

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 3 |
| 1.1 | Wheat Breeding | 4 |
| 1.2 | Wheat Genetics | 4 |
| 1.3 | Wheat Genomics | 4 |
| 1.4 | Sequencing | 4 |
| 1.5 | Sequence analysis | 8 |
| 1.5.1 | Ambiguity Codes | 9 |
| 1.5.2 | RNA-Seq | 9 |
| 1.6 | Wheat specific resources resources | 10 |
| 2 | PolyMarker: A fast polyploid primer design pipeline | 11 |
| 2.1 | Pipeline | 11 |
| 2.2 | Applications of PolyMarker | 15 |
| 2.2.1 | KASP assays for public sets of SNPs | 16 |
| 2.2.2 | SNPs in a mutant population | 16 |
| 2.2.3 | Deletions on a mutant population | 16 |
| 2.2.4 | PolyMarker public web service | 17 |
| 2.2.5 | Genotyping of <i>Puccinia striiformis</i> f. sp. <i>tritici</i> isolates. | 18 |
| 2.3 | Discussion | 18 |
| 3 | Genetic map of <i>Yr15</i> with RNA-Seq | 19 |
| 3.1 | Mapping population | 23 |
| 3.2 | Sequencing and mapping | 25 |
| 3.3 | SNP Calling | 27 |
| 3.4 | Bulk Frequency Ratios | 30 |
| 3.5 | <i>In silico</i> mapping | 33 |
| 3.6 | Assay selection | 36 |
| 3.7 | Genetic map | 41 |

| | |
|---|-----------|
| <i>CONTENTS</i> | 2 |
| 3.8 Bioinformatic methods | 41 |
| 3.8.1 Alignment reads to gene models | 41 |
| 3.8.2 Bulk Frequency Ratios | 41 |
| 3.8.3 Alignment between gene models | 42 |
| 3.9 Discussion | 42 |
| 4 Gene expression (expVIP) | 44 |
| 4.1 Expression experiments (Introduction) | 44 |
| 4.2 Database design | 44 |
| 4.3 Analysis pipeline | 44 |
| 4.4 Graphical interface | 44 |
| 4.5 Conclusions | 44 |
| 5 Conclusions and final remarks | 45 |
| A Supplemental tables | 46 |
| A.1 PolyMarker supplemental tables. | 47 |
| B Quality control | 57 |
| B.1 Sequence read quality | 58 |
| B.2 Sequence GC content | 60 |

Chapter 1

Introduction

It defines the objectives and the importance of the research. It focus on the the application of Next Generation Sequencing to molecular biology, wheat genetics and ultimately to breeding programs. It also mentions the current status of the wheat reference genome and other resources (genetic maps, markers) the need of tools to query them effectively.

1.1 Wheat Breeding

An overview of how breeding is carried on currently, the different sources of genetic diversity and the relevance of fixing agriculturally important traits.

1.2 Wheat Genetics

The section describes alleles and the concept of gene, both as a locus in the genome (Quantitative Trait Locus, QTL) and an specific transcript (central dogma of molecular biology). Finally, it discusses traditional Mendelian inheritance and the effect of polyploidy. Some of this is described in the Yr15 chapter, maybe it is not needed anymore here.

1.3 Wheat Genomics

A description of the current status of the wheat genome (Mayer et al. (2014), Chapman et al. (2015)), the different available assemblies and approaches to sort the scaffolds (Genome Zipper, the various genetic maps).

1.4 Sequencing

The Human Genome Project used Sanger sequencing Lander et al. (2001). This technology is the current gold standard in terms of quality of the sequence. It evolved from electrophoresis gels where the bands represented bases to a fully automated technique. However, the throughput is limited and doing genome wide analysis has prohibitive costs. In the second half of the 2000s high-throughput sequencing technologies emerged which had reduced the cost of sequencing. The main principle of the second-generation sequencing is to produce clusters of clones (i.e. ePCR), fix them in a plate and then add bases with a fluorescent marker. The reaction happens in parallel in millions of clusters at the same time. With each cycle, a picture is taken, showing the fluorescence of each base. Then, image processing algorithms find where in the image the clusters are and the bases are called. At this scale, the volume and complexity of the information is not trivial to manipulate, hence computing is required.

According to the objectives of the experiment and the quality and volume of the available DNA, the library can be prepared on fragments of different sizes, the classification of the available sequencing for the fragments is the following Myllykangas et al. (2012); Metzker (2010); Shendure and Ji (2008); Hutchison (2007):

Single end When the fragments are short, it is possible to just sequence from the 5'-end the read.

Read Pairs When the sample consists fragments of up to 500bp, it is possible to read the 5' end up to the read length where the quality starts to drop, the molecule can be turned upside down, reverse complemented and sequence backwards. It is not required, but ideally, the fragments sequenced with read pairs should be selected to have a homogeneous size. The reads are in opposite orientation relative to each other.

Overlapping Read Pairs are a variation to read pairs, where the size of the fragment is shorter than two times the read length. This allows an alignment between the two fragments to get a longer read with the limitations of the instrument.

Mate pairs are used to get reads separated at distances between 1kbp and 5kbp. To achieve this, the molecule is circularised and the point where the two ends of the fragment were joined a biotin marker is inserted. Then, the molecule is fragmented again and the fragments containing the biotin are sequenced in the same fashion as read pairs. The resulting reads have the same orientation.

There are several types of experiments that can be analysed with high throughput sequencing, accordingly, different protocols for the sample preparation exist. The following is a short list of some of them

Whole genome shotgun When a sample is prepared for WGS, the DNA is extracted and chopped in fragments and sequenced. The reads obtained are, in principle, randomly distributed across the whole genome

RNA-Seq . Instead of sequencing DNA, mRNA is captured and sequenced. The fragments are not amplified in any way, to enable a portrait of the gene expression levels.

ChIP-SEQ . Chromatin Immunoprecipitation is used to find relationships between proteins and DNA sequence. It is useful to find transcription factors and replication-related proteins.

