher than 'Elsanta' when fully irrigated, but not statistically distinguishable when under water stress. Similarly, 'Sonata' contains more total acid than 'Elsanta' when fully irrigated, but statistically indistinguishable when under water stress (Giné Bordonaba and Terry 2010).

However, consistent with expectations of a modern cultivar on the market, the marketable yield of 'Elsanta' and 'Sonata' were both high, with 'Sonata' being markedly higher than any other genotype. 'Sonata' scored significantly higher than any other genotype for appearance with 'Elsanta' scoring moderately. Both check varieties scored highly in glossiness, and skin strength, both of which are desirable traits for the market. Moreover, both check varieties were significantly firmer than other genotypes, a trait that breeders select for as firmness correlates with post-harvest shelf life (Salentijn et al. 2003). 'Sonata' is described as producing oblate fruits, and this trait is apparent in its low shape and neck line score. Interestingly, both check species have relatively low brix values, despite consumers indicating sweetness as an important trait and neither have low unmarketable yields. The latter observation can be explained as marketable yield and unmarketable yield are correlated and commercial cultivars are expected to have high yields. Taken together, the data distributions are consistent with known and expected traits of the check varieties, suggesting the measured phenotypes are representative of the phenotypes of strawberries.

In the case of glossiness, achene position and pH, parental phenotypes have a large rank difference (> 70), suggesting that the parents have differing alleles controlling the traits, with

offspring inheriting heterozygously, displaying intermediate phenotypes. In the case of unmarketable yield, appearance, redness, seediness, shape and brix, the difference in rank of the parents is small (< 25). The distribution of achene position, shape and internal colour shows progeny with marked extremes. This is typical of traits under the control of a few large effect QTLs as segregation pyramids those QTLs by chance in a few individuals. Thus, it may be expected to identify large effect markers using tMAP for these traits.

Fifteen statistically significant correlations between traits were identified based on analysis using all data. Of these traits, the strongest was between redness and internal colour (0.82), followed by neck line and shape ($r = 0.51$), marketable yield and unmarketable yield ($r = 0.50$) and skin strength and firmness ($r = 0.50$). 14 of the 15 identified correlations were positive, with a single negative correlation identified between firmness and neck line. In order to investigate the reliability of correlations between the validation and prediction datasets, pairwise correlations of traits were also computed for the prediction and validation datasets. The three pairs of traits that had the strongest correlations were consistently correlated across years; internal colour and redness, neck line and shape and marketable yield and unmarketable yield (Figure 1). In an experiment in Bangladesh, strong pairwise correlations were found between total fruit weight, fruit length, fruit diameter and brix values (Mehraj and Jamal Uddin 2014). However, in this study, no such correlations were observed. Consistent with the correlations observed in previous years (Antanaviciute 2016), a strong correlation was observed between redness and internal colour between redness and glossiness and but not between cap size and shape. However, no other consistent correlation was observed and many fewer correlations were observed in this experiment. This may be due to differences in the method for calculating correlations and it is unclear if multiple testing corrections were applied for the result of Antanaviciute.

Correlations and concordances were computed between BLUPs of the prediction and validation dataset (Figure 2). High correlations (> 0.7) were found for firmness, neck line and redness. Low correlations (< 0.4) were found in glossiness, marketable yield, skin strength and unmarketable yield. High concordances (> 0.5) were found for redness and shape while low concordances (< 0.2) were found for appearance, firmness, glossiness, internal colour, marketable yield, pH, skin strength, and unmarketable yield.

In the case of traits such as brix and seediness, concordance was similar to correlation, indicating the absolute values were similar where there was correlation. In cases such as marketable yield, internal colour, pH and skin strength, the concordance was much lower than correlation, indicating that even when there was correlation, the absolute values of these traits did not match. In the case of marketable fruit, this is likely a reflection of the significant effect of environment on the trait; in the case of internal colour, this may be a rater effect.

Using IM, one marker for achene position, redness and unmarketable yield each were identified, on LG4B at 53.0cM, LG7A 49.9cM and LG1C at 9.9 cM respectively (Figure 3). The LOD values for each marker were 3.91, 4.14 and 5.14 respectively, exceeding the permutation test thresholds of 3.4, 3.5 and 5.1 respectively. The mean genotypic information content for each marker were 0.99, 1.0 and 0.99 respectively, with each marker explaining 17.0%, 15.5% and 21.0% of the variance respectively.

Using the KW test, 2 markers on LG3C at 90cM, 4 markers on LG4C at 64cM and 6 markers on LG7A at 52cM were identified, all for internal colour. The significance of the markers, after controlling for false discovery were between 0.15 and 0.18. Deployment of tMAP generated predictions for plant performance. The highest prediction accuracy was for internal colour (0.45) with the lowest being for unmarketable yield (0.02). Interestingly, internal colour was a trait for which markers were found, consistent with the skewed phenotypic distribution described.

Prediction and concordances of the tMAP model were generally poor. Of the studied traits, markers were identified for only four traits. Comparison with a similar genome wide association study from a previous study on the same population found no markers for the same traits on the same chromosome as identified in this study. More markers were identified in the previous study as multiple testing correction did not appear to have been applied. Inconsistencies in identified markers across different experiments is common in genome wide association studies as markers are only reported when they exceed a critical threshold, resulting in overestimation of the effect size (Xu 2003). In order for markers identified in this experiment to be deployed in MAS, validation of their predictions should be performed in other populations and in different environments.

