

June 2019



## **Project Report No. PR610**

# **MAGIC map and go: deploying MAGIC populations for rapid development and dissemination of genetic markers for yield improvement in elite UK winter wheat**

Keith Gardner<sup>1</sup>, Benedetta Sacommano<sup>1</sup>, Phil Howell<sup>1</sup>, Ian Mackay<sup>1</sup>, Anna Sanchez<sup>1</sup>, John Jacobs<sup>2</sup>, Steve Smith<sup>3</sup>, Charlotte Hayes<sup>3</sup>, Nicholas Bird<sup>4</sup>, Ed Byrne<sup>4</sup>, Jacob Lage<sup>4</sup>, Simon Berry<sup>5</sup>, Edward Flatman<sup>5</sup>, Peter Jack<sup>6</sup>, Chris Burt<sup>6</sup>, and James Cockram<sup>1</sup>

<sup>1</sup>NIAB, Huntingdon Road, Cambridge, CB3 0LE,

<sup>2</sup>BASF, Technologiepark-Zwijnaarde 38, 9052 Gent

<sup>3</sup>Elsoms Seeds Ltd, Pinchbeck Road, Spalding, Pe11 1QG

<sup>4</sup>KWS UK Ltd, 56 Church Street, Thriplow, SG8 7RE

<sup>5</sup>Limagrain UK Ltd, Market Raisen, LN7 6DT

<sup>6</sup>RAGT Seeds Ltd, Grange Road, Saffron Walden, CB10 1TA

This is the final report of a 48-month project (21130017) that started in January 2015. The work was funded by BBSRC (BB/M008908/1) and a contract for £99,544 from AHDB.

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.

Reference herein to trade names and proprietary products without stating that they are protected does not imply that they may be regarded as unprotected and thus free for general use. No endorsement of named products is intended, nor is any criticism implied of other alternative, but unnamed, products.

AHDB Cereals & Oilseeds is a part of the Agriculture and Horticulture Development Board (AHDB).



# CONTENTS

1.	ABSTRACT .....	1
2.	INTRODUCTION .....	2
3.	MATERIALS AND METHODS .....	6
	3.1 NIAB Elite MAGIC germplasm.....	6
	3.2 90k SNP genotypic data .....	6
	3.3 Field trials and phenotyping .....	6
	3.3.1 Trial design.....	6
	3.3.2 Phenotyping .....	7
	3.4 Statistical analyses .....	12
	3.4.1 Trials analysis .....	12
	3.4.2 Genetic analysis.....	12
	3.4.3 Genomic prediction .....	13
	3.5 KASP marker development .....	14
	3.6 Bioinformatic analysis.....	14
	3.7 Identification of candidate genes and TILLING lines.....	15
	3.8 NIL development .....	15
4.	RESULTS .....	16
	4.1 Field trials, phenotyping and trials analysis .....	16
	4.1.1 Trials and phenotyping.....	16
	4.1.2 Transgressive segregation.....	16
	4.1.3 Trait heritabilities .....	18
	4.2 QTL analysis .....	18
	4.3 Development of molecular markers tagging QTL.....	24
	4.4 Development of near isogenic lines for selected QTL .....	25
	4.5 Identification of candidate genes and TILLING.....	29
	4.6 Genomic prediction .....	35
5.	DISCUSSION .....	37
	5.1 Overview.....	37
	5.2 Genetic control of agronomic traits.....	37

5.3 Candidate genes and TILLING .....	39
5.4 Suggestions for future research .....	39
5.4.1 Leveraging emerging MAGIC genome re-sequencing and RNAseq data to prioritise candidate polymorphisms within QTL intervals .....	40
5.4.2 Detailed understanding of Mendalised QTL.....	41
5.4.3 Development of hybrid wheat approaches and resources .....	41
6. REFERENCES .....	43
7. WORKPLAN.....	47
8. DISEMINATION AND PROJECT OUTPUTS TO DATE .....	48
9. ACKNOWLEDGEMENTS.....	49

## Glossary

CDS	Coding sequence or regions
GP	Genomic prediction
KASP	Kompetitive allele specific PCR
MAGIC	Multi-founder advanced generation inter-cross
NIL	Near isogenic line
PCR	Polymerase chain reaction
QTL	Quantitative trait locus
RIL	Recombinant inbred line
SSD	Single seed descent
SNP	Single nucleotide polymorphism
TILLING	Targeting induced local lesions in genomes

## Summary of appendixes

Appendixes are not published within this report. They have been embargoed until 29 June 2023 (unless permission is gained from the industrial partners to publish information beforehand).

Appendix	Appendix title
Appendix 1	KASP marker primer details
Appendix 2	Phenotypic data: predicted means and details of transgressive segregation for 2015 and 2016 season trials
Appendix 3	Summary of QTL results
Appendix 4	QTLs investigated for KASP marker development, gene content and candidate gene analysis, TILLING and NIL development
Appendix 5	MAGIC RIL selections for NIL development
Appendix 6	Wheat homologues of rice genes known to control grain size characters
Appendix 7	Candidate gene TILLING information
Appendix 8	Genomic prediction summary data
Appendix 9	Analysis of the major spikelet number QTL on chromosome 7A

## 1. Abstract

The genetic improvement of grain yield in bread wheat was targeted by this project. In collaboration with five wheat breeding companies, the high-density genotyping of a wheat multi-founder advanced generation inter-cross (MAGIC) population was exploited. This population captures high levels of genetic recombination and diversity. It was used to: **1.** Identify the genetic regions in wheat controlling yield and stability. **2.** Provide a molecular toolkit to track, within breeding programmes, regions of the wheat genome that confer increased yield/yield stability. **3.** Provide the participating breeders with analysis pipelines and resources to carry out analysis of MAGIC datasets. **4.** Exploit the structure of the MAGIC population to rapidly 'Mendelize' QTL for multiple yield/yield component traits, providing precise genetic and molecular resources for subsequent studies to fine-map to the gene/causative polymorphism level. **5.** Use the molecular breeding methodology, genomic prediction (GP), to allow selection for yield/yield stability in MAGIC lines, based on molecular data alone.

MAGIC lines (1109) were grown at five UK sites over two seasons. This delivered 4,996 2x6m plots on which 18 yield, yield component and agronomic traits were measured. This generated ~90,000 phenotypic data points. Phenotypic information was combined with genotypic data for ~20,000 SNPs for genetic analysis. This identified 376 QTL, with genetic intervals for most QTL <10 cM. A subset of 20 QTLs was prioritised for the development of 'KASP' genetic markers, based on QTL significance, allelic effect and stability across years and sites. This resulted in the design and validation of 58 co-dominant KASP markers for the 20 target QTL. We developed/initiated 31 near isogenic lines (NILs) for 17 traits, along with genetic markers with which to further exploit these materials. Additionally, we determined the gene content from the variety Chinese Spring and identified candidate genes. For one QTL, we identified a single candidate gene controlling spikelet number per ear on chromosome 7A. We searched for artificial mutants using wheat 'TILLING' populations for the tetraploid and hexaploid varieties Kronos and Cadenza, respectively. Highly deleterious mutations across multiple homoeologues in 23 candidate genes were identified. These TILLING resources will be used to help determine the specific genes and genetic variants underlying the QTLs identified. Finally, we used pre-project MAGIC phenotypic data for yield, in combination with phenotypic data collected in this project, to investigate GP for predicting phenotypic performance in a given generation, based on markers for plant height.

MAGIC reliably delineates QTL to relatively precise genetic and physical intervals. Tightly linked co-dominant genetic markers, that tag yield and yield component QTL, have been delivered to our industry partners for potential use in their breeding programmes. The work allowed candidate genes within QTL to be identified, and extensive biological (NILs, TILLING lines) and molecular (KASP markers) materials to be generated with which to further investigate phenotypic effects of QTL and genes in isolation.

## 2. Introduction

To date, bi-parental populations have been the norm for detecting quantitative trait loci (QTL) in wheat, and these commonly use low numbers of progeny (typically 100-200). This approach often suffers from a lack of statistical power and precision, due to the limited genetic variation and recombination captured. Previously, low molecular marker densities have historically meant that such problems were commonly overlooked. Such limitations have contributed to a lack of translation into breeding programmes. It has been estimated that globally, 10,000 QTL have been published (Bernardo, 2008), of which the vast majority have never been employed in marker-assisted breeding. We believe there are several reasons for this:

- (1) A focus on traits which are of minor breeder interest. Here, we focus on yield and yield stability, core breeding targets as supported by the participation of six industrial partners.
- (2) Mapping in populations created by crossing extremes for the target traits. It is easy to map QTL in a cross between the best and worst lines available, but of more immediate relevance in breeding is tagging QTL in crosses among the best lines. The use of breeder-selected founder lines in our MAGIC population has created a resource which is diverse and highly recombined, yet of immediate relevance.
- (3) Lack of precision. Bi-parental populations, deployed using population sizes >200 as is the common case, are relatively imprecise in their ability to locate QTL. Following large chromosomal tracts in breeding programmes to tag a single QTL is not good practice. The association between the flanking markers and QTL is very readily lost through recombination and large tracts are also more likely to contain alleles at QTL for detrimental traits. MAGIC locates QTL with much higher precision, facilitating rapid uptake in breeding programmes.

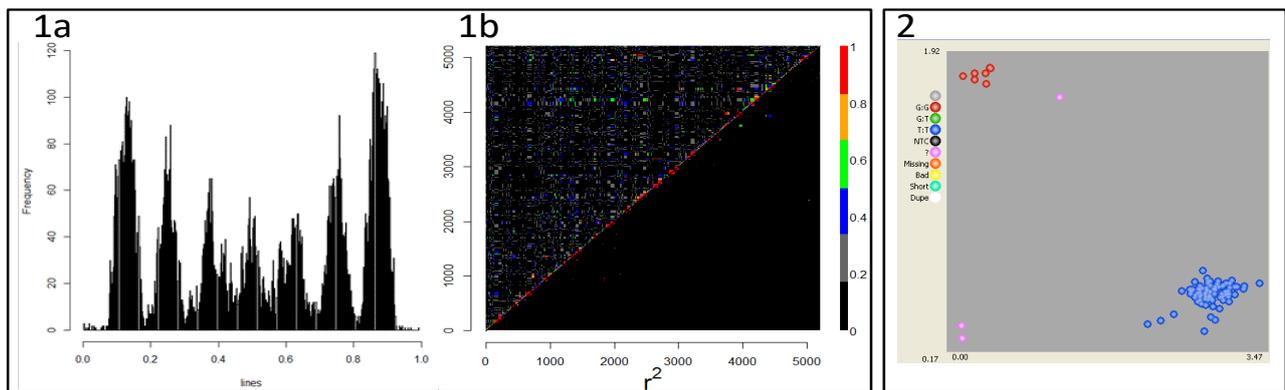
The recent development of high-density single nucleotide polymorphism (SNP) arrays has meant the availability of molecular markers is no longer the limiting factor for wheat QTL studies. However, lack of similar improvements in mapping population design has meant that much of the potential for greater mapping resolution afforded by increased marker density has remained largely unrealised. BBSRC's recent investment in developing the NIAB wheat MAGIC population constructed from eight elite wheat founders ('NIAB Elite MAGIC', **Table 1**) has helped bridge this gap, providing a biological resource ideally suited to leverage advances in wheat genomics.

Variety	Reason for choice	RL	Breeder
Alchemy	Yield, disease resistance, breeding use, soft group 4	2013/14	Limagrain
Brompton	Hard feed, 1BL/1RS, OWBM resistant, group 4	2008/09	Elsoms
Claire	Soft biscuit/distilling, slow apical development, group 3	2013/14	Limagrain
Hereward	High quality benchmark 1 breadmaking, group 1	2010/11	RAGT
Rialto	Moderate breadmaking, 1BL/1RS, group 2	2002/03	RAGT
Robigus	High yielding soft biscuit/distilling, OWBM resistant, <i>Rht1</i>	2009/10	KWS
Soissons	Very early flowering, high quality group 2 breadmaking	2009/10	Depr.
Xi19	Facultative, high quality group 1 breadmaking, group 1	2010/11	Nickersons

**Table 1.** Parents of the UK Elite MAGIC population. RL = most recent appearance in AHDB Recommended List®.

**MAGIC:** By combining controlled allelic inputs from multiple parents with high levels of recombination achieved over multiple rounds of inter-crossing, MAGIC populations overcome the specific drawbacks associated with traditional bi-parental mapping populations, as well as more recent Association Mapping approaches (Mackay & Powell, 2007; Cavanagh et al. 2007). The higher allelic diversity in MAGIC populations improves sampling of available genetic diversity and phenotypic variation, facilitating the analysis of interacting or complex traits within a single mapping population. Combined with the suitability for the generation of high-density genetic maps, these factors make MAGIC populations ideal platforms for high-resolution genetic dissection of QTL and as community-based resources for crop improvement. The MAGIC approach is a translation to crops of methods employed in mouse genetics (Mott et al. 2000). The approach (and name) was advocated by NIAB in 2007 (Mackay & Powell, 2007; Cavanagh et al. 2007). It has caught on: Cavanagh et al. (2007) has been cited 292 times and Mackay & Powell (2007) 388 times. An Australian four-parent spring wheat MAGIC population has recently been developed, results from which are now emerging (e.g. Huang et al. 2012). However, our population remains the only MAGIC resource of direct relevance to UK winter wheat. The eight parental lines represent commercially available UK winter-sown varieties, selected for use in this population in consultation with UK wheat breeders (**Table 1**). The resulting highly recombined population of >1,000 progeny provides a high-resolution platform for the genetic dissection of trait inheritance. NIAB funding has provided 90k Illumina iSelect SNP array datasets for 643 of the >1,000 MAGIC lines (**Figure 1a**). Our analyses find the population to be well suited as a platform for fine-mapping QTL and gene isolation: patterns of linkage disequilibrium (LD) show the population to be highly recombined (**Figure 1b**; Mackay et al. 2014), while comparison with a north-western European association mapping panel of 480 varieties finds the MAGIC population captures ~80% of the available SNP variation (Gardner et al. 2016). To demonstrate the potential of the resource, we used these MAGIC datasets to investigate the genetic control of awning, finding a highly diagnostic marker for awn presence absence, which was converted to the flexible KASP genotyping system (**Figure 2**). An advantage of the NIAB Elite MAGIC population and associated SNP genotypic data, is the ability to exploit the residual heterozygosity present in the genotyped lines (predicted to be at ~2% of the genome) in order to develop near isogenic lines (NILs). These

germplasm resources can be rapidly developed by identifying MAGIC progeny that are heterozygous across the target QTL region. By growing and selfing such lines, progeny can be identified using genetic markers that are homozygous for each of the two allele classes (i.e. A:A individuals and B:B individuals). Such lines will differ at the QTL locus, but will be almost identical (~98% on average) throughout the rest of the genome. Such NIL resources, in which individual QTL have been rapidly Mendelised, represent key resources with which to undertake detailed analysis of the phenotypic effects of single QTL in isolation. The NIL pair can also be crossed together to generate large numbers of F2 progeny, allowing fine-mapping of Mendelised QTL via identification of genetic recombinations within the target genetic interval.



**Figure 1.** Elite MAGIC population properties. (a) Allele frequencies in MAGIC progeny genotyped with 90k SNP array. Peak increments are as expected from an 8-parent crossing scheme. (b) Analysis of genome-wide linkage disequilibrium ( $D'$  above horizontal,  $R^2$  below). **Figure 2.** Diagnostic SNP for awning identified using 90k array, converted to the KASP genotyping platform.