Amplicon sequencing . Used primarily to do barcoding of species. A known gene is amplified (i.e. 16S) with the intention of characterising the species present in the sample.

Metagenomic capture From a mixed sample (soil, root, animal fluids) all the DNA is extracted and sequenced, this gives a snapshot of the microbial community in the sample

RAD-seq Restriction site associated DNA markers are useful to do population analysis. The technique focuses on sequencing regions around restriction sites and the variations around them can be used to genotype individuals.

Exon capture The DNA is extracted and baits are used to attract the regions with motifs common around exons. This allows to sequence only the genes and regions near them.

The different sequencing technologies available as of 2013 have different yields, advantages and disadvantages, as described below:

Illumina Each fragment is amplified using bridge amplification over and over in the same place in the plate to form clusters. After the clusters are formed, a last cycle of amplification is carried on with the bases being added to the template, with the intervention of a polymerase, have a fluorescent marker which makes the cluster glow depending on the added base. It adds one base per cycle. With a read length between 75bp and 250bp is currently the most widely adopted platform. As a de facto standard, many tools exist to cope bioinformatically with the biases of the machine. The run takes 4 or 9 days, depending on days, depending if one or two reads are generated for each fragment. It produces up to 35 gigabases per run.

SOLiD The preparation of the fragments is similar to Illumina, however, when adding the bases they are added in pairs. This technique is called sequencing by ligation as it uses a DNA Ligase, as opposed

to a polymerase, to determine the transition between bases. The resulting sequence is not in base space, but in colour space, which represents the transition state between bases. This technique is robust for finding SNPs when you have a good reference where to align the reads. However, the number of tools available and the research done to analyse sequences in colour space is low compared to the tools using base space. The runs take between one and two weeks to complete, with a yield of up to 50 gig abases per run. The read length can be up to 50 bases

Roche/454 The fragments are cloned in beads, which then fall in wells in the slide. The sequencing is done by adding nucleotides in a determined order. The next nucleotides to be added in the reaction contain a fluorescent marker. The bases are not added one by one, but all the bases that are the same are added together. The amount of glow on each well can tell how many times a base is added. As the glow is not a discrete number, when a long homopolymer appear (above 5 bases) the likelihood of having a wrong count of the homopolymer is increased. The average read length varies between 300 and 700bp. A run usually takes half a day, but it only yelds 0.45 gigabases. The cost of the reagents is relatively expensive, but if the experiment requires longer reads it is a good option.

PacBio Opposed to all the previous technologies, Pacific Biosciences has developed a sequencing technology where the molecules doesn't need to be PCR amplified before the sequencing. The glass slide used contains wells with a depth of 100nm where a polymerase lays at the bottom. The nucleotides to be added have a fluorescent marked that is freed when the polymerase adds the nucleotide, releasing a light signal, which then can be captured from the bottom of the glass. The error rate for this technology is still high (about 10% of the bases are miscalled), however reading several times the same molecule reduce the error rate. The main advantage is that the reads can be over 1kbp.

OpGen Additionally, high-througput optical mapping technologies, like OpGen, are becoming accessible. The maps are done by fixing single molecules of DNA are held on a slide. Then, restriction enzymes targeted to specific digestion sites cut the fragment and fluorescent

markers are added to the ends of the fragments. Finally, the fragments are visualised and the size of the molecules is measured by the distance between fluorescent points in the slide. This is done with several fragments at the same time. Then, the distances between restriction sizes can be compared across all the fragments to generate a consensus. Finally, if you have contigs from other technologies, it is possible to complement the information and get better assemblies. Even without the contigs, the data can be used to compare translocations within strains of different bacteria or homologous species at a chromosome level.

ION Torrent (Do some research on newer sequencing things)

1.5 Sequence analysis

This section discusses the criteria to decide analysis done after sequencing, when to do re-alignments or *de novo* assemblies, how to do SNP calling in diploid and polyploid organisms and the bulk frequency ratios.

DNA sequence alone is not enough to understand the biology behind, a context is required. There are databases like Ensembl and NCBI that act as repositories of the known public sequences.

From the computational point of view, the problem can be viewed as a string matching. The Smith-Waterman Smith and Waterman (1981) and Needleman-Wunsch Needleman and Wunsch (1970) algorithms are the gold standard interns of accuracy looking for similarity between sequences. However, the execution time for both of them is prohibitive to run in massive databases. The algorithm execution time is $O(mn)$, as it requires calculating a matrix of size mn where m is the target sequence and n is the query sequence. To scale this to a manageable problem algorithms like BLAST index the references and use heuristics to make the search more manageable, with some penalty in the accuracy. This alignments tools are useful for long stretches of DNA (like cDNA or contigs) Altschul et al. (1990).

TODO: List of global aligners -BLAST -BLAT -Exonerate -nucmer

When looking at a protein level, where the sequences may be only loosely similar, Hidden Markov Models (HMM) are used to search for

protein families. This can be useful to annotate putative proteins and their functions. HMMs require a training dataset, where proteins are previously annotated and the reference is a model encoding the characteristics of a family, with associated probabilities. Hence, this technique is something between a sequences aligner and a classifier Eddy (2004).

When analysing high-throughput sequencing, having millions of short sequences make unfeasible to try to align the data to every possible reference. However, one can take in advantage the fact that you know which organism you are looking for and, if available, use a genomic reference. For this, tools like MAQ, BWA, Bowtie, among others, provide indexed search. Once you have your reads aligned to a reference you can do more analysis, depending on the biological question being asked and the type of sequencing carried on. Fortunately, most of the Short-Read sequence alignment produce similar outputs and the SAM format is becoming a de facto standard. This is allowing to make more modularised downstream analysis where you can test different aligners with different settings and pick the algorithm that better fits your experiment Liu and Schmidt (2012); Li and Durbin (2009); Li et al. (2009).

1.5.1 Ambiguity Codes

Make a table with the ambiguity codes and why they are useful.