Deployment of GBLUP generated GEBVs for all traits (Figure 4). High (> 0.6) prediction accuracies could be achieved for cap size, firmness, internal colour, neck line, redness and shape. Glossiness, skin strength, marketable yield and unmarketable yield had low (< 0.4) prediction accuracies. The highest prediction accuracy was for redness (0.74) with the lowest being for unmarketable yield (0.08). The concordances between GEBVs estimated using the prediction and validation datasets were close to zero (data not shown). This is likely due to significant shrinkage, so slight differences in means between validations and prediction datasets would result in lack of concordance. Interestingly, no predictive model performed well for unmarketable yield, perhaps due to its low heritability.

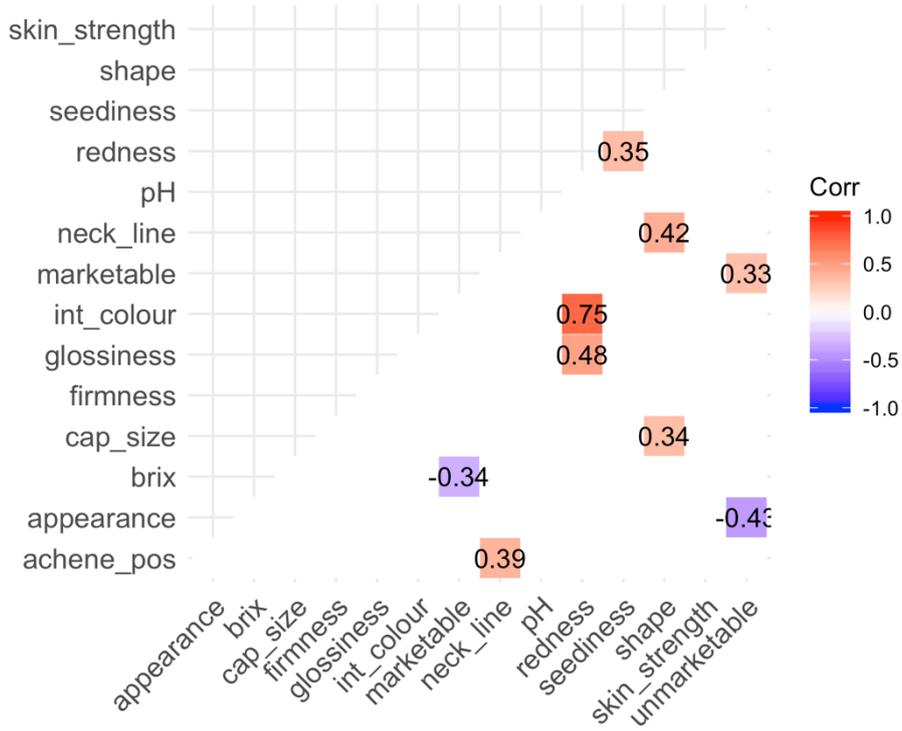
In theory, MAP in this study must be at least as good as phenotype only prediction when available information is utilised optimally. This is because the data included in the phenotype only prediction models include phenotype data only, whereas the MAP models include

identical phenotype as well as genotype information. If there is no relationship between marker data and plant performance, then optimal use of MAP ought to be identical to phenotype only prediction; any predictive power that the genotypic information has improves the model.

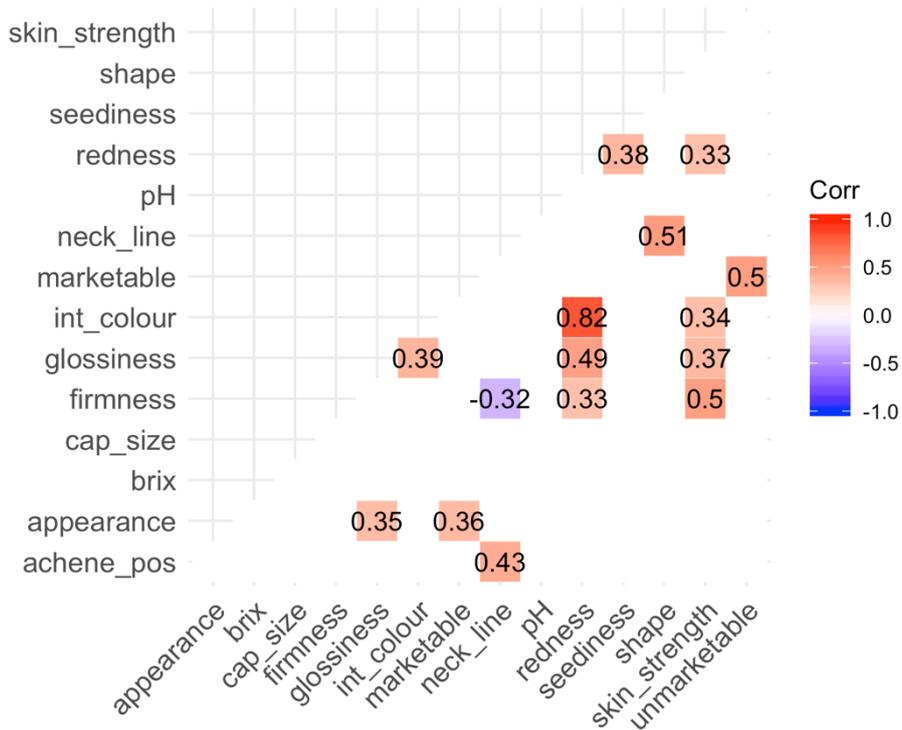
One possible reason that the tMAP models explored here perform poorer than phenotype only prediction is the introduction of bias. tMAP performs the poorest of all the models in the traits measured. One source of bias in tMAP is the accept/reject nature of marker identification, which assigns an effect to markers that fall above a critical threshold, whilst rejecting effects that fall below. In this approach, markers with small effects, which together may account for a significant proportion of the variation, may not be identified, thus biasing the effect of the markers. One expected effect of this, which was observed in this study, is that markers for some polygenic traits cannot be identified, making tMAP incapable of predicting performance.