**Yield:** Grain yield is the most important trait in wheat, and is under complex genetic control. Yield can be largely partitioned into three major components: kernel weight/size and shape, kernel number per spike and spike number per unit area. These phenotypes can be further subdivided (e.g. grain size can be subdivided into grain length, width, depth, shape and density, and is thought to be predominantly determined by the genetic control of cell proliferation and expansion) and are influenced by other traits, such as accumulation and transport of photosynthetic products, flag leaf size, plant height, biomass and flowering time. Furthermore, wheat processing requirements mean traits that affect milling performance (e.g. grain shape, size, density and uniformity) are also critical for flour yield. Numerous studies investigating the genetics of yield (e.g. Rustgi et al. 2013), grain size (Gegas et al. 2010), flowering time (Griffiths et al. 2009a), height (Griffiths et al. 2009b), flag leaf size (Xue et al. 2013) and other yield components have been undertaken. However, many have either suffered from lack of QTL precision, are based on crosses between extremes of phenotype (i.e. not the best crossed with the best), or are not relevant to UK germplasm. Furthermore, these studies have been conducted within individual bi-parental populations, which restricts the genetic variation (and therefore, the number of yield-related traits) studied in any one cross. This limits the

scope for detection of QTL-by-QTL and QTL-by-environment interactions. The advent of high-throughput image-based phenotypic data collection of yield components provides the opportunity to revolutionise their genetic dissection. However, the gains obtained by the deployment of such technology will always rely on the germplasm investigated. In this project, we combine high-throughput phenotyping and high-density genotyping with germplasm resources ideally suited to realise the advances in these technologies. The combination of complementary resources used in this project provide the experimental power, precision and phenotypic depth to undertake precise genetic dissection of all target traits within a single platform, allowing analysis of the genetic and environmental interactions between multiple yield-related traits measured over six growing seasons.

**Objectives:** This project exploits NIAB's expertise in crop genetics and quantitative analysis, breeder expertise in growing and phenotyping wheat yield plots, and NIAB's MAGIC wheat biological resource, to:

1. Deliver genetic markers for the genetic improvement of yield, yield components and yield stability in wheat.
2. Develop NILs for major yield and yield component QTLs, as well as the associated genetic markers with which to validate them.
3. Provide genotypically and phenotypically characterised MAGIC germplasm for selection for lines for possible inclusion within breeding programmes.
4. Develop a Genomic Prediction strategy for selection for yield and yield stability on markers alone.

### 3. Materials and methods

#### 3.1 NIAB Elite MAGIC germplasm

The NIAB Elite MAGIC population has been previously reported by Mackay et al. (2014). Briefly, the eight wheat founder varieties (**Table 1**) were intercrossed over three generations, using a design based on a simple replicated funnel crossing scheme of the form  $\{[(A \times B) \times (C \times D)] \times [(E \times F) \times (G \times H)]\}$ , where the matched brackets (), [] and {} delineate the four 2-way, two 4-way and one 8-way cross, respectively, and the letters denote the eight founders. All 28 possible 2-way crosses were performed (AxB, AxC, AxD etc). Similarly, all 210 possible 4-way crosses between unrelated 2-way lines were performed. At the 8-way stage, only 210 of the possible 315 crosses were made, using each 4-way twice. The outputs of the 8-way cross, which possess contributions from all eight founders, are termed F1s. These F1s were selfed through multiple rounds of single seed descent (SSD) to generate a target population of 1,000 recombinant inbred lines (RILs). The only purposeful selection imposed on the population was removal of lines with extreme short stature, due to the presence of dwarfing alleles at both the *REDUCED HEIGHT-B1* (*RHT-B1*) and *RHT-D1* semi-dwarfing loci, which segregated in the population.

#### 3.2 90k SNP genotypic data

Previously generated SNP genotypic data for 643 MAGIC lines and the 8 founders is as described by Mackay et al. (2014), and further quality controlled as described by Gardner et al. (2016). Briefly, DNA was extracted from a single F<sub>5</sub> RIL individual per line using a modified Tanksley method (Fulton et al. 1995), and genotyped using the Illumina Infinium wheat 90k SNP array (Wang et al. 2014). Genotype calls were processed using Genome Studio v2011.1 (Illumina, San Diego, USA) and quality controlled following the pipeline outlined in Appendix S1 from Gardner et al. (2016). This generated 20,639 polymorphic SNPs, 18,601 of which were successfully anchored on the NIAB Elite MAGIC genetic map (Gardner et al. 2016).

#### 3.3 Field trials and phenotyping

##### 3.3.1 Trial design

Autumn sown field trials using the NIAB Elite MAGIC population (1109 lines, 8 founders) and one commercial wheat control (cv. KWS Santiago) were undertaken over two wheat seasons (2014-2015 and 2015-2016), with five sites per season (as listed in **Table 2**). For any one year, trials were designed so that approximately 80% (2015) or 65% (2016) of lines were grown in a single rep in each of two trials and 20% (2015) or 35% (2016) of lines were grown in 2 reps in one trial and one rep in another trial. Exact details per trial are shown in **Table 2**. Within each trial, trial design was undertaken using the “Design of experiments” website (“DEW”, now obsolete but replaced by the R package “blocksdesign”). All sites except LIM 2015 had a 3 tier, Main/sub-

Trial site – harvest year	Location: lat, long (degrees)	Date sown	No. of plots (runs x rows)	Trial design (b, s1, s2) <sup>†</sup>	No. MAGIC RILs	RIL replication 1:2 reps	Reps per founder; reps per Santiago
BAY-YLD-15	N/A	N/A	504 (8x63)	2,7,4	420	376:44	4;8
BAY-YLD-16	N/A	18/10/2015	504 (14x36)	2,7,4	384	300:84	4;4
ELS-YLD-15	52.82342, -0.12293	17/10/2014	500 (10x50)	4,5,5	421	378:43	4;4
ELS-YLD-16	52.83472, -0.10886	19/10/2015	500 (10x50)	2,5,5	390	316:74	4;4
KWS-YLD-15	52.12131, 0.08278	22/10/2015	504 (36x14)	2,7,4	421	374:47	4;5
KWS-YLD-16	N/A	N/A	504 (12x42)	2,6,6	393	318:75	4;4
LIM-YLD-15	52.20274, 0.85011	24/10/2014	500 (25x20)	25 plots	422	380:42	4;4
LIM-YLD-16	52.20402, 0.88312	13/10/2015	504 (12x42)	2,6,3	395	322:73	4;4
RAG-YLD-15	52.14363, -0.16692	22/10/2014	490 (14x35)	2,7,5	420	386:34	4;4
RAG-YLD-16	N/A	14/10/2015	490 (14x35)	2,7,5	394	334:60	4;4

**Table 2.** Details of field trials, including location, trial design and traits phenotyped. The Trial name is in the format XXX-YYY-ZZ where XXX = trial site (BAY = Bayer trials, ELS = Elsoms trials, KWS = KWS trials, LIM = Limagrain trials, RAG = RAGT trials), YYY indicated the trial type (YLD = yield trial) and ZZ = harvest year (15 = 2015, 16 = 2016). N/A = data not available. <sup>†</sup>b = blocks, s1 = sub-block1, s2 = sub-block2.

block1/sub-block2 design but with variation in the details based on individual site characteristics, with a range from 5-9 plots per lowest tier. For the 2015 and 2016 season trials, F8 and F9 NIAB Elite MAGIC seed was sourced from nursery plots grown at NIAB-Cambridge in the preceding season of each trial, respectively. All trials were run using standard agronomic packages (fertilisers, pesticides and growth regulators) for the locations in which they were grown.

### 3.3.2 Phenotyping

Trials were phenotyped for a suite of 18 traits. These are divided into two categories: pre-harvest (8 traits) and post-harvest (10 traits), and are listed along with a summary of their scoring methodologies, in **Table 3**. All post-harvest traits were measured using a Marvin Grain Analyser (GTA Sensorik GmbH, Germany) and an associated balance. MARVIN analysis was undertaken as detailed in **Box1**.

Trait abbreviation	Trait description
AWN	Awn presence/absence
GS55	Growth stage 55 (days from 1 <sup>st</sup> May), half of ear emerged above flag leaf
LOD <sup>†</sup>	Lodging (% of plants leaning >45°)
PHT	Plant height (cm), at maturity, excluding awns/scurs
SKT	Spikelet number, mean of 10 ears
SPW	Specific weight (g), off combine
TILL	Tiller number, average of 3 30x30cm quadrats
YLD	Grain yield (t/ha), calculated from harvest fresh-weight and moisture %
ARE	Area of seed (mm <sup>2</sup> ), mean of seed from 30 ears
CIR	Circumference of seed (arbitrary units), mean of seed from 30 ears
FFD	Factor form density (arbitrary unity), TGW/ARE, mean of seed from 30 ears
LEN	Length of seed (cm), mean of seed from 30 ears
LWR	Seed length to width ratio, mean of seed from 30 ears
PSH	Percentage of shrivelled seeds (%), mean of seed from 30 ears
SNO	Seed number, mean of seed from 30 ears
TGW	Thousand grain weight (g), calculated using seed from 30 ears
VWT	Volume weight (g), weight of 25 ml seeds from 30 ears
WID	Seed width (mm), mean of seed from 30 ears

**Table 3.** Traits abbreviations and descriptions. Pre-harvest (top) and post-harvest (bottom) In some trials, the lodging trait was further subdivided into lodging (% of plants leaning >45°), leaning (LEA, (% of plants leaning <45°) and lodging+leaning (LLE, (% of plants leaning >45° + leaning <45°).

## **Box1. MARVIN analysis protocol.**

Threshed and cleaned grain from 30 ears per plot were used for MARVIN analysis. For each plot, grain samples were analysed by MARVIN in two batches. 2015 season analysis was undertaken using MARVIN software version 4.0. For the 2016 season analysis, version 5.0 was used, which differs in that data output is provided for each individual grain, rather than a mean value for the imaged sample.

### A. 2015 season

MARVIN data collection was undertaken by following the standard parameters for wheat.

### B. 2016 season

To control for any potential temporal drift in MARVIN data for ARE, LEN, WID and LWR over the timescale of data collection, revision of the MARVIN protocol for the 2016 season grain included the use of two controls, phenotyped at approximately 2 hr intervals (every 20 samples) throughout each day of analysing: (1) clay seeds, and (2) real seeds.

Two types of output were generated from MARVIN: (a) summary output data for all seeds in a sample, and (b) data for individual seeds within a sample. Sample summary data was used for initial QC, highlighting obvious outliers for MARVIN re-analysis. Subsequently, deviation of MAGIC sample data from the reference sample data phenotyped on the same day was calculated, by fitting a smoothing spline to the data in using a custom R script (NIAB reference: Control\_splines\_script.R). The finalised summary table and the individual seed data table were then uploaded to R and processed using a custom script (NIAB reference: Marvin\_Processing\_script.R). Critical steps are described below.

### 2. Seed classification

Using the individual seed data from (1) above, individual grains were classified into eight classes:

DOUBLES = 2 touching grains imaged as a single grain

MULTIS = more than 2 grains imaged as a single grain

JUNK = chaff, awns, dust, small fragments not automatically removed by Marvin protocol

CHAFF\_ATTACHED = grain with chaff attached

BROKEN = broken grain

SHRIVELLED = shrivelled grain

ANGLED\_SEEDS = grains not sitting flat on scanner bed

NORMAL\_SEEDS = all other imaged grains

Imaged seeds were allocated to these classes based on the thresholds and ordered process listed below:

Seed number >1 = MULTIS

Length <3.6mm = JUNK

Width <1.5mm = JUNK

L/w >3.4 = JUNK

Area <7.9mm<sup>2</sup> = JUNK

Width >5.1mm = DOUBLES

Length >8.6mm AND area >30mm<sup>2</sup> = DOUBLES

Length >8.6mm AND area <12mm<sup>2</sup> = JUNK

Length >8.6mm AND 12<area>30mm<sup>2</sup> = CHAFF\_ATTACHED

L/W <1.4 = BROKEN

L/W <1.5 AND Length <5.1mm = BROKEN

CIR ≥2 AND area <12mm<sup>2</sup> = JUNK

Area >31mm<sup>2</sup> = DOUBLES

Width <2.8mm = SHRIVELLED

L/W >2.3 = SHRIVELLED

Width >4.7mm = ANGLED\_SEEDS

Remainder = NORMAL\_SEEDS

### 3. Downstream Analysis

Traits included in downstream analysis fall into four categories

Group 1. Seed size/shape traits: LEN (average length), WID (average width), ARE (average area), LWR (average length/width), CIR (average CIR).

Group 2. Seed number (SNO) based on all seeds including shrivelled, thousand grain weight (TGW) based on all seeds including shrivelled, factor form density (FFD) based on all seeds including shrivelled.

Group 3. SNO based on TGW with non-shrivelled seed only, FFD based on non-shrivelled seed only.

Group 4. Percent shrivelled seed (PSH).

These categories were used in downstream analysis as follows: MULTIS – remove measures from Group 1 Traits, but use actual seed number in Group 2 and Group 3 traits, assuming MULTIS are non-shrivelled. DOUBLES – treat as MULTIS after correcting seed number to two. JUNK – remove from all trait categories. SHRIVELLED – remove from Group 1 traits, include in Group 2 traits, remove from seed number in Group 3 traits TGW and FFD. Use for Gp4 estimation. CHAFF\_ATTACHED – remove from Group 1 traits, use actual seed number in Group 2 and Group 3 traits. BROKEN – omit from Group 1. For Group 2 and Group 3, use following calculation: Number

of broken seeds = sum of area of broken seeds/average seed area of normal seed.  
 ANGLED\_SEEDS – drop from Group1 analysis, keep in Group 2, Group3 and Group 4 (as non-shrivelled). NORMAL – use in all analyses. These analyses result in the final list of MARVIN traits shown in **Table 3**:

*Size traits.*

ARE_NOR_MEAN	Average area of normal seeds
ARE_NORSHR_MEAN	Average area of normal plus shrivelled seed
LEN_NOR_MEAN	Average length of normal seeds
LEN_NORSHR_MEAN	Average length of normal plus shrivelled seed
WID_NOR_MEAN	Average width of normal seeds
WID_NORSHR_MEAN	Average width of normal plus shrivelled seed
LWR_NOR_MEAN	Average length-width ratio of normal seeds
LWR_NORSHR_MEAN	Average length-width ratio of normal plus shrivelled seed
CIR_NOR_MEAN	Average circularity of normal seeds
CIR_NORSHR_MEAN	Average CIRTU circularity of normal plus shrivelled seed

*Seed Number Traits*

SNO_NOR	Number of seeds in all categories EXCEPT shrivelled
SNO_TOT	Total seed number including shrivelled
PSH	Percentage Shrivelled Seed

*Weight traits*

TGW_NOR	Thousand grain weight of all non-shrivelled seed
TGW_TOT	Thousand grain weight of all seed
VWT1	Weight of 25ml seeds averaged over all sub-samples weighed
VWT2	Weight of 25ml seeds averaged over first 2 sub-samples weighed

*Derived traits*

FFD_NOR	TGW_NOR/ARE_NOR
FFD_TOT	TGW_TOT/ARE_NOR_SHR

### **3.4 Statistical analyses**

#### **3.4.1 Trials analysis**

All data was subject to rigorous quality control, including judicial removal of outliers where the weight of evidence suggested they were erroneous. All trials were analysed using linear mixed models in Genstat 18.0 (VSN International) and spatial autocorrelation and blocking models were thoroughly explored, including all blocking components, different spatial models and relevant covariates. Model selection was "informed AIC", i.e. based on Akaike information content (AIC) but if two models differed very slightly in AIC, the more biologically informative model was chosen, even if it had a fractionally larger AIC. In a few cases of unstable AIC values between similar models, the parameter causing the instability was dropped, even if the models including this parameter apparently had a lower AIC. Variograms and residual analysis were used in some cases to reject poor models.

#### **3.4.2 Genetic analyses**

MAGIC genetic analyses were undertaken using two approaches.

- Single marker analysis base on identity by descent (IBS): a simple linear model test in R/lme4 using all 20,643 SNPs. After identification of a major QTL, the analysis was repeated with the major QTL as a covariate.
- Haplotype analysis using a subset of 7,369 uniquely mapped SNPs from the MAGIC genetic map (Gardner et al. 2016).

The haplotype approaches are more likely to (a) detect QTL and (b) accurately locate QTL. However, in poorly mapped areas of the genome (e.g. around introgressions) these approaches will fail. In these cases, if markers in these regions are synchronised with the QTL, the IBS approach will be the only one able to detect QTL.