1.5.2 RNA-Seq

One way to narrow down which genes are involved in certain trait or response to the environment is to focus on studying only the expressed genes. One of the techniques involving high-throughput sequencing is RNA-Seq. This technique captures the messenger RNA in the tissue being studied and sequenced. The premise is that you will find a gene more expressed if it is being used by the organism. Some proteins with a vital role for the cell are always expressed (i.e. RuBisCO for carbon fixation in plants GM (2000)). On the simplest of the experiments you would need two datasets to compare, one with the gene being looked expressed and one where it is not. The expression can come from different environmental conditions, development stage or different genotypes. Mortazavi et al. (2008)

Depending on how much *a priori* information of the analysed organism is available different bioinformatic approaches can be used.

Transcriptome alignment The reads are aligned to a database of known cDNA. Ideally, alternative splicing sequences are available, so a simple alignment should work (i.e. BWA, bowtie).

Genomic alignment The reads are aligned to the genome. The splice junctions, introns and axons need to be accounted, so simple alignment doesn't work. Regular alignments are used, but the reads may be trimmed at fixed sizes to allow discontinuous alignments using regular tools (i.e. Stampy, tophat/cufflinkns)

De Novo transcriptome assembly If a reference of the organism is not available, it is possible to generate a draft transcriptome with the RNA-Seq reads with traditional assemblers (velvet, abyss) or with specialised assembler tools like Trinity.

Once you have the alignments it is possible to evaluate the relative expression of the genes in the sample calculating the Reads per Kilobase per Million mapped reads (RPKM) or the Transcripts per Million (TPM). This normalises the expression by the amount of sequenced data and can be used to find which genes change in expression volume across different samples.

1.6 Wheat specific resources resources

Gene models -UniGene -UCW Gene models -Gene annotation IWGSC
-Gene annotation TGACv1

Genetic maps -Wang -Chapman/PopSeq (is the same populaiton, improved)

Markers -90k -820k -MASwheat/SRR

Portal -CeralsDB -MASWheat -Ensembl -Wheat-expression

Assemblies -Chapman -IWGSC -TGACv1 -NRGene (unpublished?)
-454 Liverpool

A compilation of the currently available resource for whet genetics and genomics. MAS wheat, CeralsDB, Ensembl, etc.

Chapter 2

PolyMarker: A fast polyploid primer design pipeline

One of the main challenges of working with polyploid species is the design of genome specific molecular markers. This is particularly true when targeting conserved homoeologue regions, where a primer could bind to a pair, or triplet, of identical sequences. For that reason, designing primers for polyploids require to include bases that are specific to the target, in addition to the physicochemical properties of the primer. The traditional methodology to find primer candidates include a blast search and a local alignment, select the primer candidates manually, and finally, validate the primers with a tool, like **Primer3** (Rozen and Skaletsky, 2000). To reduce the time invested in designed primers I have developed PolyMarker (Ramirez-Gonzalez et al., 2015a), a pipeline to automate the primer design for polyploid organisms.

2.1 Pipeline

PolyMarker is an automated pipeline that takes as input a list of SNPs and a reference file and produces a list of primer triplets for SNP genotyping. The list of SNPs is first converted to a FASTA file with ambiguity codes (Cornish-Bowden, 1985). The template sequences are aligned with **exonerate** (Slater and Birney, 2005) to find the homoeologous regions to the target sequence. Then, the alignment between homoeologues is refined using **MAFFT** (Katoh and Standley, 2013). A list of candidate variations is produced and used as input for **Primer3** (Rozen and Skaletsky, 2000). Finally, the output of **Primer3** is parsed to find the best