In the cases of most traits, prediction accuracies between phenotype only predictions and GEBVs between the training and validation populations are similar. This indicated that there is little additional information that genotype data add to make predictions. The results presented may underestimate the performance of GP in a strawberry breeding population. In a real breeding population, there is potential to leverage a much larger training population including individuals that were genotyped/phenotyped in previous years and other locations because their relationship with the breeding population is known through shared markers.

2013-15 Fruit Quality Correlations



2013-16 Fruit Quality Correlations



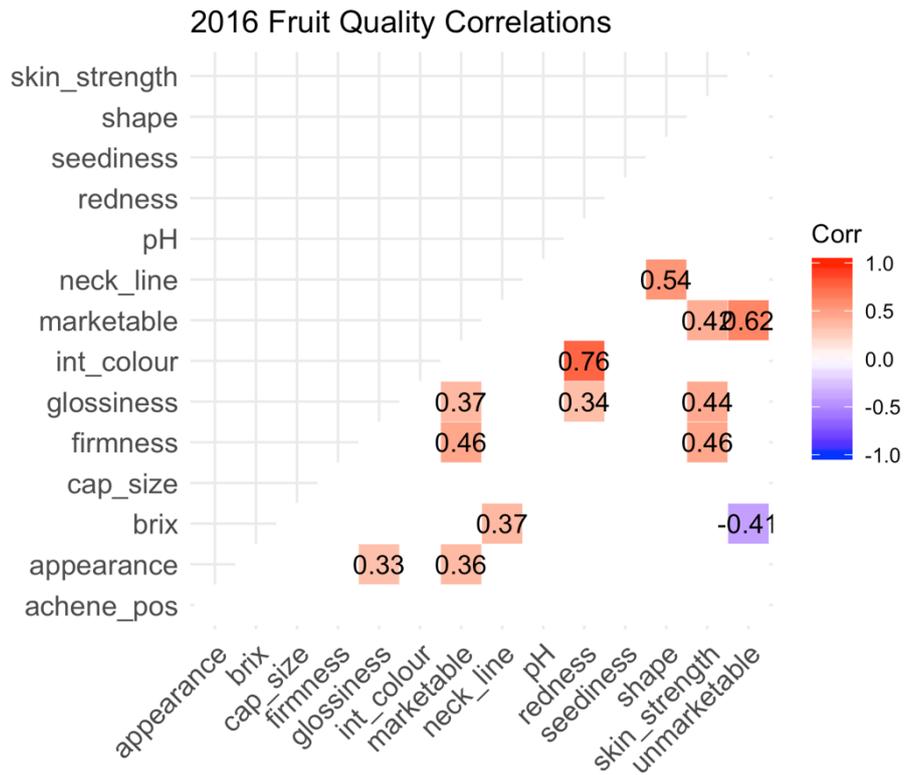


Figure 1 Correlations of 15 strawberry traits across three years

Correlation and Concordance for Phenotype only Prediction

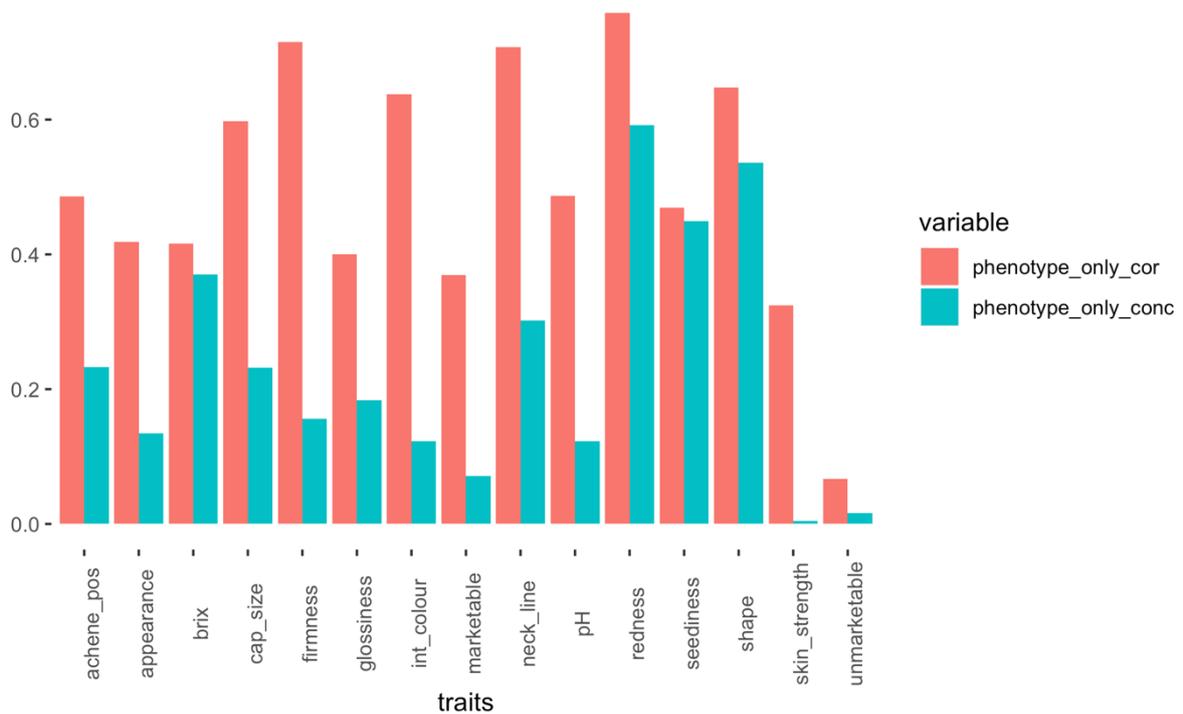


Figure 2 Phenotype only prediction between years for 15 strawberry traits

Concordance and Correlation of Conventional Marker Assisted Prediction

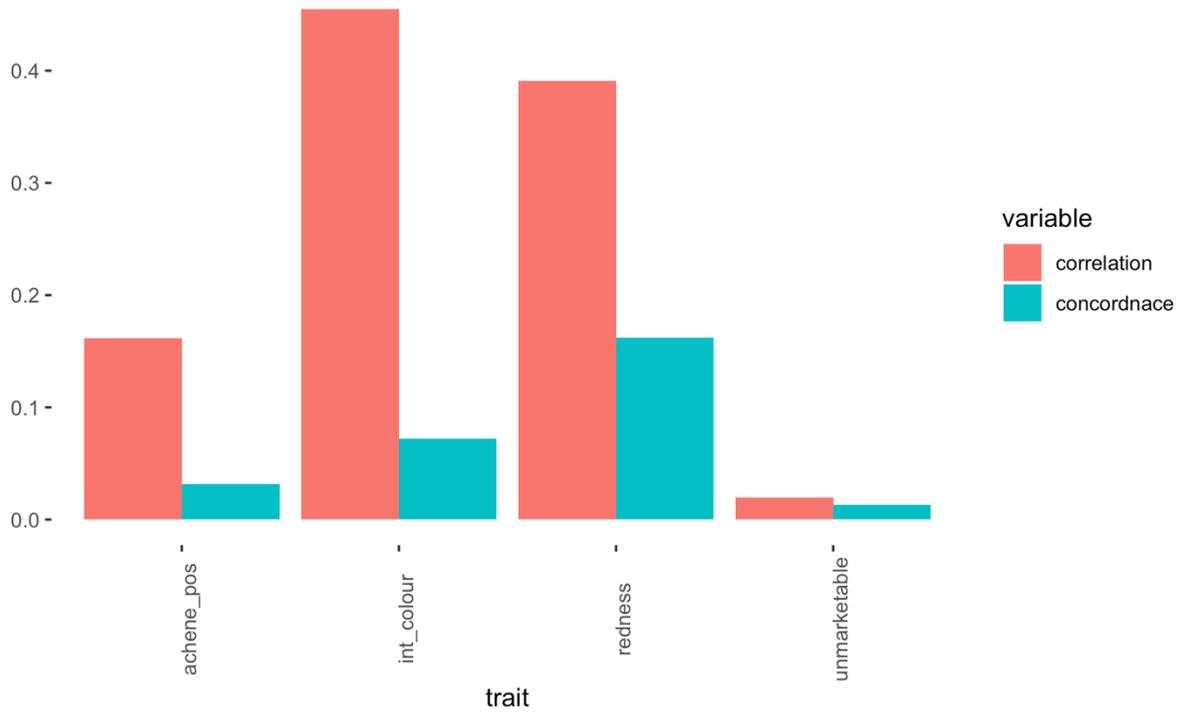


Figure 3 Concordance and Correlation from tMAP for strawberry traits

Correlations for performance prediction

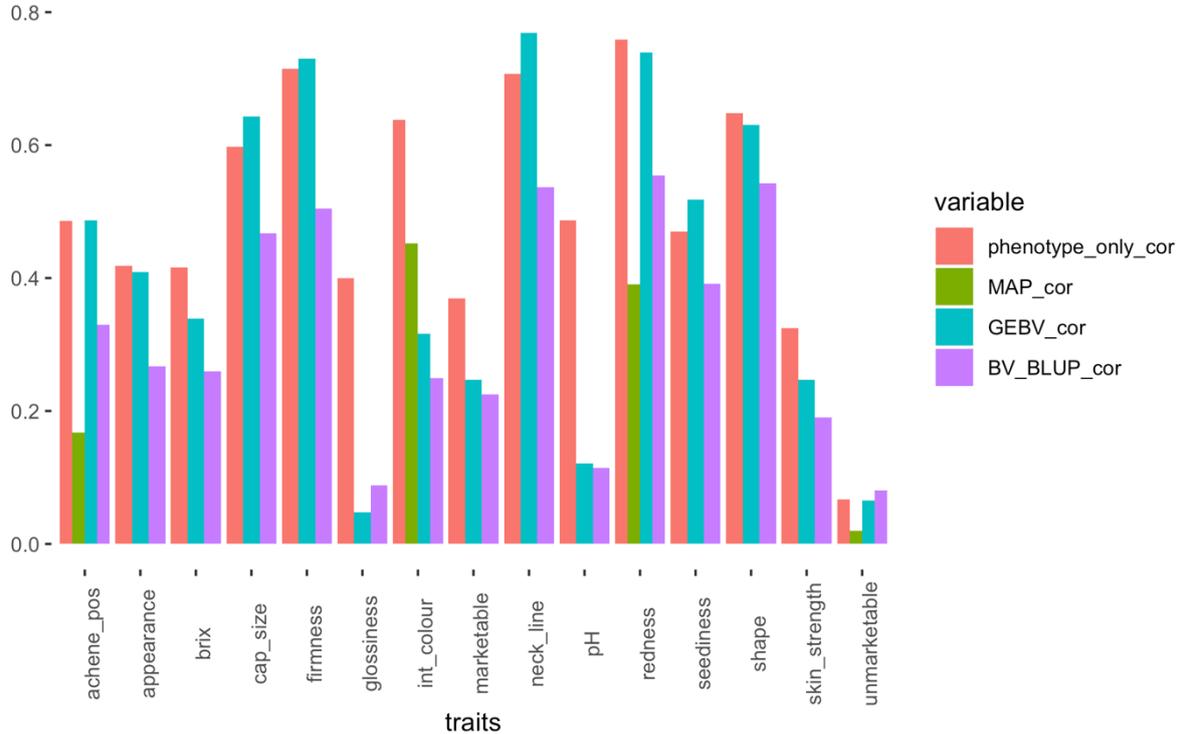


Figure 4 Correlations and concordances of three between years prediction methods for strawberry

Discussion

GP models for strawberries have been developed, which supports genetic and varietal development. A methodology for cheap and reliable genotyping for strawberry has been developed, though this approach lacks experimental validation. Additionally, an automated strawberry phenotyping platform has been developed, which aids throughput in breeding trials.