Founder haplotype probabilities were computed with the "mpprob" function in R/mpMap (Huang & George, 2011), implemented in R/qtl (Broman et al. 2003). QTL analysis with haplotypes was carried out (a) via linear mixed model using all mapped markers, called identity by descent (IBD), (b) by simple interval mapping (IM) using the mpIM function in R/mpMap and (c) by composite interval mapping (CIM) using the mpIM function in R/mpMap with either 5 or 10 covariates. A full QTL model was then fitted with all QTL using R/fit.mpQTL.

The MAGIC genetic map of Gardner et al. (2016) was used for QTL mapping. Co-ordinates for the genetic map markers on the IWGSC RefSeq v1.0 wheat reference genome were obtained by a reciprocal alignment of the 90K array marker sequence data with the reference genome, followed by a filtering step: physical map positions were added to the map only if they were <50MB from the average of the surrounding 20 adjacent markers on the genetic map, with the exception of the centromeric regions (and regions with known inter-specific introgressions), where a more subjective, chromosome specific approach was taken (Gardner et al. 2016).

For IBS and IBD analyses, multiple-test correction was carried out using R/qvalue, with a threshold of  $q < 0.05$ . QTL distinctness (i.e. where does one QTL end and the next begin) was assessed manually for the IBS/IBD analyses, based on observation of QTL continuity, map quality (based on agreement between the genetic map and IWGSC RefSeq v1.0 co-ordinates), and recombination rate in the genomic region under investigation. Once QTL were separated out, the peak marker (minimum  $q$  value) was selected as a reference marker for that QTL. For IM/CIM analyses, an initial liberal cut-off of  $-\log_{10}p < 3$  and a window size of 100 markers was used in "find.qtl". "Fit.qtl" was then applied, and QTL retained which had  $p < 0.05$  in the fitted model, as well as percentage variation explained  $> 1\%$ . Significance thresholds were also estimated by simulation using the sim.sig.thr function in R/mpMap to produce a more conservative cut-off.

QTL results for all traits and analysis approaches were compiled by trait and grouped into unique QTLs based on proximity on both the genetic map and the physical positions of the markers on the IWGSC RefSeq v1.0 genome, as well as similarity of overall QTL peak shape, i.e. if the peak is in approximately the same location, are the patterns of  $p$ -values of individual markers/haplotypes also the same? QTL were then numbered in order of significance, i.e. the QTL with the highest significance value in any analysis was considered 'QTL1', and the peak marker of this QTL was assigned as the reference marker. Given the apparent inflation of  $p$  values in the CIM analysis, only the IBS, IBD and IM significance values were used for this ordering of QTLs. QTLs that were only found in the CIM analysis were given a separate numbering system with an "X" prefix (to indicate a weakly supported QTL) and listed after all the other QTL. On rare occasions where it was unclear if two QTLs were distinct or not, QTL suffixes were used of the form 'QTL4a' where the main QTL is 'QTL4' and 'QTL4a' may or may not be the same. QTL found only in individual trials, but not the meta-analysis of all trials, were prefixed with a trial letter (e.g. B = Bayer)

### **3.4.3 Genomic prediction**

Genomic prediction for grain yield (YLD) and three other traits of interest (PHT, SPW, SKT) were made using ridge regression in the rrBLUP package v 4.6 (Endelman, 2011) using R v 3.3.3 (R Core Development Team, 2013). Training models were built following filtering of individuals for missing phenotypic data. Genotypic marker data, where absent, were imputed using the column means of each specific marker of interest. Models were trained on the full set of available phenotypic data and their counterpart genotype data. For yield, models were built using three years of pre-project data (2012, 2013, 2014), and the 2015 project-derived data; for all other traits, models were trained on the 2015 phenotypic data. Models were assessed using the Spearman rank correlation between model derived phenotype predictions and known phenotypic outcomes for the test set. This was compared against the true Spearman rank correlation between the observed phenotypes of the training and test sets. This methodology was used to emulate a true breeding programme, in which the absolute within year values may alter due to environment, but selection would be imposed based on the rank order of lines within a year.

### 3.5 KASP marker development

Selected SNPs from the 90k array were converted to the Kompetitive Allele-Specific PCR (KASP) genotyping platform (LGC Genomics, UK) SNP flanking DNA sequences were used to design KASP primers using the software PolyMarker (Ramirez-Gonzalez et al. 2015), supplemented by manual inspection and design. DNA oligo 'tails' were added to the each of the two allele-specific primers: 5'-GAAGGTCGGAGTCAACGGATT-3' (VIC), 5'-GAAGGTGACCAAGTTCATGCT-3' (FAM). Primers were ordered from Sigma-Aldrich, and suspended using PCR-grade water (Sigma-Aldrich) to a concentration of 100 µM. DNA for the NIAB Elite MAGIC founders was extracted as described above, concentrations determined using a Nanodrop 200 spectrophotometer (Thermo Scientific), and diluted to a final concentration of 10 ng/µl using sterile PCR-grade water (Sigma-Aldrich). To determine whether KASP markers were co-dominant (i.e. able to robustly detect heterozygote alleles), 50:50 by-volume DNA mixtures of founders contrasting for SNP allele call were made, and included as controls during marker validation experiments. DNA KASP amplification reactions were undertaken using KASP Master Mix (LGC Genomics, UK), following the manufacturers guidelines. For each assay, reaction volumes were: 2.5µl KASP V4.0 2x Master Mix, 0.07µl KASP primer mix (for primer details, see **Appendix 1**) and 2.5µl DNA template (or 2.5 µl PCR-grade water for negative controls). KASP products were visualised using ProFlex PCR System Thermocycler (Applied Biosystems) using the following conditions: 1 cycle at 94 °C for 15 mins; 10 cycles at 94 °C for 20 s, 65 °C for 60 s with a touchdown of -0.8 °C/cycle to 57 °C; 35 cycles at 94 °C for 20 s, 57 °C for 60 s; final hold at 10 °C. Fluorescence of VIC and FAM fluorophore 5' end labelled PCR products were subsequently read using a Scientific QuantStudio™ 12K Flex Real-time PCR System (Thermo Fisher Scientific). ROX was used as a passive fluorescent reference to allow normalisation of variations in signal caused by differences in well-to-well liquid volume, following the manufacturer's instructions (LGC Genomics). Data was further analysed and visualised using Excel (Microsoft, USA), and the resulting allele calls compared to the corresponding SNP calls from the Illumina 90k SNP array (Gardner et al. 2016).

### 3.6 Bioinformatic analysis

SNPs identified as defining the boundaries of QTL intervals based on the MAGIC Genetic map were anchored to the wheat genome physical map (cv. Chinese Spring, IWGSC RefSeq v1.0. IWGSC, 2018) by BLASTn analyses, as described by Gardner et al. (2016). IWGSC RefSeq v1.0 high- and low-confidence gene models present within the intervals defined were outputted, along with their corresponding gene annotation data, and candidate genes identified manually by reference to the literature. Preliminary gene expression information (versus cultivar, tissue type, treatment type) for candidate genes was obtained from the data summarised in public repositories (<http://www.wheat-expression.com/cite> and [http://bar.utoronto.ca/efp\\_wheat/cgi-bin/efpWeb.cgi](http://bar.utoronto.ca/efp_wheat/cgi-bin/efpWeb.cgi)).

### 3.7 Identification of candidate genes and TILLING lines.

Candidate genes were identified after analysis of relevant published literature for wheat and related cereal species. Identification of wheat homologues of known genes from other plant species were identified by using their predicted coding regions (CDS) as queries for BLASTn searches of the wheat genome (assembly: RefSeq v.10; annotation: RefSeq v1.0. IWGSC, 2018). In order to identify artificially induced mutants in candidate genes, CDS were used for BLASTn searches of publicly available Targeting Induced Local Lesions IN Genomes (TILLING) mutants created in *T. aestivum* cv. Cadenza and *T. durum* cv. Kronos (Krasileva et al. 2017), undertaken using an online search tool (<http://www.wheat-tilling.com/>). Mutations in all homoeologues of each candidate gene were sought, and mutations ranked by predicted effect on the protein model, as: premature stop codon > splice-acceptor/donor mutation > non-synonymous mutation in a conserved protein domain (with consideration of SIFT score, which summarises the predicted effect of a given change in amino acid). Selected TILLING mutants were ordered from the SeedStor, JIC, UK (<https://www.seedstor.ac.uk/>), grown in 1 litre pots under glasshouse conditions, leaf tissue sampled for DNA extraction, and the developing ears bagged allowing selfed seed to be collected.

### 3.8 NIL development

For any given target QTL, NIAB Elite MAGIC progeny were identified that were heterozygous across the QTL interval, using the existing 90k SNP data for 643 lines. By comparing the pattern of allele calls (AA, A:B, B:B) for each SNP within the QTL region with the allele calls of the founders, the parental origin of the two alleles present in the region of heterozygosity was determined, where possible. These parental contributions to the heterozygous target region in a given line was then compared to the predicted allelic effects at the QTL peak, as outputted from CIM analysis. MAGIC lines carrying founder alleles with the greatest predicted contrast in phenotypic effect were prioritised for NIL development. Sib F5 seed for each of these MAGIC lines was germinated, DNA extracted, and genotyped with 2-3 codominant KASP markers (developed and validated as described above), allowing individuals to be classified as homozygous (A:A or B:B) or heterozygous (A:B). Each individual was grown to maturity in glasshouse conditions, ears bagged, and selfed seed collected. Where lines carrying contrasting homozygous alleles at a QTL were identified, these were sown in autumn 2018 in 1x1 m plots (consisting of 6x1 m rows per plot) for field bulking and subsequent preliminary phenotyping in summer 2019 (post project).

## 4. Results

### 4.1 Field trials, phenotyping and trials analysis

*Delivery of Milestones: M1.1, M1.2, M1.3, M1.4, M2.1, M2.2, M2.3, M2.4.*

#### 4.1.1 Trials and phenotyping

Across both the 2015 and 2016 season, 1109 NIAB Elite MAGIC lines, along with the eight founders and one control variety (KWS Santiago), were trialled in yield plots at five UK sites across, as detailed in **Table 2**. At each site, between 390 and 464 MAGIC lines were grown in partial replication, along with four replicates of each of each MAGIC founder and 4-8 replicates of KWS Santiago (e.g. **Figure 3**), resulting in the delivery of 4,996 yield plots across 2 years. Eighteen traits targeting yield, yield components and agronomic traits were measured across the trials and predicted means generated using block and spatial analyses, with one model per analysis taken forward, resulting in ~90,000 phenotypic data points. The full list of MAGIC lines used and the corresponding BLUPs for all phenotypes are listed in **Appendix 2**.

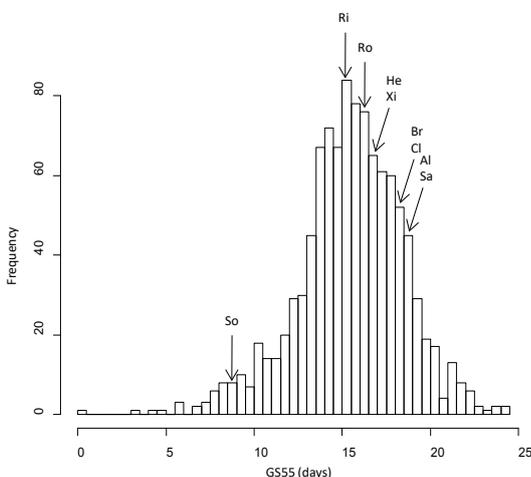


**Figure 3.** Example of one of the 10 project field trials: establishment and early growth in the BAY-YLD-15 trial located at NIAB Cambridge in 2015. The trial consisted of 504 plots of 2x6 m each (2x4 m harvested area), arranged in a randomised block design within a matrix of 8 rows of 63 plots each. In this trial, 420 MAGIC lines were grown in partial replicate: 44 lines in 2 replicates and 376 lines in 1 replicate. Four replicate plots of each of the 8 MAGIC founders were included, along with 8 replicates of the control variety KWS Santiago.

#### 4.1.2 Transgressive segregation

A key advantage of MAGIC mapping populations is the extensive mixing of the alleles underlying quantitative traits, helping to ensure transgressive segregation of phenotypes. Considering the pre-harvest data, almost all traits showed significant transgressive segregation in both directions (**Appendix 2, Figure 4**): all but one trait had  $\geq 5\%$  of lines transgressive in at least one direction, and

for 50% of traits,  $\geq 5\%$  of lines were transgressive in both directions. Of the pre-harvest traits, and across all sites and seasons, flowering time (as measured by growth stage 55, GS55) and spikelet number per ear (SKT) had the lowest number of lines showing transgressive segregation. On average, GS55 showed 18% total transgressive segregation, due largely to the low numbers of lines ( $<4\%$ ) reaching GS55 before the founder Soissons, which carries the early flowering allele at the photoperiod response locus *PPD-D1*; on average, SKT traits showed 22% total transgressive segregation, due largely to the low numbers of lines ( $<3.5\%$ ) with lower spikelet number than the founder Soissons. The highest levels of transgressive segregation (50%) were observed for plant height (PHT). This was as expected since all the founders carry exactly one semi-dwarfing allele at either the *RHT-B1* or *RHT-D1* locus which have a major effect on plant stature, but the lines can carry 0 or 1 semi-dwarfing allele (i.e. they can be 'tall' due to the absence of *RHT1* semi-dwarfing alleles). Note: lines with 2 semi-dwarfing alleles have been consciously removed from the population due to their excessively short stature. As a likely consequence of this, consistently higher proportions of lines across all sites and years were higher than the tallest founder (mean = 35%), compared to the proportion of lines that were shorter than the shortest founder (mean = 15%). Spikelet number (SKT) also showed a similar trend across sites and years, with 19% lines possessing more spikelets than the founder with the highest number of spikelets, compared to 3% of lines possessing fewer spikelets compared to the founder with the lowest number of spikelets. Grain yield (YLD) showed relatively high mean levels of transgressive segregation (35%), which was predominantly driven by lines with lower yield than the lowest founder (Hereward), and with 3% of lines on average yielding more than the highest yielding founder (Alchemy). Examining the data in more detail, it is apparent that differences in transgressive segregation between years could be found for a number of traits. For example, specific weight (SPW) in 2015 and 2016 (21% versus 27%), which is driven by the proportion of lines with increased SPW compared to the maximum founder in 2016 compared to 2015 (12% versus 3%), while the proportion SPW below the minimum founder remained unchanged between the years ( $\sim 17\%$ ). Overall, this suggests that substantial inter-annual differences for key yield component traits may have occurred during this project.



**Figure 4.** Example of transgressive segregation in the NIAB Elite MAGIC population. Trait: predicted means for growth stage 55 (GS55) from the meta-analysis of all 2016 season trials. The trait scores for the 8 founders (Al = Alchemy, Br = Brompton, Cl = Claire, He = Hereward, Ri = Rialto, Ro = Robigus, So = Soissons, Xi = Xi19) and the control variety KWS Santiago (Sa), are indicated.

### 4.1.3 Trait heritabilities

Heritabilities for all traits, per site, per year, and meta-analysis, are listed in **Table 4**. Heritability for plant height (PHT) and flowering time (GS55) was ~0.90. For specific weight (SPW), heritability was slightly lower at ~85%, but with notably high heritabilities in the Limagrain 2015 and RAGT 2016 trial sites ( $h^2 = 0.98$  and  $0.94$ , respectively). Heritability was around 0.80 for spikelet number (SPK), and ~0.85 for yield (range: 0.80 for meta-analysis 2016, to 0.93 for the RAGT trial in 2015). The heritability for tiller number (TIL) was lower, and more variable between sites and years (mean = 0.27, minimum = 0 for Elsoms 2015 and RAGT 2015, maximum 0.67 for KWS 2016), reflecting the difficulty in accurately phenotyping this trait at scale in the field. All of the post-harvest grain traits showed high heritability, with thousand grain weight (TGW), seed length (LEN), width (WDT), length to width ratio (LWR) and area (ARE) all possessing mean  $h^2 > 0.82$ . Seed number per ear (SNO) and factor form density (FFD) had mean heritabilities of ~0.70. FFD is a measure of seed density, and is a derived trait calculated by dividing seed area by TGW, possibly explaining its slightly lower heritability compared to the majority of other seed traits.