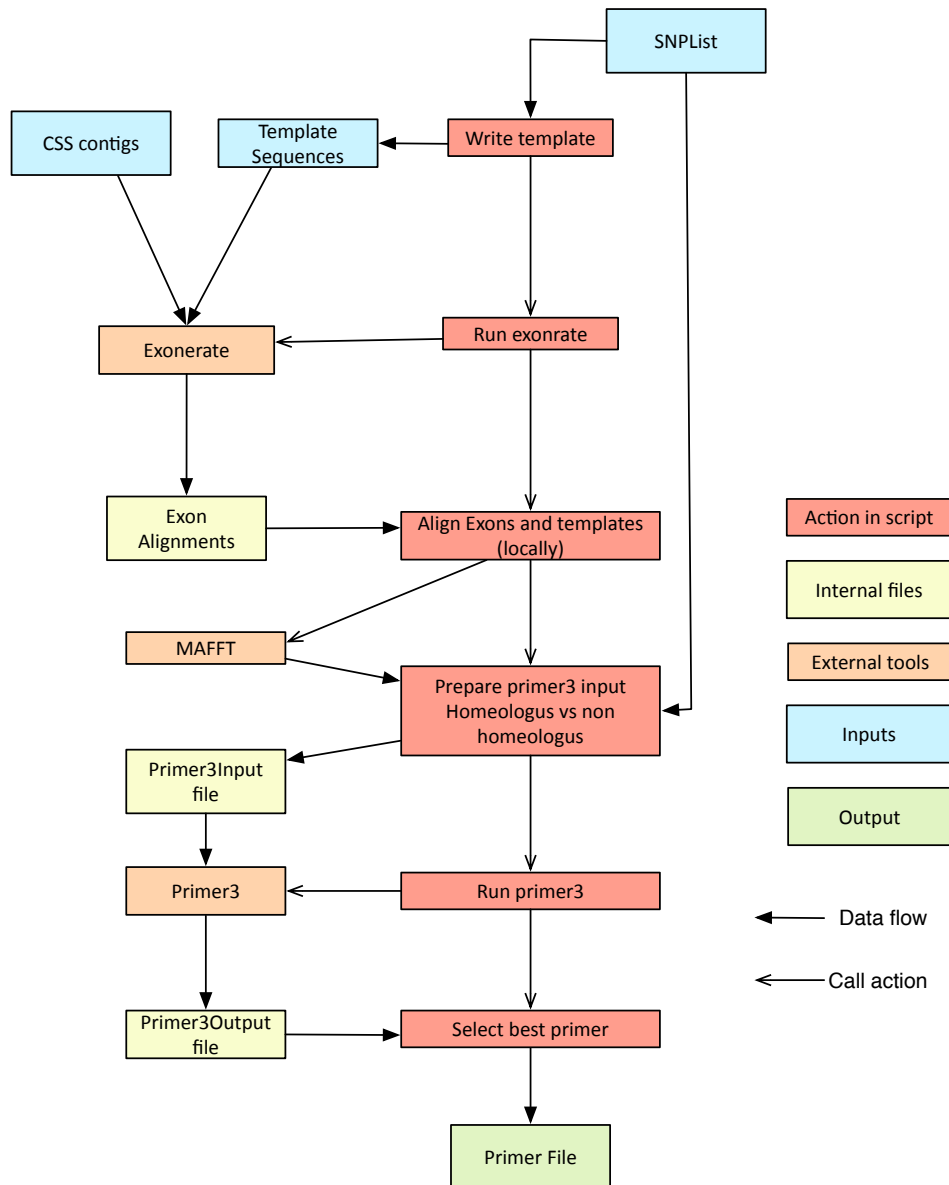


Figure 2.1: Steps and tools called by PolyMarker. The colour of the boxes represent: the step is an action inside the script (red); actions of the script (orange); temporary files (yellow); inputs (blue) and; output (green)

primer pair that contains a the targeted SNP and a base that is specific to the target genome (Figure 2.1). The pipeline is written as a Ruby script, using parsers and wrappers from BioRuby (Goto et al., 2010) and bio-samtools (Etherington et al., 2015; Ramirez-Gonzalez et al., 2012). The software is open source and released as a biogem (Bonnal et al., 2012), `bio-polyploid-tools`, the source code is available in github: <https://github.com/TGAC/bioruby-polyploid-tools>.

The PolyMarker input consist on SNP list with: unique name for the marker, the target chromosome and the sequence for the marker. The alternative alleles are surrounded by square brackets within the sequence. PolyMarker can take a list of several markers and design them in batch (Figure 2.2a). A FASTA file is produced with all the template sequences, with the alternative alleles substituted by the IUAPC ambiguity codes (Cornish-Bowden, 1985). The flanking sequence surrounding the SNP is limited by default to 100bp to reduce the search time and avoid missing regions that diverge near the SNP, as when the variation is near an intron-exon junction.

The template sequences are aligned to the reference using `exonerate` (Slater and Birney 2005; Figure 2.2b). The alignment is refined with the `--model est2genome` option, to allow the search of sequences coming from transcripts, a common source of SNPs (Allen et al., 2011). The `exonerate` output is formatted with the `--ryo` (roll your own format) to get an output easy to parse. All the hits that contain the SNP are extracted from the reference with a flanking sequence that extend out of the hit, by default, to 100bp on each side of the SNP, Figure 2.2c. The size of the flanking sequence can be set to different sizes to allow the design of different types of primers. Different homoeologues may contain small indels, Figure 2.2d. To enable a comparasion base-per-base, a local alignment with `MAFFT` (Kato and Standley, 2013) is produced, Figure 2.2e.

PolyMarker searches across each base in the local alignment to identify the variations across homoeologues and the target marker. A mask is produced to highlight the bases with a variations, Figure 2.2f, on the following categories:

| | |
|-----------------|--|
| Specific | Homoeologous polymorphism which is only present in the target genome (upper case). |
|-----------------|--|