Conclusions

With the global population projected to reach nearly 10 billion in 2050, changes in the composition of human calorie consumption is needed to ensure sufficient food for all, including significant transition to plant based calorie intake (Berners-Lee et al. 2018). GS has the potential to increase the rate of genetic gain in breeding efforts and thus contribute towards feeding the projected population growth. Strawberries are a valuable commodity that are a source of nutrients, notably vitamin C and manganese as well as associated with protective effects for cancer and cardiovascular disease. Production and consumption in the UK, and the rest of the world are on an upwards trend.

Commercial strawberry breeding must generate money from its breeding efforts, and while the amounts that are generated depend on the nature of the breeding program, it is also dependant on the quality of their output, novel strawberry varieties. In this thesis, research is presented in three areas to improve deployment of GP in strawberry breeding and thus potentially improve the quality of novel strawberry varieties. The automated high-throughput 3D phenotyping platform increases the precision and reliability of measurement of seven strawberry fruit quality traits, which should allow for more accurate GP and reduction of labour costs associated phenotyping. Use of objective measurement scales also allows for easier interpretation of results. The rational design of amplicon sets for GP attempts to integrate parameters likely to be informative for GP in an open, scalable, genotyping system. This reduces cost by measuring only the most informative markers and allows scaling of genotyping effort dependant on the resources of the breeding program. The deployment of predictive models in a biparental strawberry population serves as experimental validation of the efficacy of GP in strawberry breeding compared to MAP and phenotype only prediction approaches, offering models for breeders to utilise.

It appears from simulation studies that genotyping may not be cost effective in saving costs associated with maintaining plants in the field, even assuming perfect Mendelian inheritance and accuracy, for MAS. This is due to the relatively low cost of maintaining strawberries and the high labour and reagent costs associated with performing genotyping (Edge-Garza, Luby, and Peace 2015; Wannemuehler et al. 2020). In the case of GS, a large

marker panel is required, of which most markers have small or zero effect; indeed, in the case of some models such as BayesB, a prior with some proportion of markers with zero effect is assumed. The cost of genotyping an individual for GS is thus higher than the cost of genotyping for tMAS. Attempts to rationally design a marker set multiplexing the large number of individuals in a breeding population to reduce costs is challenging due to poor predictive programs for multiplex primer design. Current published experimental validations of GP in strawberry utilise one of the available SNP arrays (Gezan et al. 2017; Osorio et al. 2021; Pincot et al. 2020).