## 4.2 QTL analysis

*Delivery of Milestones: M3.1, M3.2, M3.3, M3.4*

**Aim:** to identify QTL for yield and yield components in the NIAB Elite MAGIC population.

**Process:** For the subset of 643 MAGIC lines with SNP data, QTL analyses were undertaken using 4 approaches (IBS, IBD, IM and CIM, as described in the Methods section), identifying 376 significant ( $P > 0.05$ ) QTL loci across all 21 wheat chromosomes for 18 traits (**Appendix 3a**), and as summarised in **Table 5**. For each trait, all QTL detected by any analysis methods were sorted by map location and grouped into unique QTLs based on proximity on both the genetic map and the physical positions of the markers on the wheat reference genome assembly for cultivar 'Chinese Spring 42', IWGSC RefSeq v1.0 (IWGSC, 2018), as well as similarity of overall QTL peak shape. They were then numbered in order of significance, i.e. the QTL with the highest significance value in any analysis was labelled as 'QTL1'. Given the apparent inflation of  $p$  values in the CIM analysis, only the IBS, IBD and IM significance values were used for the ordering of QTLs. The mean and median number of QTL per trait was 19 and 12, respectively. The highest and lowest significance value of QTL per trait were 8 (for awn presence/absence, AWN) and 42 (specific weight), respectively. Chromosome 6A had the highest number of QTL, with 34 loci identified across 14 traits. The chromosomes with the lowest number of QTL were 1D and 5D, both of which had 9 QTL.

Trait	Year	META	BAY	ELS	KWS	LIM	RAG
PHT	2015	0.93	0.95	0.92	0.92	0.94	0.98
	2016	0.92	0.93	0.85	0.93	0.93	0.97
GS55	2015	0.79	0.96	0.58	0.98	0.86	0.94
	2016	0.93	0.89	0.93	0.98	0.88	0.94
YLD	2015	0.72	0.81	0.86	0.91	0.89	0.93
	2016	0.80	0.89	0.83	0.89	0.83	0.92
SPW	2015	0.72	0.73	0.86	0.82	0.98	0.87
	2016	0.76	0.89	0.86	0.89	0.90	0.94
LLE	2015	NA	NA	NA	NA	NA	NA
	2016	0.66	0.48	0.69	0.78	NA	0.94
LOD	2015	NA	NA	NA	NA	0.80	0.80
	2016	0.41	NA	0.51	0.95	NA	NA
TIL	2015	0.38	0.62	0	0.56	NA	0
	2016	0.28	0.14	0.13	0.67	0.14	0.11
SKT	2015	0.89	0.89	0.92	0.9	0.89	0.84
	2016	0.81	0.72	0.63	0.72	0.62	0.91
TGW	2015	0.87	0.8	0.86	0.93	0.92	0.79
	2016	0.89	0.92	0.90	0.93	0.87	0.69
LEN	2015	0.88	0.83	0.96	0.86	0.91	0.65
	2016	0.93	0.94	0.94	0.99	0.83	0.67
WID	2015	0.79	0.74	0.82	0.83	0.88	0.75
	2016	0.74	0.92	0.85	0.95	0.89	0.65
ARE	2015	0.84	0.75	0.90	0.77	0.93	0.60
	2016	0.91	0.93	0.90	0.97	0.86	0.68
LWR	2015	0.86	0.84	0.93	0.89	0.84	0.63
	2016	0.94	0.96	0.94	0.98	0.87	0.59
FFD	2015	0.59	0.72	0.78	0.65	0.77	0.25
	2016	0.81	0.93	0.87	0.85	0.87	0.63
SNO	2015	0.79	0.80	0.74	0.82	0.51	0.37
	2016	0.71	0.76	0.75	0.85	0.71	0.44

**Table 4.** Trait heritabilities. Traits: PHT (plant height), GS55 (growth stage 55), YLD (grain yield), SPW (specific weight), LLE (lodging x leaning score), LOD (lodging), TIL (tiller number), SKT (spikelet number per ear), TGW (thousand grain weight), LEN (seed length), WID (mean seed width), ARE (seed area), LWR (seed length to width ratio), FFD (factor form density), SNO (seed number per ear), PSH (percent shrivelled seed). Trials: BAY (BAY-YLD-16), ELS (ELS-YLD-16), KWS (KWS-YLD-16), LIM (LIM-YLD-16), RGT (RGT-YLD-16). NA = not applicable (trait either not expressed, or not measured).

Trait	No.																						
	QTLs	1A	1B	1D	2A	2B	2D	3A	3B	3D	4A	4B	4D	5A	5B	5D	6A	6B	6D	7A	7B	7D	U
AWN	8				1	2								1			1	2			1		
GS55	41	2	3	2		3	3	3			2	3	1	3	2	1	2	2	3	1	4		1
LEA	4		1	1									1										1
LLE	6		1			1						1	1		1					1			
LOD	7		1	1						1		1	1				1			1			
PHT	13			1		1	2	3	1	1		1	1				2						
SKT	30			1	1	1	2	1	1	1	2	2	1	1	1	2	3	1	2	3	3	1	
SPW	42	1	1	1	1	5	1	2	1	1	3	3	2	3	1	2	2			5	5	1	1
TILL	11	3	2											3			2					1	
YLD	23			1		1		2		1	4	3	2		1	1	1	2	1	1	1	1	
ARE	29	2	3		1	1	1		2	1	1	4		3	4	1	2			1		1	1
CIR	25		3			2	1					2	1	3	2		6	1	1	1			2
FFD	9				1			1	1	1		1	1	2								1	
LEN	30	2	2		3			2	3	1	2	1	2	1	3	1	3	1		2		1	
LWR	26	1	5		1		3	3				1	1	4	1		3		1	1			1
PSH	11	1			1		1			1	1	1	1	1					1	1		1	
SNO	12	1			1		2				1	1	1	1				1	1	2			
TGW	26	1	2	1	1	1		1	2	1	2	2	1	1	3	1	3					2	
VWT	12		1		1		1	1	1	1		2	1	2	1								
WID	11	1	1		1			1				2	1				3		1				
	<b>Totals:</b>	<b>15</b>	<b>26</b>	<b>9</b>	<b>14</b>	<b>18</b>	<b>17</b>	<b>20</b>	<b>12</b>	<b>11</b>	<b>18</b>	<b>31</b>	<b>20</b>	<b>29</b>	<b>20</b>	<b>9</b>	<b>34</b>	<b>10</b>	<b>11</b>	<b>20</b>	<b>14</b>	<b>10</b>	<b>7</b>

**Table 5.** Summary of distribution of QTLs across chromosomes, detected for the 20 yield/yield component traits investigated in 10 sites over 2 years. Here, the lodging trait (LOD) has also been separated into leaning (LEA) and lodging + leaning (LLE) traits.

To identify chromosomal regions containing QTL shared between traits, all 2015 and 2016 season QTLs were sorted based on the genetic map (Gardner et al. 2016), and then by physical map location (IWGSC, 2018). Predicted founder allelic effects at the QTL peaks were also compared, in order to help determine QTL clusters. To help determine the expected resolution of QTL mapping, the highly non-recombining regions spanning the centromeres, as determined by comparison of SNP genetic map position versus physical map position, were identified. Where QTL mapped to these regions, they were generally grouped together into a single cluster of QTL, due to the lack of genetic recombination present. A total of 95 chromosomal regions were found to share QTL for two or more traits (**Table 6; Appendix 3b**). Of these, the following QTL number per chromosomal location combinations were found: 2 QTLs (33 locations), 3 QTL (23 locations), 4 QTL (11 locations), 5 QTL (11 locations), 6 QTL (10 locations), 7 QTL (4 locations), 8 QTL (1 location), 12 QTL (1 location) and 15 QTL (1 location). The major semi-dwarfing loci *RHT-B1* and *RHT-D1* possessed the highest number of co-localizing QTL (12 and 15, respectively). Additionally, the major photoperiod response locus *PPD-D1* had 6 co-localising QTL. The major QTL on chromosome 7A known to affect spikelet number (e.g. Quarrie et al 2005, 2006; Boeven et al. 2016) co-localised with four other grain traits (PSH, SNO, SPW and LWR). A total of 34 of the 95 multi QTL regions contained QTL for flowering time (GS55), which is known to often have pleiotropic effects on other traits. For grain yield, 24 QTL that co-localised with one or more QTL for other traits were detected. Of these, seven co-localised with QTL for GS55, and one for lodging (chromosome 2B). An additional three yield QTL co-localised with plant height, and one yield QTL co-localised with a tiller number QTL on chromosome 7D. The remaining co-localising yield QTL were associated with QTL for grain traits, sometimes in conjunction with QTL for spikelet number. Finally, excluding the *RHT-B1* and *RHT-D1* loci, a subset of five chromosomal regions contained notably high ( $\geq 7$ ) numbers of QTL. Of these, one was located within a highly non-recombining region of the chromosome, and so was excluded from further analysis here (chromosome 5A, location 3. Termed 5A-3). Location 3A-1 contained QTL for six grain traits (CIR, FFD, LWR, SPW, TGW, WID), as well as a QTL for yield. Location 3D-1 contained QTL for five grain traits (ARE, FFD, LEN, SPW, VWT), as well as WTL for spikelet number and plant height. Location 5B-2 contained QTL for five grain traits (ARE, LEN, TGW, VWT, SPW) as well as QTL for spikelet number and lodging. Location 5A-4 contained QTL for six grain traits (CIR, FFD, LWR, PSH, SNO, VWT) as well as QTL for spikelet number and flowering time.

Notes	Chr.	Locus No.	No. QTL	Traits
C	1A	1	3	PSH, SNO, TIL
C	1A	2	2	LWR, SNO
	1A	3	3	GS55, LEN, SPW
	1A	4	3	GS55, LEN, TIL
	1A	5	2	ARE, WID
	1A	6	3	ARE, LEN, TGW
	1B	1	2	GS55, LOD
	1B	2	3	CIR, LOD, LWR
C	1B	3	6	ARE, CIR, GS55, LEN, LWR, TGW
	1B	4	2	CIR, LWR
	1B	5	2	LWR, LWT
	1B	6	4	CIR, LWR, SPW, VWT
	1B	7	5	CIR, LEN, LWR, TGW, WID
	1D	1	3	LEA, LOD, TGW
	1D	2	2	GS55, YLD
C	2A	1	4	LEN, LWR, SKT, SNO
C	2A	2	3	PSH, SNO, TIL
C	2A	3	2	LWR, PSH
	2A	4	5	SKT, SNO, SPW, TIL, VWT
	2A	5	4	PSH, SKT, SNO, WID
	2A	6	4	ARE, FFD, LEN, TGW
	2B	1	2	LOD, YLD
	2B	2	2	AWN, SPW
	2B	3	2	ARE, GS55
PPD-D1	2D	1	6	ARE, CIR, GS55, PHT, SKT, SNO
	2D	2	4	GS55, LWR, PHT, SKT
	2D	3	3	CIR, LWR, SPW
	3A	1	7	CIR, FFD, LWR, SPW, TGW, WID, YLD
	3A	2	3	LWR, TGW, YLD
	3A	3	3	LWR, SPW, YLD
C	3A	4	2	PHT, VWT
C	3A	5	3	GS55, LEN, PHT
	3A	6	2	GS55, PHT
	3A	7	3	CIR, LEN, LWR
	3A	8	4	CIR, GS55, LEN, LWR
	3B	1	2	SKT, VWT
C	3B	2	4	ARE, FFD, LEN, TGW
	3B	3	3	ARE, LEN, SPW
	3D	1	7	ARE, FFD, LEN, PHT, SPW, SKT, VWT
	3D	2	6	FFD, PHT, PSH, SKT, TGW, YLD
	4A	1	3	SPW, TGW, YLD
	4A	2	2	GS55, SPW
	4A	3	2	LEN, YLD
	4A	4	6	GS55, LEN, PSH, SKT, SNO, YLD
	4A	5	2	SPW, YLD
	4A	6	6	ARE, GS55, LEN, SKT, TGW, YLD
RHT-B1	4B	1	4	FFD, LWR, SKT, VWT
	4B	2	12	ARE, CIR, FFD, LWR, PHT, PSH, SNO, SPW, TGW, WID, YLD, VWT
	4B	3	3	SPW, VWT, YLD
	4B	4	2	SPW, VWT
	4B	5	6	GS55, LLE, LOD, SKT, VWT, YLD
	4B	6	5	WID, ARE, GS55, LEA, SPW
	4B	7	5	ARE, CIR, GS55, LEN, TGW
RHT-D1	4D	1	3	GS55, SPW, TGW
	4D	2	15	CIR, FFD, LEA, LEN, LLE, LOD, LWR, PHT, PSH, SKT, SNO, SPW, VWT, WID, YLD
	5A	1	4	CIR, CIR, LWR, TIL
	5A	2	5	ARE, CIR, GS55, LEN, LWR
C	5A	3	7	CIR, ARE, CIR, GS55, LWR, TGW, TIL
	5A	4	8	CIR, FFD, GS55, LWR, PSH, SKT, SNO, VWT

	5A	5	3	TIL, SPW, VWT
	5A	6	2	ARE, SPW
	5A	7	3	AWN, FFD, GS55
	5B	1	3	LEN, TGW, YLD
	5B	2	7	ARE, LEN, LLE, SKT, SPW, TGW, VWT
	5B	3	4	SPW, ARE, TGW, VWT
	5B	4	5	LEN, ARE, CIR, GS55, LWR
	5D	1	5	ARE, LEN, TGW, SKT, YLD
	6A	1	2	CIR, GS55
	6A	2	2	CIR, LEN
	6A	3	5	CIR, LWR, SKT, SPW, YLD
	6A	4	4	CIR, PHT, TGW, WID
	6A	5	2	CIR, LWR
GW2 C	6A	6	6	ARE, LEN, PHT, TGW, TIL, WID
	6A	7	2	SKT, TGW
	6A	8	6	CIR, GS55, LEN, LWR, SKT, TIL
	6A	9	5	ARE, LEN, LOD, SPW, TGW
	6A	10	2	AWN, WID
	6B	1	2	GS55, YLD
	6B	2	2	GS55, YLD
	6B	3	2	SKT, SNO
	6D	1	6	CIR, GS55, LWR, PSH, SNO, YLD
	6D	2	2	GS55, SKT
	6D	3	2	GS55, WID
	7A	1	2	ARE, SPW
	7A	2	5	ARE, LEN, SKT, TGW, YLD
7A-SKT	7A	3	5	PSH, SKT, SNO, SPW, LWR
	7A	4	6	GS55, LLE, LOD, SKT, SNO, SPW
VRN3	7B	1	2	GS55, SKT
	7B	2	2	GS55, SPW
	7B	3	2	SKT, SPW
	7B	4	3	SKT, SPW, YLD
	7B	5	3	GS55, SKT, SPW
	7D	1	3	FFD, TGW, TIL
	7D	2	3	PSH, TIL, YLD
	7D	3	2	ARE, TGW

**Table 6.** Summary of co-localised QTL identified in the MAGIC population. Trait abbreviations as listed in the Methods section. Locations of known major effect genes are indicated. 7A-SKT = the major effect QTL for spikelet number previously reported (e.g. Quarrie et al. 2005, 2006; Boeven et al. 2016). C = QTLs that lie within the highly non-recombining regions surrounding the centromere, based on comparison of the genetic map position (Gardner et al. 2016) versus the physical map (IWGSC, 2018).