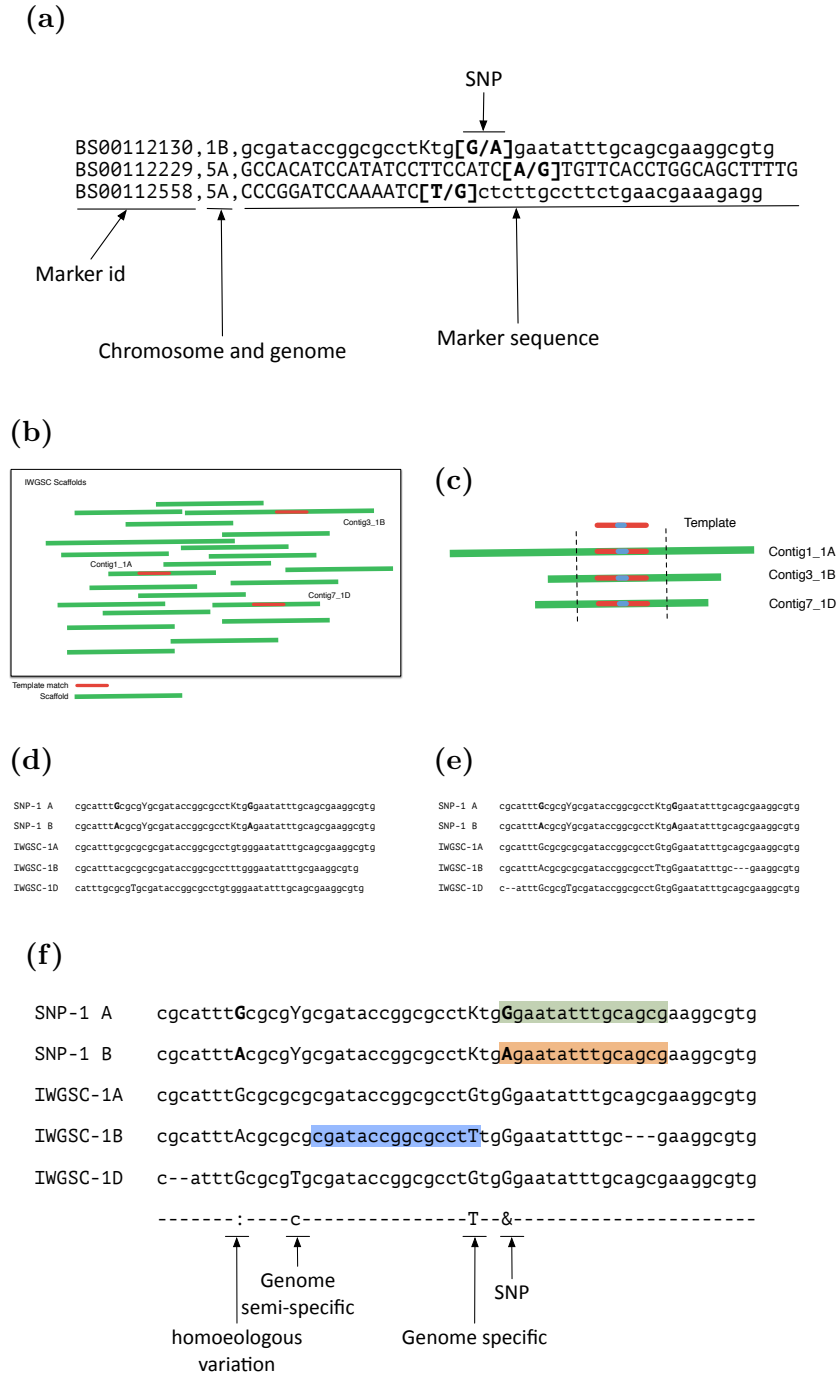


Figure 2.2: Alignments done by PolyMarker. (a) input. The alternative alleles are surrounded by brackets. (b) Global search of templates in the reference contigs. (c) Selected regions around the SNP on every chromosome. (d) Sequence of found regions around the SNP. (e) Local alignment on regions around the SNP detects indels. (f) Alignment with mask and primer candidates.

| | |
|-------------------------|--|
| Semi-specific | Homoeologous polymorphism which is found in 2 of the 3 genomes, hence it discriminates against one of the off-target genomes or when not all the homoeologous sequences were found (lower case). |
| Non-specific | No variation is found across homoeologues (-). |
| Homoeologous | The target SNP is present across different chromosomes, so candidate SNP markers on this category are not expected to be reliably identify the allele (:). |
| Non-homoeologous | The target SNP is not present across chromosomes, so it can be used to identify an allele (&). |

PolyMarker was designed to produce SNP assays for KASP genotyping (LGC Genomics, 2013), which requires a common primer and two allele-specific primers. The common primer is selected to start on a position from a: Specific; Semi-specific or; Non-specific, on that priority. This means that the common primer will be as specific as possible in the region. For the allele-specific primers, the starting position of the primer is on the base with the SNP. To ensure that the stability of the candidate primers will be met, the putative starting positions are tested with *Primer3* (Rozen and Skaletsky, 2000).

PolyMarker was designed and validated with the markers described in section 3.7. For wheat, PolyMarker uses the contigs from Mayer et al. (2014), as deposited in Ensembl. As new releases of the wheat genome are made available, different parsers to assign the chromosome to each sequence can be added with little effort to PolyMarker.

2.2 Applications of PolyMarker

PolyMarker is not restricted to wheat or to KASP assays, the source code is flexible and can be extended for other types of analysis. On each of the following projects, PolyMarker has been adapted to design primers in species where KASP hasn't been used before, the primers are used for regular PCR amplification, or the use of KASP is not the conventional SNP calling.

2.2.1 KASP assays for public sets of SNPs

PolyMarker was used to design KASP assays for the 81,587 markers from (Wang et al., 2014), available on the PolyMarker website and in CerealsDB (Wilkinson et al., 2012). Of those markers, 40,267 were designed using the target chromosome using the genetic map published by the genetic map. Genes without a genetic position were aligned to scaffolds sorted by chromosome from the International Wheat Genome Sequencing Consortium (Mayer et al., 2014) with BLAT (Kent, 2002) and the best hit was selected as putative location. 97.5% of the assays were designed and 76% of them are semi-specific or specific, thereby improving their expected performance with respect to randomly designed primers (Table A.1). A subset of the designed assay was used to genotype a mapping population to find resistance to Fusarium head blight (Burt et al., 2015).

2.2.2 SNPs in a mutant population

PolyMarker was used to design primers to validate SNPs in a Targeted Induced Local Lesions in Genomes (TILLING) population, an approach to identify the function of genes by mutating them. To calibrate the SNP calling, KASP assays were designed to get the mutations from M_2 , M_3 and, M_5 mutants (King et al., 2015). Then primers were designed for the whole mutant population, consisting of 1,200 Cadenza (Hexaploid) and 1,535 Kronos (Tetraploid) wheat lines (Krasileva et al., submitted 2016). Genome-specific primers 172 and 80 SNP assays on 19 and 8 M_4 Cadenza and Kronos lines respectively. Of those, 71(85.5%) Kronos and 147(88.8%) of the Cadenza primers were valid assays (Tables A.4 and A.5).

2.2.3 Deletions on a mutant population

On some of the TILLING mutant lines long deletions were detected (Krasileva et al., submitted 2016). To validate the deletions it is possible to use KASP assays to produce primers that amplify homoeologues. PolyMarker was modified to search for variations across homoeologues to select a common primer that will amplify two genomes (Figure 2.3a, b). On lines without the targeted deletion, the amplification corresponds to an homozygous assay (Figure 2.3c). When a deletion is present the results of the assay look like an homozygous sample, with the intensity