In order for a commercial breeding program to benefit from GS, selections at the seedling stage are required. This requires almost all traits of agronomic importance to be amenable to GP as if some traits cannot be predicted on genomic data, maturation of the plant for assessment of those traits is needed. Depending on the breeding programme, over 40 traits may be of interest to breeders. Fortunately, there appears to be no experimental evidence of a trait that cannot be predicted using GP, though some traits have low accuracies of prediction. Note that MAS would not be suitable for a similar speed breeding scheme by genotyping and selection at the seedling stage as many traits do not have large effect markers that can be identified. Traits of interest in strawberry breeding can be split into fruit quality and plant habit traits. The GP model presented in this thesis encompasses almost all the traits of interest to strawberry breeders (Mathey et al. 2013), potentially allowing for selections to be made prior to fruit development. This would reduce the selection cycle by a few months if plant habit traits can be adequately assessed before fruit maturation.

If GP allows for increased genetic gain per unit time through reduction of the duration of the breeding cycle, then adoption of GP is a binary choice for strawberry breeders. Gradual introduction of GP for some traits cannot be performed (as can be envisaged for MAS) as shortening of breeding cycle cannot be achieved in these cases. Thus, experimental evidence for efficacy of GP in strawberry, as well as availability of datasets must be convincing before commercial adoption. Research focus should be targeted towards demonstration of speed breeding by ensuring that almost all traits of interest to breeders are experimentally modelled with GP.

Selection of models typically has a smaller effect on selection accuracy than number of markers or trait to be predicted, and is not expected to be the primary mechanism by which GP improves selection in strawberries. Moreover, if the duration of selection cycles are reduced to months (from planting of breeding population to selection), then the days or weeks of computation required for the more computationally intensive models of GP (such as BayesA) may become significant in speed breeding. Models for cost efficiency of MAP should

be extended to account for the potential to remove the growth time for the first stage selections under speed breeding and GP.

Deployment of GP requires most traits of strawberry to be predictable by GP. However, due to correlation between traits, selection on individual traits may be inefficient; selection indices allow for appropriate weighing of traits that may be correlated. Correlation between traits is primarily due to LD between controlling genes and pleiotropy (Pedruzzi and Rouzine 2019). Effort should be made to ensure multiple traits of agronomic importance are measured in the same strawberry population so these correlations can be quantified to inform construction of selection indices. Establishment of selection indices are dependent on the goals of the breeding program, but typically combination of traits of agronomic import are non-linear. Often there are thresholds that must be met for some characteristics to ensure that a potential novel cultivar exceeds performance of some check species, making unpredictable traits particularly disadvantageous for GP. If there is a trait that cannot be selected for using GP, or has a low accuracy compared to phenotypic assessment methods, loss of prediction for it must be compensated for by increases in genetic gain from speed breeding.

Knowledge and Technology Transfer

AHDB Student Industry Visit (July 2018)

Soft Fruit walk, Kent, UK (June 2018)

AHDB Studentship Conference, UK (November 2017) – Poster Presentation on Genotyping-in-Thousands as a cost-effective method of genotyping strawberry

NIAB Student Outreach Event, Histon, UK (November 2017) – Oral and poster presentation on 3D strawberry phenotyping

Current and future applications of phenotyping for plant breeding, Novi Sad, Serbia (September 2017) – Poster and oral presentation

Crops Group Student Symposium, Reading, UK (Nov 2017)

AHDB Studentship Conference, UK (November 2017) – Oral presentation (TBC)

NIAB Student Outreach Event, Histon, UK (November 2017) – Oral and poster presentation on 3D strawberry phenotyping
4th International Horticultural Conference, East Malling, UK (July 2017) – Oral Presentation on 3D imaging in strawberry; poster presentation on cost-effective genotyping for strawberry breeding

Plant and Animal Genome XXV, San Diego, USA (January 2017) – Received Travel Award from AHDB and GCRI to attend conference

Tuscon Plant Breeding Institute, Tuscon, USA (January 2017) - Received Travel Award from AHDB and GCRI to attend course

AHDB studentship Conference, Stratford, UK (November 2016) – Oral presentation on PhD overview

Soft Fruit Day, East Malling, UK (November 2016) – Poster presentation on PhD overview

Grand Challenges in Plant Pathology, Oxford, UK (September 2016)

Software Carpentry, Norwich, UK (June 2016)

References

Akond, Zobaer, Md. Jahangir Alam, Mohammad Nazmol Hasan, Md. Shalim Uddin, Munirul Alam, and Nurul Haque Mollah. 2019. “A Comparison on Some Interval Mapping Approaches for QTL Detection.” *Bioinformatics* 15(2):90–94.

Antanaviciute, Laima. 2016. “Genetic Mapping and Phenotyping Plant Characteristics , Fruit Quality and Disease Resistance Traits in Octoploid Strawberry (*Fragaria x Ananassa*).” University of Reading.

Bates, D., M. Mächler, B. Bolker, and S. Walker. 2015. “Fitting Linear Mixed-Effects Models Using Lme4.” *Journal of Statistical Software* 67(1):1–48.

Berners-Lee, M., C. Kennelly, R. Watson, and C. N. Hewitt. 2018. “Current Global Food Production Is Sufficient to Meet Human Nutritional Needs in 2050 Provided There Is Radical Societal Adaptation.” *Elementa Science of the Anthropocene* 6(52):1–14.

Broman, Karl W. 2003. “Mapping Quantitative Trait Loci in the Case of a Spike in the Phenotype Distribution.” *Genetics* 163(3):1169–75.

Churchill, G. A. and R. W. Doerge. 1994. “Empirical Threshold Values for Quantitative Trait Mapping.” *Genetics* 138(3):963–71.