### 4.3 Development of molecular markers tagging QTL

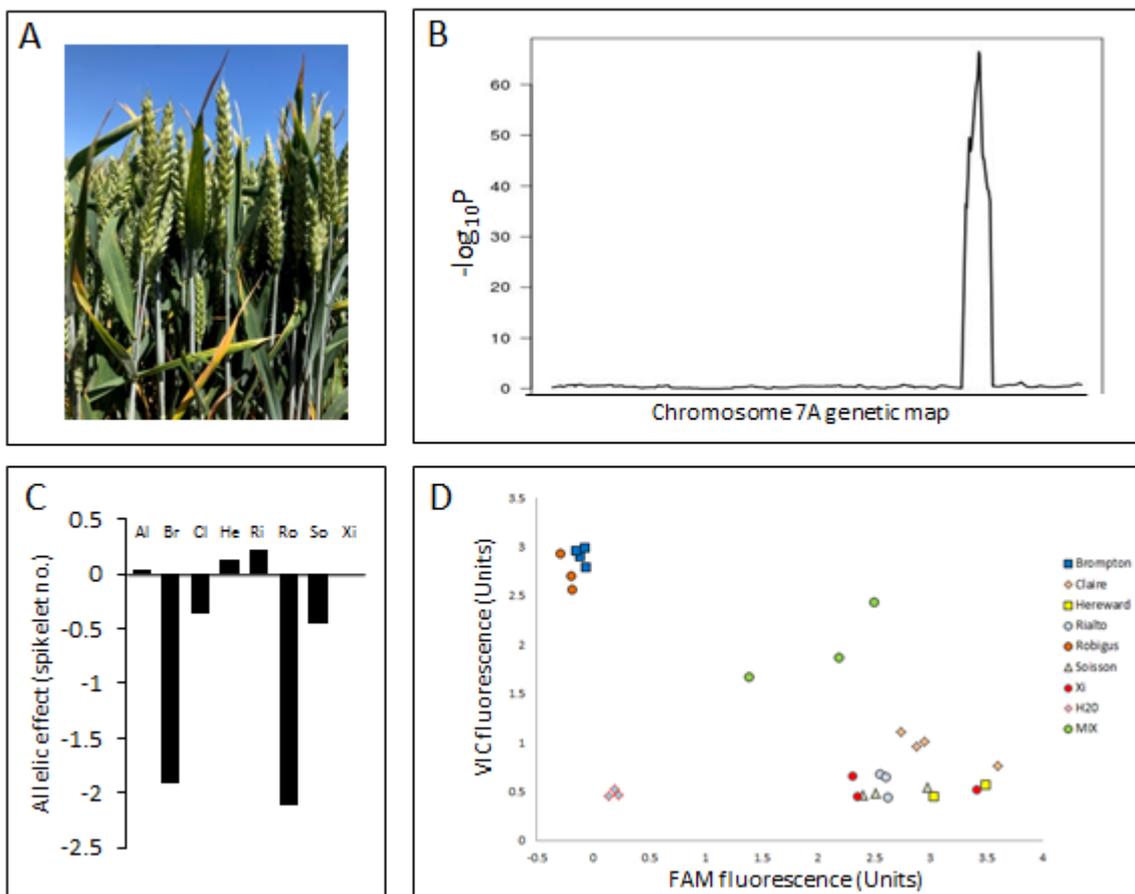
*Delivery of Milestones: M4.1, M4.2*

**Aim:** to develop simple, cheap molecular markers for genetic variants identified using the 90k SNP array as tagging the most promising yield and yield component QTL identified in the MAGIC population.

**Process:** Of the 376 QTL identified, 20 were selected for the development of KASP genetic markers (primers listed in **Appendix 1**, prioritised QTL listed in **Appendix 4**). KASP markers allow flexible and cheap genotyping of individual SNPs identified using the 90k SNP array, and are commonly used by the project partner wheat breeding companies. For each target QTL, 3-5 of the most significant SNPs within the QTL peak were selected for conversion to KASP. For a given QTL, SNPs were prioritised for KASP conversion where allele calls (A:A or B:B) matched the pattern of predicted allelic effects of the founders for the target QTL. For example, where A:A alleles at a target SNP were carried by founders predicted to confer increased grain size, and B:B alleles were carried by founders predicted to confer decreased grain size. Note: due to QTL overlap, the 20 QTL selected were represented by 18 genomic regions

KASP markers were validated by genotyping DNA extracted from each of the eight founders, as well as a 50:50 mix of two founders known to carry contrasting alleles at the SNP – creating an ‘artificial heterozygote’ (allele A:B), making it possible to determine if the KASP marker tested was co-dominant (i.e. able to detect heterozygotes). For each SNP, KASP genotyping results from the founders were compared to those derived from the 90k SNP array. In this way, a total, 58 co-dominant KASP markers were validated for the 20 QTL targeted, following the methods listed in Section 3.5. These markers displayed the following ‘KASP number per QTL’ summary statistics: median = 3, mode = 3, minimum = 1 (QTLs for seed weight on chromosome 4B and yield on 7D), max = 8 (for the co-locating QTLs TIL\_6A, FFD\_6A, WID\_6A, TGW\_6A.1, TGW\_6A.2). An example of the process of converting SNPs from the 90k Illumina SNP array to the KASP genotyping platform is illustrated in **Figure 5**.

**Deliverables:** 58 ‘breeder friendly’ co-dominant KASP markers tagging 20 yield/yield component QTL were successfully validated, and the associated information on how to use these in practice distributed to all project partners. Based on the KASP genotyping platform used by all industrial project partners, these resources provide the partnering wheat breeding companies the ability to explore the tracking and manipulation of beneficial alleles for multiple QTL within their breeding programmes.



**Figure 5.** Illustration of the process for development of KASP markers tagging targeted QTL. Trait: spikelet number per ear (SKT). (A) Field trials and phenotypic data collection. (B) Genetic mapping to identify QTL, genetic map of chromosome 7 shown. (C) Predicted allelic effects at the QTL determined. Al = Alchemy, Br = Brompton, Cl = Claire, He = Hereward, Ri = Rialto, Ro = Robigus, So = Soissons, Xi = Xi19. Brompton and Robigus are predicted to carry alleles with strong effect on reducing spikelet number per ear. (D) Conversion of a SNP at the peak of the QTL to the KASP genotyping platform. Brompton and Robigus (A:A) have contrasting alleles to all other MAGIC founders (B:B). A 50:50 mix of DNA from founders with contrasting allele calls (MIX) was used to determine whether the marker was capable of detecting heterozygote alleles, and so represent a co-dominant marker. VIC and FAM fluorescence is shown on the X and Y axis, respectively.

#### 4.4 Development of near isogenic lines for selected QTL

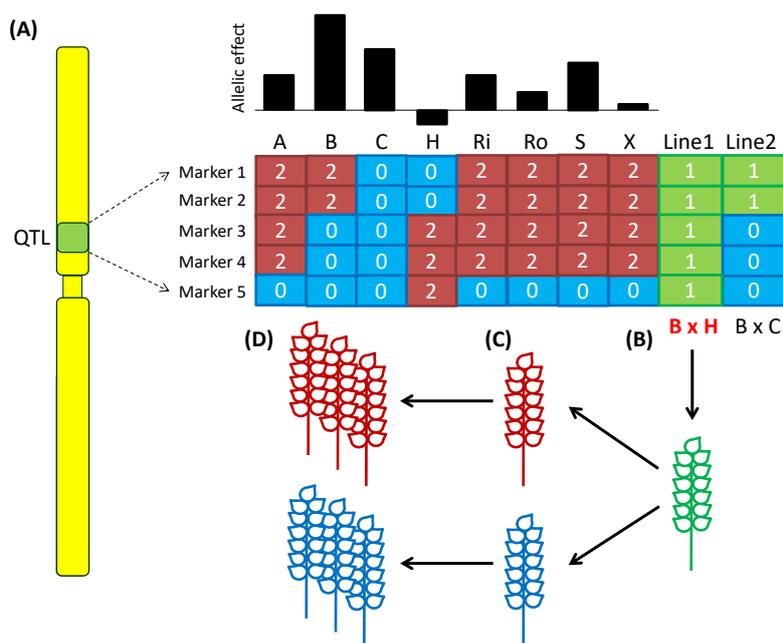
Delivery of Milestones: M4.4, M4.5 (for M4.5, see also Section 4.5)

**Aim:** to develop wheat 'nearly isogenic lines' (NILs) for target QTL, providing the resources for future studies to undertake detailed characterisation of QTL in isolation, and provide the basis for map-based cloning - the identification of the specific gene and genetic variants underlying each QTL.

**Process:** For any given QTL identified in the MAGIC population, a NIL pair represents a pair of lines that are on average 98% identical at the DNA level, but contrast for alleles carried at the target QTL. As the genetic background in NIL pairs has essentially been fixed, NILs can be used for precise characterisation of QTL effects largely independent of the effects of interacting alleles present in other regions of the genome. For example, NILs can be used for: (1) detailed phenotypic characterisation of a QTL, at resolutions ranging from evaluation in field experiments, right down to investigations at the sub-cellular level, (2) detailed analysis of differences in gene expression associated with each allele, and (3) a NIL pair can be inter-crossed to generate genetic recombination within the QTL interval that may ultimately result in allowing the underlying genetic variant to be identified. Therefore, NILs represent powerful resources for precise evaluation of QTL phenotypic effects, and ultimately the development of 'perfect' genetic markers for use in breeding programmes. Such perfect markers differ from 'linked markers' (such as those developed here in Section 4.3) in that they directly assay for the specific genetic variants in the underlying gene, rather than genetic variants at locations closely linked to the underlying gene. NIL development normally takes many generations to complete. However, we used the residual heterozygosity present in the MAGIC recombinant inbred lines (RILs) to identify NILs in just one generation, following the methodology listed in Section 3.8. This was possible because the RILs used for 90k SNP genotyping were at the F5 stage of inbreeding, and so still expected to possess low levels of heterozygosity (~2%). Therefore, of the 643 MAGIC F5 RILs, on average ~13 lines are expected to be heterozygous for any given chromosomal location. The process of NIL development is illustrated in **Figure 6**.

We used the existing 90k SNP data (NIAB unpublished data, curated from the data published by Gardner et al. 2016) to identify F5 MAGIC RILs heterozygous at each of our target yield/yield component QTLs. Between 2 and 4 RILs were selected for each of the 20 QTL prioritised in Section 4.3 (**Appendix 4**). As the F5 RIL seed is old, germination was often low. Despite this, we were able to germinate  $\geq 1$  individuals for 34 of the 46 RILs selected for 18 of the 20 QTL (**Appendix 5a**). Five RILs were heterozygous across more than one target QTL: RIL10 for WID\_3B / SPW\_3B and TGW\_3B, RIL33 for FFD\_6A and WID\_6A, RIL36 for LEN\_1B.2 and TGW\_6A, RIL42 for WID\_6A and TIL\_6A, and RIL44 for TGW\_6A and WID\_6A. All individuals from each RIL that germinated were genotyped with the relevant co-dominant KASP markers developed in Section 4.3, in order to identify individuals carrying either homozygous A:A or homozygous B:B alleles. Where one or more individual of each homozygous allele class were identified (i.e.  $\geq 1$  A:A individual and  $\geq 1$  B:B individual), then the NIL pair was identified, and the relevant individuals grown to maturity and selfed seed collected. For RILs where it was not possible to identify  $\geq 1$  individual from each homozygous allele class, if  $\geq 1$  heterozygous F5 RIL individual was identified, these were selfed, the resulting seed grown, and homozygous A:A and B:B individuals sought via KASP genotyping. Following this approach, 22 of the 34 RILs that resulted in germination of  $\geq 1$  seed were found to either possess individuals for both of the homozygous allele classes (14 RILs, representing 12 QTL), or failing that,

possess heterozygotes - allowing individuals homozygous for each of the two alleles to be searched for via KASP genotyping in the subsequent generation.



**Figure 6.** Illustration outlining the process of developing a near isogenic line (NIL) pair for a target QTL. (A) QTL localised to a specific chromosomal region. (B) Allele calls at genetic markers at the QTL are analysed in the MAGIC recombinant inbred lines (RILs). RILs found to be heterozygous at the target QTL using the genotypic data derived from the 90k SNP array are identified (Line1, Line2 here). For each of these RILs, on a marker-by-marker basis, SNP allele calls present in the RIL is compared to SNP calls in the founders in order to determine, where possible, the two founders that contributed to the region of heterozygosity. RILs heterozygous across the QTL interval predicted to carry alleles from founders with contrasting allelic effects at the QTL are selected for NIL production - here, the region of heterozygosity within the QTL interval is predicted to come from the founders Brompton (B) and Hereward (H), themselves predicted to have strongly contrasting allelic effects for the target QTL. (C) Sib F5 seed for the selected RIL is grown, the plants genotyped using KASP markers assaying for SNPs at the QTL, and individuals carrying contrasting homozygous alleles (A:A or B:B) identified. These two individuals represent a NIL pair. They are grown to maturity, and selfed seed collected. (D). This seed is then grown and selfed seed collected, providing seed bulks for the NIL pair, to be used for downstream R&D, e.g. conformation of phenotype via field and glasshouse experiments.

These resources provided the potential to develop one or more NIL pair for all but two of the 18 QTL for which F5 seed was available: LEN\_3B and TIL\_5A. Post project, we are currently finalising the following NIL resources to underpin subsequent R&D:

(A) NIL field multiplication and preliminary phenotyping: 1x1 m nursery plots for seed bulking were sown in October 2018 for all 14 RILs found to possess individuals from both homozygous classes (**Appendix 5b**). Preliminary phenotyping of yield and yield component traits will be carried out in these nursery plots in summer 2019. For many RILs, the seed stocks available allowed more than one plot per homozygous allele to be sown (**Table 7; Appendix 5b**), which will allow phenotypic data to be collected from replicated plots. Interestingly, MAGIC RIL36 segregates for QTL at two locations: LEN\_1B.2 and TGW\_6A. Of the four possible allelic combinations at LEN\_1B.2 and TGW\_6A, three are present in the germplasm sown: B:B + B:B (2 plots), A:A + B:B (1 plot), B:B + A:A (5 plots). This germplasm will allow interactions at the phenotypic and gene expression level between the seed length QTL on chromosome 1B and the thousand grain weight QTL on chromosome 6A to be determined in fine detail within a single genetic background.

(B) Recovery of additional NILs: for 7 RILs, although both homozygous allele classes were not identified in our first screen, we were able to identify lines which were heterozygous at the target QTLs, which when selfed, will allow the identification of A:A and B:B individuals, and thus generate NIL pairs. We are currently screening this germplasm using the relevant KASP markers developed in Section 4.3. If successful, the project will have developed NILs for 15 yield/yield component QTL, available for downstream R&D. It would also be possible to develop NILs for additional QTL.

**Deliverables:** NILs for 15 yield/yield component QTL were either developed, or are close to finalisation, with 14 NILs for 12 QTL currently being grown for field seed bulking and to verify whether NIL pairs contrast for target phenotype. This germplasm, and the associated MAGIC molecular resources, provide the raw materials for future work to identify the underlying genes/genetic variants, with the aim of providing 'perfect markers' for exploitation in wheat breeding programmes.

Target QTL	MAGIC RIL code	Number of 1x1m field plots (allele A:A, allele B:B)
LEN_1B.2	RIL39	3, 4
LEN_1B.2	RIL36 <sup>†</sup>	12, 3
YLD_2B	RIL25	3, 1
YLD_3A	RIL31	1, 1
TGW_3B, WID_3B, SPW_3B	RIL10	2, 1
LEN_5A	RIL11	8, 5
LEN_5A	RIL26	2, 4
LWR_5D	RIL28	1, 1
TGW_6A	RIL36 <sup>†</sup>	3, 1
TIL_6A	RIL21	8, 6
SKT_7A	RIL9	12, 7
YLD_7D	RIL12	1, 1
YLD_7D	RIL13	5, 1

**Table 7.** NIL germplasm sown in October 2018 at NIAB, Cambridge for field multiplication and preliminary phenotyping. <sup>†</sup>MAGIC RIL36 segregates for QTL at two locations: LEN\_1B.2 and TGW\_6A. Of the four possible allelic combinations at these two QTL, three are present in the germplasm sown: B:B + B:B (2 nursery plots), A:A + B:B (1 plot), B:B + A:A (5 plots).

## 4.5 Identification of candidate genes and TILLING

*Milestones delivered: M4.4, M4.5 (see also Section 4.4)*

**Aim:** Identify candidate genes within selected QTL for the identification and development of TILLING resources, as well as for other future gene validation approaches.