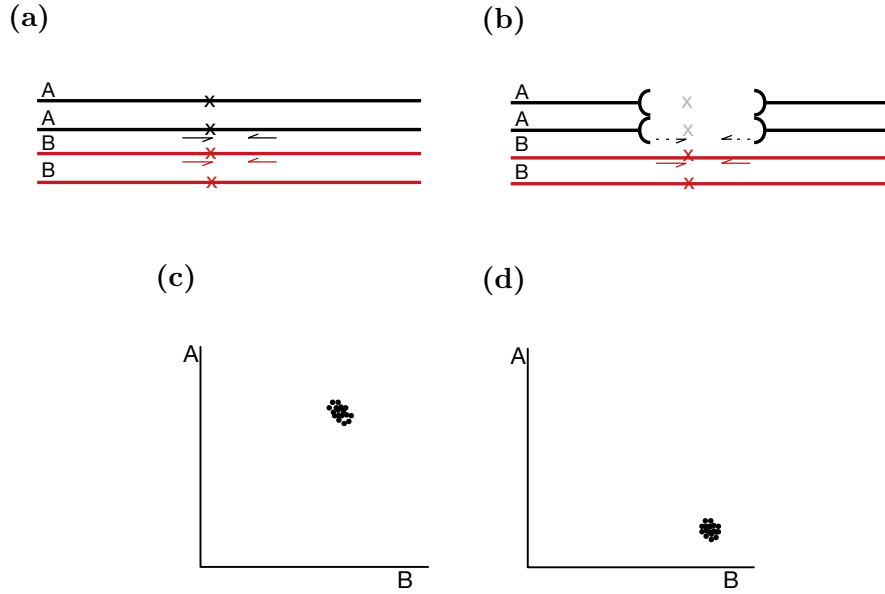


Figure 2.3: KASP assays to validate homozygous deletions. (a) Primer positions for wildtype. (b) Primer positions on homozygous deletion on M_4 (c) Heterozygous amplification on wildtype, including both homoeologues. (d) Homozygous amplification on deletion line, only the non-deleted homoeologue is amplified.

of the assay towards the the conserved homoeologue (Figure 2.3d). A set of KASP assays for the the deletions and mutations located on the same chromosome where designed to validate 11 homozygous deletions on M_4 plants. In all cases the segregation of the mutations was as expected, except for a predicted heterozygous mutation that was called as homozygous. Also, all the KASP assays that contained a deletion were called homozygous, as expected. To ensure that the calls didn't come from a single cluster, 4 wildtype plants were genotyped and the markers for deletions where called as heterozygous. An example of a validated deletion, with the calls for each individual is shown on Table A.3.

2.2.4 PolyMarker public web service

To make PolyMarker accessible to the community, a web server that allow the submission of SNPs was developed. The web interface consists on two virtual machines, one with a web facing interface that stores the queries, and a dedicated node to submit jobs to an HPC cluster. The on-line interface further simplifies the design of KASP assays, a process

that used to take a couple of weeks now is done in a couple of hours. Since the release of the public service in July 2014 until August 2016, 1,739 requests to PolyMarker have been done.

2.2.5 Genotyping of *Puccinia striiformis* f. sp. *tritici* isolates.

In Hubbard et al. (2015), *Puccinia striiformis* f. sp. *tritici* (PST) isolates were sequenced and assigned to clusters, according to their genotype. The clusters are useful to monitor the changes in the pathogen population, which can be used to predict if certain wheat lines will be resistant to the isolates in the field. PolyMarker was used to design primers for PST, using the assembly PST-130 (Cantu et al., 2011). Out of 15 assays 11 can be used to identify to which cluster of isolates a sample is likely to belong, Supplemental Table A.2.

2.3 Discussion

PolyMarker is a tool that was born as part of the validation of the SNPs found in Chapter 3. Originally, the primer design was ought to be done manually, a slow, error-prone and, repetitive process. The steps require the use of several bioinformatics tools, but once I figured out the steps I decided to automate the process. Since designing genome-specific primers is a common task in wheat research and breeding, the community showed interest on the tool and I decided to refine it and make it open source. PolyMarker has been used successfully in several projects and it even allowed the novel use of KASP assays to validate long deletions in polyploids.

The current web interface of PolyMarker is limited to KASP assays, however the command line version is more flexible and has been used to design primers for PCR amplicons, capillary sequencing and on other organisms. The ideas behind PolyMarker had been taken by other projects like the scripts described in Ma et al. (2015) and the corresponding web interface, GSP (Wang et al., 2016). As new references of wheat come available, PolyMarker should be updated to work with pseudomolecules and the web interface updated accordingly.

Chapter 3

Genetic map of *Yr15* with RNA-Seq

Wheat breeding programs aim to improve the wheat lines available for production. One of the traits desired in an elite line is the resistance to pathogens, such as *Puccinia striiformis* f. sp. *tritici*, the fungi responsible of yellow rust. A source of resistance genes is are introgressions from other species, such as *Triticum diccoides*. In the University of Sydney a collection of Near Isogenic Lines (NILs) with introgressions to several Yellow Rust resistance genes on a susceptible background were developed (Wellings and McIntosh, 1998). On this chapter the NIL for the *Yr15* locus is used to produce a mapping population to improve diagnostic markers.

Line selection can be done with molecular markers that can be used to test if certain allele is present in a line, without the need to do a phenotype. To find which regions are linked to a trait the use of F_2 mapping populations is a common practice. The population is produced by crossing two homozygous parents (P_1 and P_2) with different alleles, A/A (dominant) and a/a (recessive). When the trait is dominant and has a mendelian segregation, the F_1 population show the dominant trait, as it has a copy of each allele (A/a). The F_1 is then self-crossed to and the population segregates with a ration 1:2:1, dominant:heterozygous:recessive respectively. This generates a population with a phenotype ratio of 3:1 (dominant:recessive), since the effect of the recessive allele is masked by the dominant gene (Van Ooijen and Jansen 2013; Figure 3.1).