Edge-Garza, Daniel A., James J. Luby, and Cameron Peace. 2015. “Decision Support for Cost-Efficient and Logistically Feasible Marker-Assisted Seedling Selection in Fruit Breeding.” *Molecular Breeding* 35(12):1–15.

Endelman, Jeffrey B. 2011. “Ridge Regression and Other Kernels for Genomic Selection with R Package RrBLUP.” *The Plant Genome* 4(3):250–55.

Endelman, Jeffrey B. and Jean Luc Jannink. 2012. “Shrinkage Estimation of the Realized Relationship Matrix.” *G3: Genes, Genomes, Genetics* 2(11):1405–13.

Foulley, J. L. 1993. “A Simple Argument Showing How to Derive Restricted Maximum Likelihood.” *Journal of Dairy Science* 76(8):2320–24.

- Gezan, Salvador A., Luis F. Osorio, Sujeet Verma, and Vance M. Whitaker. 2017. "An Experimental Validation of Genomic Selection in Octoploid Strawberry." *Horticulture Research* 4(October 2016):1–9.
- Giné Bordonaba, J. and L. A. Terry. 2010. "Manipulating the Taste-Related Composition of Strawberry Fruits (*Fragaria* × *Ananassa*) from Different Cultivars Using Deficit Irrigation." *Food Chemistry* 122(4):1020–26.
- Habier, D., R. L. Fernando, and J. C. M. Dekkers. 2007. "The Impact of Genetic Relationship Information on Genome-Assisted Breeding Values." *Genetics* 177(4):2389–97.
- Habier, David, Rohan L. Fernando, Kadir Kizilkaya, and Dorian J. Garrick. 2011. "Extension of the Bayesian Alphabet for Genomic Selection." *BMC Bioinformatics* 12.
- He, Dan, Subrata Saha, Richard Finkers, and Laxmi Parida. 2018. "Efficient Algorithms for Polyploid Haplotype Phasing." *BMC Genomics* 19.
- He, Sang, Albert Wilhelm Schulthess, Vilson Mirdita, Yusheng Zhao, Viktor Korzun, Reiner Bothe, Erhard Ebmeyer, Jochen C. Reif, and Yong Jiang. 2016. "Genomic Selection in a Commercial Winter Wheat Population." *Theoretical and Applied Genetics* 129(3):641–51.
- Henderson, C. R. 1975. "Best Linear Unbiased Estimation and Prediction under a Selection Model Published by : International Biometric Society Stable." *Biometrics* 31(2):423–47.
- Henderson, Charles R. 1984. "Applications of Linear Models in Animal Breeding Models." *University of Guelph* 384.
- Heslot, Nicolas, Hsiao Pei Yang, Mark E. Sorrells, and Jean Luc Jannink. 2012. "Genomic Selection in Plant Breeding: A Comparison of Models." *Crop Science* 52(1):146–60.
- Kruijer, Willem, Martin P. Boer, Marcos Malosetti, Pádraic J. Flood, Bas Engel, Rik Kooke, Joost J. B. Keurentjes, and Fred A. Van Eeuwijk. 2014. "Marker-Based Estimation of Heritability in Immortal Populations." *Genetics* 199(2):379–98.
- Lande, R. and R. Thompson. 1990. "Efficiency of Marker-Assisted Selection in the Improvement of Quantitative Traits." *Genetics* 124(3):743–56.
- Mathey, Megan M., Sonali Mookerjee, Kazim Gündüz, James F. Hancock, Amy F. Iezzoni, Lise L. Mahoney, Thomas M. Davis, Nahla V. Bassil, Kim E. Hummer, Philip J. Stewart, Vance M. Whitaker, Daniel J. Sargent, Beatrice Denoyes, Iraidia Amaya, Eric Van De Weg, and Chad E. Finn. 2013. "Large-Scale Standardized Phenotyping of Strawberry in RosBREED." *Journal of the American Pomological Society* 67(4):205–16.
- Mehraj, H. and A. F. M. Jamal Uddin. 2014. "Correlation Pathway for Phenotypic Variability