**Process:** For each QTL, the identification of candidate genes (i.e. genes thought using existing published knowledge to possibly represent the underlying gene, due to involvement in the control of similar phenotypes in other plant species) was based on two steps:

(1) Definition of gene content within QTL regions: A subset of 53 yield and yield component QTLs, defining 44 chromosomal locations, was selected (**Appendix 4**). The physical map locations of the genetic markers marking the QTL boundaries were used to define QTL physical intervals and to identify the predicted gene content - based on IWGSC RefSeq v1.0 gene models for the reference wheat cultivar Chinese Spring. The smallest number of predicted gene models identified within a given QTL was 6 (for the seed area QTL on chromosome 5B), the largest was 549 (seed shape QTL on 4B) and the median was 65 (seed weight QTL on 3B). Predicted gene numbers and further information for all 53 QTL are listed in **Appendix 4**.

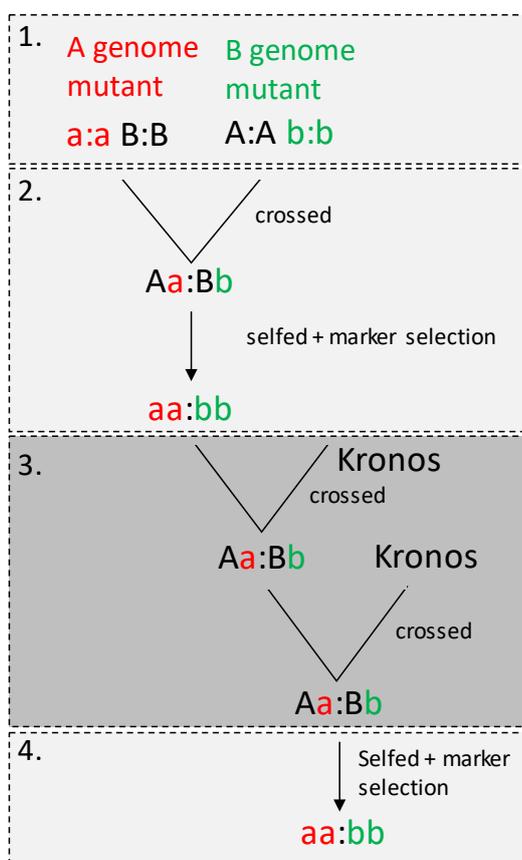
(2) Identification of candidate genes: using the publicly available functional annotations of the genes identified in step 1 above, and a review of the available literature for wheat and related cereal species, we searched for candidate genes within each QTL interval. For example, focusing on grain size characters, a list of 26 rice genes known to control grain size was identified in the literature, of which gene sequence was publicly available for 25 genes (**Appendix 6**). Using the CDS of these 25 rice genes as queries for BLASTn searches of the wheat genome identified 125 wheat homologues (**Appendix 6**). Based on established collinearity between the genomes of the related cereal species rice and wheat, 24 of the 25 rice genes were found to contain putative orthologues in wheat. The only exception was the rice gene *GRAIN SIZE 5* (GS5, MSU gene model LOC\_Os05g06660, Li et al. 2011), a putative serine carboxypeptidase for which wheat is predicted by collinearity to possess orthologous genes on the group 1 chromosomes, but for which we identified wheat homologues on chromosomes 3A, 3B and 3D (RefSeq v1.1 gene models TraesCS3A02G212900LC, TraesCS3B02G277100LC and TraesCS3D02G172900, e-values  $\leq 2e-27$ ). In total, across all 53 QTL, 41 candidate genes were identified (listed in **Appendix 4**). The number of candidate genes per QTL ranged from 0 (for 19 QTL) to 3 (for 7 QTL: seed width on 3B; seed length 3B; tiller number 5A; seed length on 5A; seed shape on 5D; TGW on 6A; spikelet number on 7A).

Determining candidate genes is a key step towards identifying the gene and genetic variant at a given QTL. In addition to seeking to capture additional genetic recombination, to help determine whether natural allelic variation at a given gene underlies a QTL, functional validation is commonly sought using germplasm carrying artificial mutations at the candidate gene of interest. In this way, genes located within a QTL interval can be prioritised or discounted as useful candidates. Termed 'Targeting Induced Local Lesions IN Genomes' (TILLING) populations, such resources represent collections of seed that have been artificially mutated using the chemical ethyl methanesulfonate (EMS). Treatment with EMS predominantly induces single nucleotide changes in the DNA, termed 'point mutations'. In wheat, TILLING populations are available in both tetraploid (cv. Kronos) and hexaploid (cv. Cadenza) wheat. These TILLING populations have been sequenced at the DNA level via a technique called exome capture, and the results databased in a way to allow users to search for EMS derived mutations in their target genes (Krasileva et al. 2017). Mutations within a target gene can be broadly categorised into four classes, based on the type of effect they are predicted to have on the protein encoded by the gene:

- (i) Premature stop – result in truncation of the protein, and so are predicted to have a strong effect on protein function. Mutations that truncate large regions of the protein are predicted to have the greatest effect on function.
- (ii) Splice mutation – may cause the RNA transcribed from the gene to be 'miss-spliced', leading either to a premature stop codon in the protein (and so result in a truncated protein) or to deletions/additions of amino acid residues within the protein.
- (iii) Missense mutation – result in a change of a single amino acid within the protein, which may have an effect on protein function depending on how critical the amino acid is to protein function.
- (iv) Sense mutation – while there is a single nucleotide change at the DNA level, no amino acid change is predicted at the protein level. This class of mutation is not commonly expected to have an effect on gene/protein function.

Generally, the aim when exploiting TILLING populations is to identify mutations most likely to have an adverse effect on protein function, and therefore have a strong effect on phenotype. Accordingly, type i or ii mutations are preferable to type iii mutations. In addition, for tetraploid wheat (where we expect two homoeologous copies of each gene, one on each of the A and B sub-genomes) it is preferable to identify mutations in both homoeologous copies of the target gene, i.e. from the A and B sub-genomes. Similarly in hexaploid wheat (where we expect three homoeologous copies of each gene), we aim to identify mutations in all three homoeologous copies of the target genes, i.e. from the A, B and D sub-genomes. As the presence of one functional homoeologue can buffer the effect of mutation in another homoeologue, to maximise the observable phenotypic effect of a given candidate gene, deleterious mutations at all homoeologues should ideally be combined into a single genetic background. For example, in the tetraploid variety Kronos, a line carrying a deleterious

mutation at the A homoeologue must be crossed with a line carrying a deleterious mutation at the B homoeologue to create an F1 individual (**Figure 7**). This F1 is then grown and selfed F2 seed collected, and F2 progeny selected that are homozygous for the mutations at both the A and B homoeologues, via the use of genetic markers. Subsequently, backcrossing combined with marker selection for the candidate gene homoeologues is normally undertaken to remove background mutations that may affect phenotype, followed by one round of selfing to recover line(s) homozygous for the target mutations, followed by selfing to generate bulked seed for subsequent phenotypic evaluation. Note: it is also possible to cross tetraploid and hexaploid wheat, so if for example no suitable mutation was identified in hexaploid Cadenza for the B homoeologue, a B homoeologue mutation from tetraploid Kronos could be crossed into a hexaploid background.



**Figure 7.** Example of a crossing scheme to generate a line containing homozygous TILLING mutant alleles at a target gene located on the A and B sub-genomes of the tetraploid wheat variety, Kronos. (1) TILLING lines carrying homozygous mutations at homoeologues on the A (a:a) and B (b:b) sub-genomes are identified. (2) The two TILLING lines are crossed to generate an F1 heterozygous for mutant alleles at both homoeologues (Aa:Bb), the F1 selfed, and an F2 individual homozygous for mutant alleles at both homoeologues (aa:bb) selected using genetic markers. (3) This individual is backcrossed to Kronos over two generations to remove background mutations. Genetic markers are used to ensure an individual carrying mutations at both homoeologues (Aa:Bb) is selected at the end of the second backcross. (4) This individual is selfed, and progeny homozygous for the target mutation at both homoeologues selected using genetic markers.

Here, we used DNA sequence from the coding regions (CDS) of 30 of our candidate genes to search the exome-capture sequenced TILLING populations of cv. Kronos and cv. Cadenza for lines carrying EMS induced mutations predicted to affect protein function. Considering both Kronos (A and B sub-genomes) and Cadenza (A, B and D), class i and ii mutations predicted to have an extreme impact on protein function were found for 45 of the estimated 90 homoeologues (i.e. 30 target candidate genes x 3 sub-genomes = 90 homoeologues) (**Table 8. Appendix 7**). Of these, highly deleterious mutations of all homoeologues within both Kronos (A and B sub-genome homoeologues) and Cadenza (A, B and D) were identified for 7 candidate genes: 5 from Kronos (genes 8, 13, 16, 22,

and 30) and 4 from Cadenza (genes 1, 3, 8, 16). It is notable that for 5 genes, no homoeologues were identified by BLASTn analyses on one of the sub-genomes (e.g. gene 4 appeared to lack the D homoeologue. **Table 8**). The lack of a homoeologue means that it is not necessary (or possible) to identify a TILLING mutant for that sub-genome. Taking these 5 instances into account, the number of genes with highly deleterious mutations on all homoeologues present in the genome rose from 7 to 9 (see gene 15 and 33, **Table 8**). Therefore, highly deleterious mutations for all of the homoeologues present, either within Kronos or within Cadenza, were identified for 30% of all the candidate genes investigated (9/30: genes 1, 3, 8, 13, 15, 16, 22, 30, 33). Of the remaining 21 candidate genes, an additional 12 genes were identified that lacked all but one of the homoeologues present in the genomes of either Kronos (i.e. 1 of the 2 homoeologues mutated/absent) or Cadenza (i.e. 2 of the 3 homoeologues mutated/absent) (genes 2, 6, 9, 10, 11, 12, 14, 19, 23, 27, 28 and 34). As gene 14 represented the well-characterised gene *Rht-B1*, this was excluded from further analysis. For the remaining 11 candidate genes, the presence of just one remaining putatively functional homoeologue increases the chance of observing phenotypic effect. Thus, the total number of candidate genes for which we predict we would be able to identify strong phenotypic effect, due to the mutation of all or all but one of the homoeologue in the genome, was 20 of the 30 genes investigated (67%).

TILLING lines containing highly deleterious mutations for 17 candidate genes (targeting 36 homoeologues via 33 TILLING lines) were ordered from the JIC SeedStor (**Table 8**). Of these, no seed was available for one line (for gene 35), and so was not progressed. Note: some TILLING lines contain deleterious mutations in two of the homoeologues we selected to target, so the number of TILLING lines ordered was less than the total number of homoeologues we targeted. Seed for the 32 TILLING lines received was sown, grown to maturity under glasshouse conditions, the developing ears selfed, seed harvested, and the bulked seed stored in the long term seed store for future use. Additionally, leaf tissue samples were also collected from each line, and genomic DNA extracted for future use. Subsequently, TILLING lines were also ordered for 6 genes for which severe TILLING mutations were identified in either all homoeologues, or in all but one of the homoeologues present (for genes 2, 6, 11, 23, 27 and 28). Thus, in total we have progressed TILLING germplasm resources for a total of 23 candidate genes.

It should be noted that as well as the highly deleterious mutations (class i and ii) discussed above, missense mutations (class iii) were also searched for across all candidate genes (**Appendix 4**). The effect of missense mutations is harder to predict, and so are normally exploited where no highly deleterious mutations are available. In this study, where highly deleterious mutations were either (i) not found at all in Kronos, or (ii) found for just one homoeologue in Cadenza, it is less likely we would be able to detect notable phenotypes using highly deleterious mutations alone, given the buffering effect of the remaining functional homoeologues. This was the case for 7 candidates: genes 2, 5,

20, 24, 25, 26, 31. In these cases, it could be beneficial to explore the available missense mutations. However, given the difficulty in predicting which missense mutations are likely to affect phenotype, and the work involved in crossing and phenotyping TILLING mutants, it would be more prudent to attempt such work in the tetraploid Kronos genetic background, rather than hexaploid Cadenza, as there are fewer numbers of homoeologues involved. All missense mutations identified for these 7 genes, and for all other candidate genes, are listed in **Appendix 4**, and available for future study. We note that even when considering missense mutations, three candidate genes still lacked mutations in all homoeologues: genes 2, 5 and 31. For these, if functional validation was required, alternative methods such as gene editing or RNA silencing, would need to be explored.

**Deliverables:** TILLING resources carrying highly deleterious mutations for either all homoeologues, or all but one homoeologue, were identified for 23 of the 30 candidate genes investigated, and seed and DNA resources developed. Additionally, TILLING resources for all but 3 of the remaining candidate genes were identified. Collectively, these resources provide the foundation with which future studies can functionally validate the candidate genes. Accordingly, these artificial mutants can potentially be used to either, (i) help identify the gene and natural genetic variants underlying the QTLs identified, or (ii) be used directly in breeding programmes.

Candidate gene <sup>1</sup>	A genome highly deleterious mutation	B genome highly deleterious mutation	D genome highly deleterious mutation	Kronos A+B highly deleterious mutation <sup>2</sup>	Cadanza A+B+D highly deleterious mutation <sup>3</sup>	Kronos 1 of 2 highly deleterious mutation	Cadanza 2 of 3 highly deleterious mutation	All homoeologues via Kronos + cadanza <sup>4</sup>	TILLING lines ordered <sup>5</sup>
1	C	C/K	C		Y	Y		Y	C-stopA, C-stopB, C-stopD
2	X	K	X			Y			<i>K-stopB</i>
3	C/K	C	C		Y	Y		Y	C-mis, C-stopB, C-stopD
4	X	C	No D				(Y)		C-stopB
5	X	X	X						
6	X	K	X			Y			<i>K-stopB</i>
7	ni	ni	ni						
8	C/K	C/K	C	Y	Y			Y	K-stopA, K-stopB
9	C/K	C	X			Y	Y		C-spliceA, C-stopB
10	X	K	X			Y			K-stopB
11	X	C/K	X			Y			<i>K-stopB</i> , <i>K-stopB</i>
12	X	C/K	X			Y			K-stopB
13	C/K	K	C	Y			Y	Y	C-stopA, K-stopB, C-stopD
14	X	C	C				Y		Not investigated ( <i>Rht-B1</i> )
15	No A	C/K	C	(Y)	(Y)			(Y)	C-stopB, C-stopD
16	C/K	C/K	C	Y	Y			Y	K-missA, C-missA, C-stopB, C-stopD
17	ni	ni	ni						
18	ni	ni	ni						
19	X	K	X			Y			K-stopB
20	X	X	X						
21	ni	ni	ni						
22	C/K	C/K	X	Y			Y		K-stopA, K-stopB
23	X	C/K	C			Y	Y		<i>K-stopB</i> , <i>K-spliceD</i>
24	X	C	X						C-stopB
25	X	C	X						
26	C	X	X						
27	C/K	X	X			Y			<i>K-stopB</i>
28	C	No B	X				(Y)		<i>C-stopA</i>
29	ni	ni	ni						
30	C/K	C/K	X	Y			Y		K-stopA, K-stopB, C-stopA, C-stopB, C-stopD
31	X	C	X						
32	ni	ni	ni						
33	C	No B	C		(Y)				C-stopA, C-stopD
34	C	No B	X				(Y)		C-stopA
35	K	X	X				Y		K-stopA <sup>†</sup> , K-misB, K-misB
36	ni	ni	ni						
37	ni	ni	ni						

**Table 8.** Summary of highly deleterious (i.e. premature stop or splice site mutations) in the candidate genes identified within yield/yield component QTL, considered at the sub-genome (A, B, D) level. 'C' =  $\geq 1$  highly deleterious mutations identified in Cadenza. 'K' =  $\geq 1$  highly deleterious mutations identified in Kronos. 'X' = no highly deleterious mutations identified. 'No A' or 'No B' or 'No D' = no homoeologue identified for the A, B or D sub-genome, respectively. 'ni' = not investigated. <sup>1</sup>Candidate gene numbering as listed in **Appendix 7**. <sup>2</sup>Highly deleterious mutations identified in all Kronos homoeologues. <sup>3</sup>Highly deleterious mutations identified in all Cadenza homoeologues. <sup>4</sup>All homoeologues mutated in either Kronos or Cadenza, when considering the possibility of using a mutated homoeologue from one species to the other (e.g. Using A or B highly deleterious mutation from Cadenza and crossing into Kronos). <sup>5</sup>Summary of TILLING lines ordered, in the format X-yZ, where X is represented by 'C' (Cadenza) or 'K' (Kronos), -y is represented by '-stop' (premature stop codon), '-splice' (splice site mutation) or '-mis' (missense mutation), and Z is represented by 'A', 'B' or 'D' (A, B or D sub-genome). TILLING lines listed in *italic* represent those most recently ordered (and for which bulked seed is not yet available). †No seed available from JIC SeedStor.