Bulk Segregant Analysis (BSA) consists on pooling the DNA of individuals with contrasting phenotypes (Michelmore et al., 1991) on a

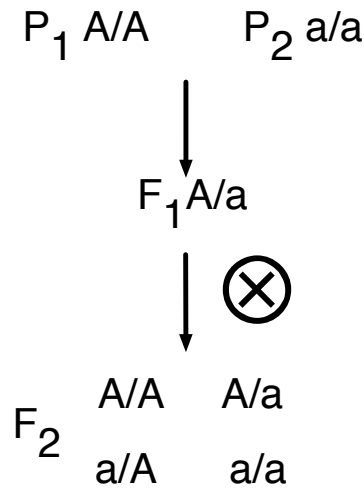
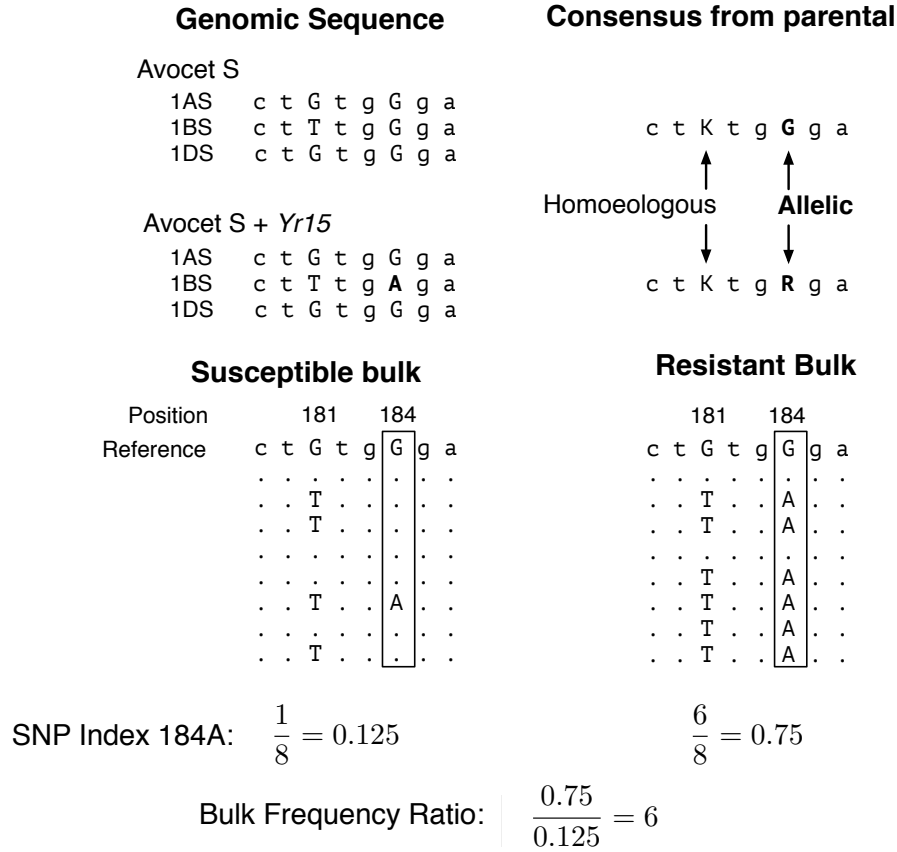


Figure 3.1: The cross of two homozygous parents, P_1 and P_2 , with a dominant and a recessive allele of a gene produce an heterozygous F_1 . The F_1 crossed with itself produce a segregating F_2 population with a 1:2:1 ratio (A/A:A/a:a/a). The upper and lower cases represent dominant and recessive alleles

segregating population. The bulks show as heterozygous except for the region that is linked to the trait of interest. This approach can be used to identify SNPs using High Throughput Sequencing, such as: exome capture (Hodges et al., 2007), RNA-Seq (Pickrell et al., 2010), whole genome resquencing (Schneeberger et al., 2009), among others.

To Call for SNPs from RNA-Seq a reference transcriptome is used as target when aligning the reads. The Bulk Frequency Ratio (BFR) methodology can work on organisms has more than one pseudo genome and that the genes are not necessarily fully characterised independently among homoeologues or paralogues, you can have in a single reference collapsing similar regions. The UniGenes database, from NCBI, contains the genes of each species, with all the variations of each gene automatically collapsed and represented as with the longest cDNA (Pontius et al., 2002). The UCW genes described in Krasileva et al. (2013) contains 94,177 models from tetraploid and hexaploid wheat, assembled and phased to separate different homoeologues. Both gene sets are complement each other, however, the UCW gene models should provide an improved alignment, since the different homoeologues aren't merged in a single model, a possible side effect of the UniGene pipeline.

Homoeologous variants, as exemplified by the G>T variant at position 181; K in consensus (Figure 3.2), will produce the same ambiguity code for both parental consensus sequences and can therefore be excluded. Real allelic SNPs between the parental genotypes, exemplified by the G>A variant at position 184; R in consensus, are distinguished by the presence in one, but not the other parental consensus sequence. The



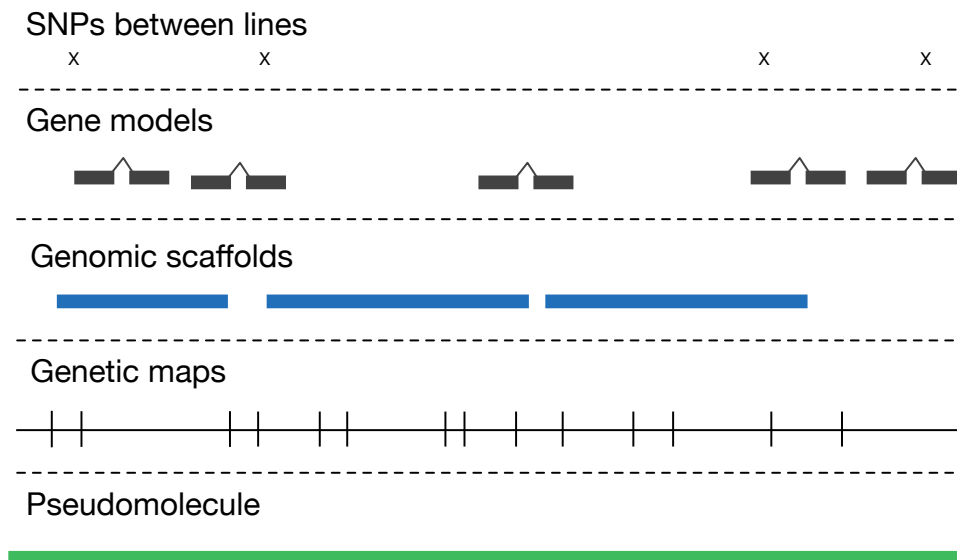


Figure 3.3: Layers of information to do *In Silico* mapping. SNPs are called from gene models. The genes and markers from genetic maps are aligned to scaffolds. The order of the markers in a genetic map can be used to sort the scaffolds.