- Study in Strawberry.” *International Journal of Sustainable Agricultural Technology* 10(109):1815–1272.
- Meuwissen, T. H. E. , Hayes, B. J., & Goddard, M. E. 2001. “Prediction of Total Genetic Value Using Genome-Wide Dense Markers Maps.” *Genetics* 157:1819–29.
- Meyer, J. Patrick and Michael A. Seaman. 2008. “A Comparison of the Exact Kruskal-Wallis Distribution to Asymptotic Approximations for N.” in *American Educational Research Association*.
- Molenaar, Heike, Robert Boehm, and Hans Peter Piepho. 2018. “Phenotypic Selection in Ornamental Breeding: It’s Better to Have the BLUPs than to Have the BLUEs.” *Frontiers in Plant Science* 871(November):1–14.
- Ogutu, Joseph O., Torben Schulz-Streeck, and Hans-Peter Piepho. 2012. “Genomic Selection Using Regularized Linear Regression Models: Ridge Regression, Lasso, Elastic Net and Their Extensions.” *BMC Proceedings* 6(Suppl 2):S10.
- van Ooijen, J. W. 2009. “MapQTL 6.” *Genome* (April).
- Osorio, Luis F., Salvador A. Gezan, Sujeet Verma, and Vance M. Whitaker. 2021. “Independent Validation of Genomic Prediction in Strawberry Over Multiple Cycles.” *Frontiers in Genetics* 11(January):1–13.
- Pedruzzi, Gabriele and Igor M. Rouzine. 2019. “Epistasis Detectably Alters Correlations between Genomic Sites in a Narrow Parameter Window.” *PLoS ONE* 14(5):1–16.
- Piepho, H. P., J. Möhring, A. E. Melchinger, and A. Büchse. 2008. “BLUP for Phenotypic Selection in Plant Breeding and Variety Testing.” *Euphytica* 161(1–2):209–28.
- Pincot, Dominique D. A., Michael A. Hardigan, Glenn S. Cole, Randi A. Famula, Peter M. Henry, Thomas R. Gordon, and Steven J. Knapp. 2020. “Accuracy of Genomic Selection and Long-Term Genetic Gain for Resistance to Verticillium Wilt in Strawberry.” *Plant Genome* 13(3):1–19.
- Salentijn, Elma M. J., Asaph Aharoni, Jan G. Schaart, Marjan J. Boone, and Frans A. Krens. 2003. “Differential Gene Expression Analysis of Strawberry Cultivars That Differ in Fruit-Firmness.” *Physiologia Plantarum* 118(4):571–78.
- Signorell, Andri, Ken Aho, Andreas Alfons, Nanina Anderegg, Tomas Aragon, Chandima Arachchige, Antti Arppe, Adrian Baddeley, Kamil Barton, Ben Bolker, Hans W. Borchers, Frederico Caeiro, Stephane Champely, Daniel Chessel, Leanne Chhay, Nicholas Cooper, Clint Cummins, Michael Dewey, Harold C. Doran, Stephane Dray, Charles Dupont, Dirk Edelbuettel, Claus Ekstrom, Martin Elff, Jeff Enos, Richard W.

Farebrother, John Fox, Romain Francois, Michael Friendly, Tal Galili, Matthias Gamer, Joseph L. Gastwirth, Vilmantas Gegzna, Yulia R. Gel, Sereina Graber, Juergen Gross, Gabor Grothendieck, Frank E. Harrell, Richard Heiberger, Michael Hoehle, Christian W. Hoffmann, Soeren Hojsgaard, Torsten Hothorn, Markus Huerzeler, Wallace W. Hui, Pete Hurd, Rob J. Hyndman, Christopher Jackson, Matthias Kohl, Mikko Korpela, Max Kuhn, Detlew Labes, Friederich Leisch, Jim Lemon, Dong Li, Martin Maechler, Arni Magnusson, Ben Mainwaring, Daniel Malter, George Marsaglia, John Marsaglia, Alina Matei, David Meyer, Weiwen Miao, Giovanni Millo, Yongyi Min, David Mitchell, Franziska Mueller, Markus Naepflin, Daniel Navarro, Henric Nilsson, Klaus Nordhausen, Derek Ogle, Hong Ooi, Nick Parsons, Sandrine Pavoine, Tony Plate, Luke Prendergast, Roland Rapold, William Revelle, Tyler Rinker, Brian D. Ripley, Caroline Rodriguez, Nathan Russell, Nick Sabbe, Ralph Scherer, Venkatraman E. Seshan, Michael Smithson, Greg Snow, Karline Soetaert, Werner A. Stahel, Alec Stephenson, Mark Stevenson, Ralf Stubner, Matthias Templ, Duncan Temple Lang, Terry Therneau, Yves Tille, Luis Torgo, Adrian Trapletti, Joshua Ulrich, Kevin Ushey, Jeremy VanDerWal, Bill Venables, John Verzani, Pablo J. Villacorta Iglesias, Gregory R. Warnes, Stefan Wellek, Hadley Wickham, Rand R. Wilcox, Peter Wolf, Daniel Wollschlaeger, Joseph Wood, Ying Wu, Thomas Yee, and Achim Zeileis. 2021. "Package 'DescTools.'"

Sood, Salej, Vinay Bhardwaj, S. K. Kaushik, and Sanjeev Sharma. 2020. "Prediction Based on Estimated Breeding Values Using Genealogy for Tuber Yield and Late Blight Resistance in Auto-Tetraploid Potato (*Solanum Tuberosum* L.)." *Heliyon* 6(11):e05624.

Viana, José Marcelo Soriano, Vinícius Ribeiro Faria, Fabyano Fonseca e Silva, and Marcos Deon Vilela de Resende. 2011. "Best Linear Unbiased Prediction and Family Selection in Crop Species." *Crop Science* 51(6):2371–81.

Wannemuehler, Seth D., Chengyan Yue, Wendy K. Hoashi-Erhardt, R. Karina Gallardo, Vicki McCracken, and R. Karina Gallardo. 2020. "Cost-Effectiveness Analysis of a Strawberry Breeding Program Incorporating DNA-Informed Technology." *HortTechnology* 30(3):365–71.

Xu, Shizhong. 2003. "Theoretical Basis of the Beavis Effect." *Genetics* 165(4):2259–68.

Appendices

Code and data for analysis can be found at the NIAB EMR github repository ([www.github.com/ organizations/eastmallingresearch/](http://www.github.com/organizations/eastmallingresearch/)).