## 4.6 Genomic prediction

*Delivery of Milestones: M5.1, M5.2, M5.3, M5.4.*

**Aim:** Use MAGIC genotypic and phenotypic datasets to develop and validate Genomic Prediction (GP) methodologies targeting key traits.

**Process:** GP methods are commonly tested in collections of cultivars or advanced breeding lines from selection programmes. These often have strong population subdivision and variation in kinship relationships. The structure of MAGIC is different. Though diverse, its balanced crossing scheme ensures that each founder contributes uniformly to each line. Testing methods of prediction within MAGIC therefore compares prediction methods in the absence of strong kinship relationships. In addition, multi-parent crosses may prove to be a better origin for genomic selection schemes, so testing prediction methods in MAGIC has direct practical use. Prior to the start of this work package, we investigated various GP methods. We found ridge regression to perform best, but found Markov blanket approaches to suffer from problems with reproducibility. Accordingly, we used ridge regression implemented in the rrBLUP package v 4.6 (Endelman, 2011) to undertake GP using the phenotypic data for the traits grain yield (YLD), plant height (PHT), spikelet number (SKT) and specific weight (SPW), combined with the genotypic data generated using a 90k SNP array (Mackay et al. 2014; Gardner et al. 2016; NIAB unpublished). Results are summarised in **Appendix 8**.

As expected given the heritabilities and underlying genetic determinants identified in Sections 4.1 and 4.2, the best genomic prediction accuracies were achieved for PHT (0.69 to 0.79), followed by SKT (0.57 to 0.76), SPW (0.40 to 0.58) and YLD (0.12 to 0.48). For PHT the lowest prediction

accuracy was found for height data from the 2016 trial held at the RAGT site (RGT-YLD-16), though surprisingly this was not the least correlated site for the true phenotypic datasets (BAY-YLD-16). For SKT the lowest prediction accuracy was found for the ELS-YLD-16 data, in fitting with this site's low correlation against the true phenotypic training data. For SPW this same trend was observed for the ELS-YLD-16 data, being both the weakest correlation to the training set and also the poorest prediction accuracy. For YLD, we find that the KWS-YLD-16 trial gave consistently the most accurate predictions, regardless of training set used. Yet the KWS-YLD-15 trial had the weakest prediction based on a training set of pre-project yield data from NIAB (2012-2014). This suggests a strong environmental factor is influential on final yield values, therefore requiring a larger training data set to improve genomic prediction accuracy for yield relative to other traits. Overall we see that all genomic predictions are linearly related to true phenotypic data correlations, as would be anticipated. For grain yield, the best genomic predictions were achieved when combining phenotypic data from all four years of trial.

**Deliverables:** Establishment of GP workflows for the MAGIC population, and establishing prediction accuracies for key traits.

## 5. Discussion

### 5.1 Overview

A total of 376 QTL were identified for 18 traits. For a subset of 20 QTL, 58 co-dominant KASP markers were developed. These provide molecular tools to the participating breeding companies to explore the tracking of beneficial alleles segregating within ongoing breeding programmes. In addition to these markers, 31 NILs and TILLING mutants targeting 36 homoeologues from 17 candidate genes were developed. Together, these provide tools and resources to underpin future map-based gene cloning of QTL from UK-relevant germplasm. The resources for this subset of MAGIC QTL represents the first tranche of a pipeline established in this project that can provide continued staggered delivery of resources to underpin cloning of yield and yield component QTL of direct relevance to UK germplasm.

### 5.2 Genetic control of agronomic traits

Known major effect genes were identified as playing a role in controlling several of the traits investigated. As expected, the *RHT-B1* and *RHT-D1* semi-dwarfing loci that segregate in the population were found to have large effect on the control of plant height (explaining up to 11% and 28% of the phenotypic variation, respectively). In addition to plant height, both loci were found to have pleiotropic effects on  $\geq 11$  additional traits, illustrating the impact of semi-dwarfing alleles on crop performance. The photoperiod response gene *PPD-D1*, for which early flowering alleles originate from the MAGIC founder Soissons, had an expected major effect on flowering time (up to 31% phenotypic variation explained), as well as on plant height ( $\leq 5\%$ ), spikelet number ( $\leq 19\%$ ), and three grain traits (seed area, seed circumference and seed number per ear). Additionally, the major spikelet number QTL identified in the MAGIC population on chromosome 7A ( $-\log_{10}P = 66.4$ ) has previously been reported by numerous groups over the last  $\sim 15$  years (e.g. Quarrie et al. 2005; 2006). Our recent analysis of selection over the wheat pedigree indicates that this locus has likely been under strong breeder selection over the last  $\sim 15$  years (Fradgley et al. 2019), indicating breeder selection for this genetic locus of major effect. Of the predicted gene models found to lie within the 7A QTL, three candidate genes were identified. Interestingly, we also identified a homoeologous spikelet number QTL on chromosome 7B. Cross-referencing the gene content common to both the 7A and the 7B QTLs narrowed down the interval to a region containing just one of the three candidate genes. Analysis of the DNA sequences of this gene, using publicly available genomic sequence for Chinese Spring and the MAGIC founders Claire and Robigus, did not identify mutations predicted to lead to premature stop codons or miss-splicing of the CDS, although non-synonymous mutations were identified (**Appendix 9**). Further investigation, such as RNA expression analysis, population-based resequencing, and reverse genetics approaches such as the use of sequenced wheat TILLING populations (Krasileva et al. 2017), gene editing and/or gene silencing, can confirm the remaining candidate as the underlying gene.

In addition to these major effect genes, other MAGIC QTLs were identified which have previously been reported. For example, the homoeologous flowering time QTL we identified close to the long arm telomere on chromosomes 1B and 1D have been previously genetically mapped (Zikhali et al. 2014; 2015). Similarly, the MAGIC flowering time QTL located towards the short arm telomere of chromosome 7B corresponds to allelic variation at the *VERNALIZATION3* (*VRN3*) locus, previously shown to control flowering time in wheat (Yan et al. 2006; Bentley et al. 2014; Dixon et al. 2018). Additionally, we found the chromosomal region containing *VRN3*, 7B-1 (**Table 6**), to contain a QTL for spikelet number, indicating possible pleiotropic effect of allelic variation at *VRN3* on yield component traits - as is common for genetic loci controlling flowering time. Recent reports investigating grain size via artificial mutation of the wheat orthologue of the rice gene *GRAIN WIDTH2* (*GW2*) found *GW2* to control seed size and weight, as well as spike traits in wheat. These include mutation via gene editing (Wang et al. 2018; Zhang et al. 2018) and TILLING mutants (e.g. Simmonds et al. 2016). Additionally, natural variation at *TaGW2-6A* has been associated with modulation of grain weight and number (Zhai et al. 2018). *TaGW2-6A* is located on chromosome 6 at 237,759 Mbp (gene model TraesCS6A02G189300) within a region reported to contain extensive haplotype blocks due to its location within the low-recombining regions surrounding the centromeric region (e.g. reviewed by Brinton & Uauy, 2019). In the MAGIC population, *TaGW2-6A* is located close to the start of the highly non-recombining region spanning the 6A centromere, and lies within the MAGIC 6A-6 chromosomal location containing QTL for seed area, seed length, seed width, thousand grain weight and tiller number. Nevertheless, the relatively high levels of genetic recombination captured in the MAGIC population over three generations of intercrossing and subsequent selfing means that recombinations in the chromosomal region harbouring *TaGW2-6A* are present in our population. Further detailed analysis of our datasets is needed to determine whether *TaGW2-6A* represents a candidate gene for one or more of the grain size QTL at this location. Finally, the recently reported thousand grain weight chromosome 5A QTL *Qtgw-cb.5A* (Brinton et al. 2017) was identified in the MAGIC population within chromosomal region 5A-3 (**Table 6**). Using NILs for *Qtgw-cb.5A*, Brinton et al (2017) found the chromosomal region containing the QTL to affect grain weight via changes in grain length, along with pleiotropic effects on grain width. Increased grain length was found to be associated with longer maternal pericarp cell length (Brinton et al. 2017), with subsequent analyses of differential gene expression using grain tissue from NILs identifying differentially expressed genes, thus helping to refine candidate genes (Brinton et al. 2018). This QTL also lies within the highly non-recombining region, and so will likely be challenging to progress to the identification of the underlying gene.

Some genomic regions contained multiple QTLs. Of these, four were highlighted for possessing six or more QTL, were outside of the highly non-recombining regions spanning the centromere, and did not represent known genes of known effect such as *PPD-D1* and the *RHT* loci. These represent

priority loci for future focus, especially the 3A-1 region on chromosome 3A which includes a yield QTL, as well six grain trait QTLs.

### **5.3 Candidate genes and TILLING**

In Section 4.5, 53 yield and yield component QTLs defining 44 chromosome regions were prioritised for analysis of gene content within QTL intervals. These QTL were predominantly located within the more highly recombining regions of the wheat genome. In total, 41 candidate genes were identified. Considering the wheat grain size QTL alone, these included 7 wheat homologues of rice genes previously shown by map-based cloning to underlie grain size QTL in rice, including *GW2*. The candidate genes identified were used to search for TILLING mutants for future validation of gene function. The TILLING resources generated can be used in parallel with the NILs initiated within this project: where future fine mapping of Mendalised QTL via crossing a NIL pair is undertaken to generate recombinations within the region, and where a candidate gene for which we have developed TILLING resources remains within the fine-mapped region, investigation of the TILLING mutants may provide evidence of the identity of the gene underlying the QTL. Where candidate genes are found to fall outside of the fine-mapped region, the TILLING resources may provide novel, and phenotypically characterised, functional variation for potential inclusion within breeding programmes. Accordingly, the extensive NIL and TILLING resources for yield and yield component traits developed here provide valuable resources for future genetic improvement of agronomic traits in wheat.

### **5.4 Suggestions for future research**

Large volumes of genotypic, phenotypic, genetic and biological, and genetic data, resources and know-how have been generated, with the potential to further underpin multiple aspects of wheat genetic, molecular genetic and genomic-assisted R&D. The forward analysis process as demonstrated here for the 7A spikelet number QTL is applicable to the other QTL identified in this study. The simultaneous progression of large numbers (>10) of QTL to characterised NIL pairs, and ultimately to the underlying genes requires application of an analysis pipeline that integrates the reverse genetics approaches listed above alongside additional resources, including:

1. Genome resequencing data for all 8 MAGIC founders: the variety from which the wheat reference genome has been developed, Chinese Spring, is not particularly representative of UK wheat germplasm, and is thought to have been a selection from a Chinese landrace (Liu et al. 2018). Chinese Spring has become the global reference for wheat genome sequence, largely due to its early use in the development of genetic stocks, such as the aneuploidy germplasm developed by Earnie Sears in the 1930s (Sears, 1939). More recently, the reducing costs of genome sequencing combined with advances in genome assembly software and approaches, allows researchers to move beyond reliance on a single reference towards a 'pan-genome' era. For example, the 10+ Wheat

Genomes Project (<http://www.10wheatgenomes.com/>) is an international consortium re-sequencing 15 wheat varieties.

2. Gene expression and gene network resources: The utility of the genome sequence and gene model resources in Chinese Spring, and other varieties, is enhanced by additional resources, such as transcription atlases for gene expression (e.g. Borrill et al. 2016; Ramirez-Gonzalez et al. 2018). These provide information on gene expression in multiple tissues collected at different stages throughout development, often under contrasting environmental treatments. These resources can be integrated with the genomic sequence datasets for specific QTL regions, allowing genes to be prioritised based on their expression profiles. Additionally, genome-scale functional networks are now being developed in wheat (e.g. Lee et al. 2017). These incorporate gene expression data, and allow functional modules underlying complex traits to be determined. The network-based functional hypotheses to be generated can be cross referenced with the gene content of QTL, or with datasets generated from NILs generated from these QTL, further aiding the identification of candidate genes for functional validation.

Below, three broad future R&D opportunities are listed in more detail. The first two (5.4.1 and 5.4.2) aim to identify the genes/genetic polymorphisms underlying QTL detected in this study, precisely quantify their effects, and provide perfect markers assaying underlying polymorphisms. They benefit from the speed and throughput afforded the resources generated in this project, and in other ongoing projects on wheat genomics, and the two approaches would likely be exploited in combination with each other. We plan to develop an academic-industrial project proposal targeting these aspects in the coming year. The third approach (5.4.3) uses the phenotypic data generated in this project, in order to develop approaches and resources for potential use in the further development of approaches for hybrid wheat breeding. NIAB, along with the Scotland's Rural College (SRUC) and industrial partners KWS, Limagrain, RAGT and Asur have recently submitted a BBSRC proposal addressing this topic.

#### ***5.4.1 Leveraging emerging MAGIC genome re-sequencing and RNAseq data to prioritise candidate polymorphisms within QTL intervals***

In a separate recently funded BBSRC project, NIAB, in collaboration with Earlham Institute, the John Innes Center and the Natural History Museum, is re-sequencing the genomes of the eight MAGIC founders (BBSRC projects BB/P010741/1, BB/P010733/1 and BB/P010768/1). This, combined with the recently released wheat reference genome for cv. Chinese Spring (IWGSC, 2018) will allow future investigation to establish a comprehensive catalogue of genic polymorphisms present in the genetic intervals of all QTL identified. Subsequent cross-comparison of this catalogue with gene expression atlases generated from public databases, as well as from each of the eight founders, will allow rapid prioritisation of candidate genes/ polymorphisms within QTL intervals. A multi-tissue

RNASeq catalogue for two of the MAGIC founders, Claire and Robigus, is already underway as part of the 10+ Wheat Genomes Project, (<http://www.10wheatgenomes.com/>), and we plan to supplement this with similar resources for the remaining six MAGIC founders. The benefit of this approach is that it is timely (exploiting resources that are only now becoming available) and simultaneously targets all of the QTL identified in within this report.

#### **5.4.2. Detailed understanding of Mendalised QTL**

The identification of NILs for multiple yield and yield component traits provides a potentially rapid route to (i) quantifying the phenotypic effect of Mendalised QTL in near isogenic backgrounds (at the organ, tissue and cellular level), and (ii) fine-mapping Mendalised QTL to single gene resolution - achieved by generation of large numbers of F2 progeny derived from crossing a given NIL pair together, and using genetic markers to identify lines carrying genetic recombinations within the QTL region. Candidate genes within refined QTL intervals can then be updated, and functionally validated via the TILLING resources initiated in this report. The benefit of this approach is that it targets specific QTL for which NIL germplasm has been generated in this project.

In practice, approaches 5.4.1 and 5.4.2 are complementary, and would be undertaken simultaneously. The ultimate aim would be to identify the genes and genetic polymorphisms underlying genetic variation currently in use in elite wheat germplasm. The approach here is notable in that it simultaneously targets large numbers of QTL for map-based cloning. Such parallelised approaches to map-based gene cloning have the potential to rapidly advance our understanding and exploitation of genetic variation controlling yield. Such an approach is now possible due to the resources and expertise generated in this project, and their alignment with emerging genomics and gene expression resources in wheat, as well as with ongoing work in other wheat research groups in the UK and beyond. The ultimate goal of such an approach would be the development of 'perfect' genetic markers that allow precise tracking within breeding programmes of the genetic variant associated with conferring beneficial expression of each target trait.

#### **5.4.3. Development of hybrid wheat approaches and resources**

The extensive yield and yield component phenotypic data generated in this project using the MAGIC population provide resources with which to develop new hybrid wheat breeding strategies and resources. Hybrid crop production has been one of the critical technological developments in modern crop breeding, providing rapid gains in yield. While wheat is one of the most important crops globally (and the most important in the UK), it is almost exclusively grown as inbred varieties. Switching to F1 wheat hybrid production would rapidly improve sustainable wheat production: annual genetic improvement in UK wheat yield due to breeding is ~1%; the average yield advantage of wheat hybrids over their parents is ~10% (equating to ~10 years of inbreeding-based varietal development). This gain translates to about £160 million p.a. at the farm gate in the UK alone. Current hybrid wheat

breeding practice follows the successful model of maize. However, the genetic determinants of yield in wheat, an inbreeding species, interact in different ways to those in maize and other outbreeding cereals. Simply mimicking maize breeding strategies will not be optimal. Furthermore, wheat hybrid seed production requires a means to render plants male-sterile at scale. Systems used in other crops, such as restoration of cytoplasmic male sterility, are not successful in wheat, largely due to the complexity of the hexaploid genome. In Europe, commercial hybrid wheat production is still reliant on chemical gametocide application. Because the increased yield of wheat hybrids has not historically compensated for the increased cost of hybrid seed production, development/uptake of hybrid wheat has not been widely successful to date. The introduction of hybrid wheat varieties therefore has two challenges:

(1) Reliable hybrid seed production systems must be established: these are being developed in the private sector and are at the point where some wheat hybrid varieties have recently been commercialised, most notably in France and Germany.

(2) The need for refinement of wheat-specific hybrid breeding approaches.

Challenge 1 is now being met via recent availability of efficient chemical gametocides, and improved wheat genomic resources have renewed private research into genetic sterility/restoration systems. There is now wide industrial interest in determining improved hybrid wheat breeding methods. To address this need, NIAB, in collaboration SRUC and industrial partners KWS, Limagrain, RAGT and ASUR have recently submitted a proposal titled '*HyBreed: exploring heterosis, transgressive segregation and epigenetic inheritance to underpin and develop efficient wheat hybrid breeding approaches*' under the BBSRC Industrial Partnership Award (IPA) scheme. Submitted in January 2019, decisions on project funding are expected in autumn 2019.

## 6. References

- Bentley AR, Scutari M, Gosman N, Faure S, Bedford F, et al. (2014). Association mapping and genomic selection for genetic dissection of key traits in elite European wheat. *Theoretical and Applied Genetics*, 127: 2619-2633.
- Bernardo R (2008). Molecular markers and selection for complex traits in plants: learning from the last 20 years. *Crop Science*, 48: 1649-64.
- Boeven PHG, Longin FH, Leiser WL, Kollers S, Ebmeyer E, Würschum T (2016). Genetic architecture of male floral traits required for hybrid wheat breeding. *Theoretical and Applied Genetics*, 129: 2343-2357.
- Borrill P, Ramirez-Gonzalez R, Uauy C (2016). expVIP: a customizable RNA-seq data analysis and visualization platform. *Plant Physiology*, 170: 2172-2186.
- Brinton J, Simmonds J, Minter F, Leverington-Waite M, Snape J, Uauy C (2017). Increased pericarp length underlies a major quantitative trait locus for grain weight in hexaploid wheat. *New Phytologist*, 215: 1026-1038.
- Brinton J, Simmonds J, Uauy C (2018). Ubiquitin-related genes are differentially expressed in isogenic lines contrasting for pericarp cell size and grain weight in hexaploid wheat. *BMC Plant Biology*, 18: 22.
- Brinton J, Uauy C (2019). A reductionist approach to dissecting grain weight and yield in wheat. *Journal of Integrative Plant Biology*, doi: 10.1111/jipb.12741.
- Broman KW, Wu H, Churchill GA (2003). R/qtl: QTL mapping in experimental crosses. *Bioinformatics*, 19: 889–890.
- Cavanagh C, Morell M, Mackay I, Powell I (2007). From mutations to MAGIC: resources for gene discovery, validation and delivery in crop plants. *Current Opinions in Plant Biology*, 11: 1-7.
- Dixon LE, Farré A, Finnegan EJ, Orford S, Griffiths S, Boden SA (2018). Developmental responses of bread wheat to changes in ambient temperature following detection of a locus that includes *FLOWERING LOCUS T1*. *Plant Cell Environment*, 41: 1715-1725.

Endelman JB (2011). Ridge regression and other kernels for genomic selection with R package rrBLUP. *Plant Genome*, 4: 250-255.

Fradgley N, Gardner KA, Cockram J, Elderfield J, Hickey JM et al. (2019). A large-scale pedigree resource of wheat reveals evidence for adaptation and selection by breeders. *PLoS Biology*, e-publication ahead of print. <https://doi.org/10.1371/journal.pbio.3000071>.

Fulton TM, Chunwongse T, Tanksley S D (1995). Microprep protocol for extraction of DNA from tomato and other herbaceous plants. *Plant Molecular Biology Reports*, 13: 207–209.

Gardner KA, Wittern LM, Mackay IJ (2016). A highly recombined, high-density, eight-founder wheat MAGIC map reveals extensive segregation distortion and genomic locations of introgression segments. *Plant Biotechnology Journal*, 14: 1406-1417.

Gegas VC, Nazari A, Griffiths S, Simmonds J, Fish L et al. (2010). A genetic framework for grain size and shape variation in wheat. *Plant Cell*, 22: 1046-1056.

Griffiths S, Simmonds J, Leverington M, Wang Y, Fish L et al. (2009a). Meta-QTL analysis of the genetic control of ear emergence in elite European winter wheat germplasm. *Theoretical and Applied Genetics*, 119: 383-395.

Griffiths S, Simmonds J, *Leverington* M, Wang Y, Fish L et al. (2009b). Meta-QTL analysis of crop height in elite European winter wheat germplasm. *Molecular Breeding*, 29:159-171.

Huang BE, George AW (2011). R/mpMap: a computational platform for the genetic analysis of multiparent recombinant inbred lines. *Bioinformatics*, 27: 727–729.

Huang BE, George AW, Forrest KL, Kilian A, Hayden MJ et al. (2012). A multiparent advanced generation inter-cross population for genetic analysis in wheat. *Plant Biotechnology Journal*, 10: 826-839.

The International Wheat Genome Sequencing Consortium (IWGSC) et al. (2018). Shifting the limits in wheat research and breeding using a fully annotated reference genome. *Science*, 361: 6403.

Krasileva KV, Vasquez-Gross HA, Howell T, Bailey P, Paraiso F, et al. (2017). Uncovering hidden variation in polyploid wheat. *Proc Natl Acad Sci USA*, 114: E913-E921.

Lee T, Hwang S, Kim CY, Shim H, Kim H et al (2017). WheatNet: a genome-scale functional network for hexaploid bread wheat, *Triticum aestivum*. *Molecular Plant*, 10, 1133-1136.

Li Y, Fan C, Xing Y, Jiang Y, Luo L, Sun L, Shao D, Xu C, Li X, Xiao J, He Y, Zhang Q (2011). Natural variation at *GS5* plays an important role in regulating grain size and yield in rice. *Nature Genetics*, 43: 1266-1269.

Liu D, Zhang L, Hao M, Ning S, Yuan Z, et al. (2018). Wheat breeding in the hometown of Chinese Spring. *The Crop Journal*, 6: 82-90.

Mackay I, Bansept-Basler P, Barber T, Bentley AR, Cockram J, Gosman N, et al. (2014). An eight-parent multiparent advanced generation intercross population for winter-sown wheat: creation, properties and validation. *G3 Genes Genomes Genetics*, 4: 1603-1610.

Mackay I, Powell W (2007). Methods for linkage disequilibrium mapping in crops. *Trends in Plant Sciences*, 12: 57-63.

Mott R, Talbot CJ, Turri MG, Collins AC, Flint J (2000). A method for fine mapping quantitative trait loci in outbred animal stocks. *Proc Natl Acad Sci USA*, 97: 12649-12654.

Quarrie SA, Steed A, Calestani C, Semikhodskii A, Lebreton C, Chinoy C, et al. (2005). A high-density genetic map of hexaploid wheat (*Triticum aestivum*) from the cross Chinese Spring x SQ1 and its use to compare QTLs for grain yield across a wide range of environments. *Theoretical and Applied Genetics*, 110: 865-880.

Quarrie SA, Quarrie SP, Radosevic R, Rancic D, Kaminska A, Barnes JD, et al. (2006). Dissecting a wheat QTL for yield present in a range of environments: from the QTL to candidate genes. *Journal of Experimental Botany*, 57: 2627-2637.

R Core Development Team, 2013. R: A language and environment for statistical computing

Ramirez-Gonzalez RH, Borrill P, Lang D, Harrington SA, Brinton J et al. (2018). The transcriptional landscape of polyploidy wheat. *Science*, 361: 6403.

Ramirez-Gonzalez RH, Uauy C, Caccamo M (2015). PolyMarker: a fast polyploid primer design pipeline. *Bioinformatics*, 31: 2038–2039.

Rustgi S, Shafqat MN, Kumar N, Baenziger S, Ali ML et al. (2013). Genetic dissection of yield and

its component traits using a high-density composite map of wheat chromosome 3A: bridging gaps between QTLs and underlying genes. PLoS One, e70526.

Sears ER (1939). Cytogenetic studies with polyploid species of wheat. I. Chromosomal aberrations in the progeny of a haploid of *Triticum vulgare*. Genetics, 24: 509-523.

Simmonds J, Scott P, Brinton J, Mestre TC, Bush M et al. (2016). A splice acceptor site mutation in *TaGW2-A1* increases thousand grain weight in tetraploid and hexaploid wheat through wider and longer grains. Theoretical and Applied Genetics, 129: 1099-1112.

Wang W, Pan Q, He F, Akhunova A, Chao S et al. (2018). Transgenerational CRISPR-Cas9 activity facilitates multiplex gene editing in allopolyploid wheat. CRISPR Journal, 1: 65-74.

Wang, S., Wong, D., Forrest, K., Allen, A., Chao, S., Huang, B. E., et al. (2014). Characterization of polyploid wheat genomic diversity using a high-density 90,000 single nucleotide polymorphism array. Plant Biotechnology. Journal, 12: 787–796.

Xue S, Xu F, Li G, Zhou Y, Lin M et al (2013). Fine mapping *TaFLW1*, a major QTL controlling flag leaf width in bread wheat (*Triticum aestivum* L.). Theoretical and Applied Genetics, 126: 1941-1949

Yan L, Fu D, Li C, Blechl A, Tranquilli G, et al. (2006). The wheat and barley vernalization gene *VRN3* is an orthologue of *FT*. Proc Natl Acad Sci USA, 103: 19581–19586.

Zhai H, Feng Z, Du X, Song Y, Liu X et al. (2018). A novel allele of *TaTGW2-6A1* is located in a finely mapped QTL that increases grain weight but decreases grain number in wheat (*Triticum aestivum* L.). Theoretical and Applied Genetics, 131: 539-553.

Zhang Y, Li D, Zhang D, Zhao X, Cao X, Dong L, Liu J, Chen K, Zhang H, Gao C, Wang D (2018). Analysis of the functions of *TaGW2* homoeologs in wheat grain weight and protein content traits. The Plant Journal, 94: 857-866

Zikhali M, Griffiths S (2015). The Effect of Earliness *per se* (*Eps*) Genes on Flowering Time in Bread Wheat. In: Ogihara Y., Takumi S., Handa H. (eds), Advances in Wheat Genetics: From Genome to Field. Springer, Tokyo.

Zikhali M, Leverington-Waite M, Fish L, Simmonds J, Orford S, et al. (2014). Validation of 1DL earliness *per se* (*eps*) flowering QTL in bread wheat (*Triticum aestivum*). Molecular Breeding, 34: 1023-1033.

## 7. Workplan

		2014	2015			2016				2017				2018
		Q-2	Q-1	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8	Q9	Q10	Q11
<b>WP1</b>	<b>MAGIC data management</b>													
1.1	MAGIC germplasm management	M1.1			M1.2									
1.2	Year 2 trial seed processed					M1.3								
						M1.4								
<b>WP2</b>	<b>Yield trials and phenotyping</b>													
	Year 1 trials	M2.1				M2.2								
	Year 2 trials					M2.3				M2.4				
<b>WP3</b>	<b>QTL mapping</b>													
	QTL mapping							M3.2			M3.3			
	Meta-QTL mapping						M3.1					M3.4		
<b>WP4</b>	<b>Marker conversion, fine mapping</b>													
4.1	Breeder friendly markers for yield								M4.1			M4.2		
4.2	Further investigation of selected QTL					M4.4								M4.5
<b>WP5</b>	<b>Genomic Prediction and Selection</b>													
	Genomic models								M5.1			M5.2	M5.3	M5.4
<b>Management</b>														
	Project meetings		PM		PM		PM		PM		PM		PM	PM

**Q-1:** Undertaken before project start.

**Milestones:** **M1.1** Seed processing for year 1 trials complete. **M1.2** Seed processing for year 2 trials complete. **M1.3** Databasing of 90k genotype data for 1,000 lines. **M1.4** Databasing of historic phenotype data. **M2.1** Year 1 trials sown. **M2.2** Year 1 phenotype data collated. **M2.3** Year 2 trials sown. **M2.4** Year 2 phenotype data collated. **M3.1** QTL and meta-QTL analysis of historical phenotype data. **M3.2** QTL analysis of year 1 data. **M3.3** QTL analysis of year 2 data. **M3.4** Meta QTL analysis of all phenotype data. **M4.1** Conversion of selected 90k SNPs to KASP – round 1. **M4.2** Conversion of selected 90k SNPs to KASP – round 2. **M4.3** Development of additional KASP markers within target QTL intervals. **M4.4** Identification of NILs. **M4.5** Genotypic evaluation of NILs and seed bulking. **M5.1** Initial GS model developed. **M5.2** Bayesian Network and additional approaches tested. **M5.3** Final models developed. **M5.4** Markov blanket derived KASP markers. **D5.1** Resources (KASP markers and MAGIC line subsets) for GS-assisted selection for yield/yield stability in MAGIC.

## 8. Dissemination and project outputs to date

<b>Dissemination</b>	<b>Date</b>	<b>Location</b>
Field demonstration plots of MAGIC germplasm	10-11 June 2015	Cereals 2015
Field demonstration plots of MAGIC germplasm	23 June 2015	NIAB Open Day
Field demonstration plots of MAGIC germplasm	26 June 2015	NIAB Director's Day
Invited Talk, WGIN Network Meeting	20 Nov 2015	RRES, UK
Field demonstration plots of MAGIC germplasm	15-16 June 2016	Cereals 2016
Field demonstration plots of MAGIC germplasm	June 2016	NIAB Open Day
Field demonstration plots of MAGIC germplasm	June 2016	NIAB Open Day
AHDB 6 month undergrad studentship (Yeorgia Argirou)	Sept 2016 – Feb 2017	NIAB
Research presentation	10 April 2017	EMR
Invited Talk	10 Nov 2017	Gregor Mendel Institute, Austria
Article for CPM magazine	May 2017	CPM magazine
Talk to Italian society, non-specialist audience	Oct 2017	Trieste, Italy
Poster presentation	13-17 Jan 2018	PAG, San Diego, USA
<b>Post-project:</b>		
BBSRC IPA grant proposal, submitted	Jan 2019	Collaborators: NIAB, SRUC, EI, Asur, KWS, Limagrain, RAGT

## 9. Acknowledgements

This project was funded by AHDB and the Biotechnology and Biological Sciences Research Council (BBSRC) grant BB/M008908/1, and in-kind contributions from BASF, Elsoms Seeds Ltd, KWS UK Ltd, Limagrain UK Ltd and RAGT Seeds Ltd.



*Note: at the time of writing, the data generated within this project was confidential under the terms of the BBSRC approved consortium agreement, and as signed by all academic and industrial project partners. This has placed restrictions on the details of the information presented in the main body of this report. Accordingly, the data contained in the Appendices is embargoed until 28<sup>th</sup> Feb 2023, unless permissions are gained from the industrial partners to publish beforehand.*