allelic SNPs are then examined further with the alignments of the bulks to identified the SNPs that are enriched on the resistant plants. The SNP index is the proportion of times an alternative allele is observed over the coverage at certain, in the example the the susceptible bulk has an SNP index of $1/8 = 0.125$ and $6/8 = 0.75$ for the resistant bulk (Takagi et al., 2013b). traditional The BFR are then calculated by dividing the SNP Index of sample containing the target phenotype (resistance) over the sample without the trait (susceptible), on the example is $0.75/0.125 = 6$. A high BFR suggests that the SNP is linked to the target trait (Trick et al., 2012). The results of the BFR analysis on the F_2 population are discussed in Section 3.4.

There are several layers of information that can be used to add a context to the SNPs. When the SNPs are called from genes like the Uni-Genes (Pontius et al., 2002) or the UCW gene models (Krasileva et al., 2013), the location of the genes can be assigned by aligning them to a genomic reference, even if it is fragmented. A source to get the order of the scaffolds are genetic maps previously published, such as the genetic map described in Wang et al. (2014), which has the sequence of the markers available. The markers and the genes can be aligned to the scaffolds with a high percentage of identity (over 98%), to avoid them being assigned

to an homoeologue or paralogue region in a different chromosome. The use of genetic maps to sort genomic sequence is frequently used to produce pseudo-chromosomes on genome wide projects, usually with ad-hoc tools (Tang et al., 2015). Since the CSS assembly is quite fragmented the genetic maps don't have enough resolution to produce a pseudomolecule, however it is enough to sort the scaffolds in bins when several markers map to the same location. In this way, it is possible to use the scaffolds as a proxy to map the genes to their genetic position (Figure 3.3). The results of mapping the genes with SNPs to the CSS assembly and the genetic map are described in Section 3.5. For a longer description of resources available for wheat see Section 1.6.

Finally, the best candidate SNPs were selected to produce a genetic map which lead to a triplet of markers diagnostic to the target locus.

The steps described in this chapter were first published in Ramirez-Gonzalez et al. (2015b) and the results of this chapter are published in Ramirez-Gonzalez et al. (2015c).

3.1 Mapping population

The population was developed by crossing the resistant line Avocet + *Yr15* (*Yr15*) (Wellings and McIntosh, 1998), Figure 3.4a, to the susceptible line Avocet S (AVS), Figure 3.4b. *Yr15* is a NIL of a 6th generation Back-cross (BC) and the AVS background is highly susceptible to yellow rust, hence the resistance is conferred by the *Yr15* locus. F_2 seeds from three independent F_1 plants were sown and tissue was collected, before the fungal inoculation to avoid the effect of the response on the gene expression. The plants were challenged at the three leaf stage as it is known that *Yr15* confers resistance in seedlings (Gerechter-Amitai et al., 1989). The expected segregation on an F_2 population is 3:1 (resistant:susceptible), since *Yr15* is a dominant gene. From the 232 plants in the F_2 population that germinated, 187 were resistant and 45 were susceptible, which deviates slightly from the expected ratio ($\chi^2 = 0.049$). Segregation distortion has been shown for the same *Yr15* donor (Randhawa et al., 2009), however the decreased number of susceptible plants can be explained by escapes in the virulence essays (i.e. plants scored as resistant without the *Yr15* locus). For this study we extracted DNA from individual plants in the F_2 population and we bulked RNA on 6

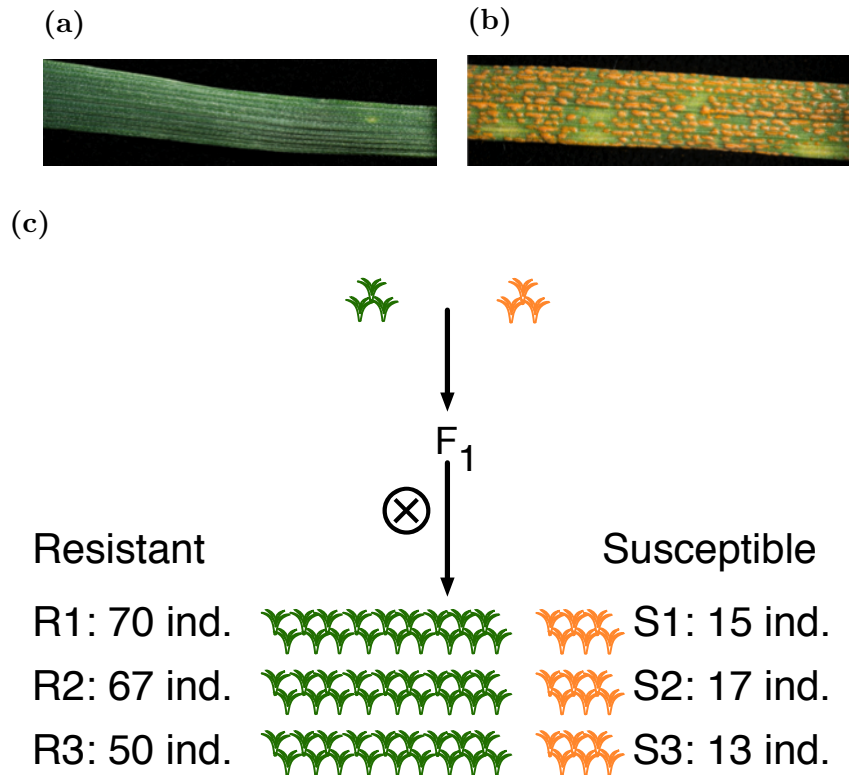


Figure 3.4: Response of (a) Avocet + Yr15 and (b) Avocet when inoculated with *Puccinia striiformis* f. sp. *tritici* at the three leaf stage. (c) The phenotype of the F_2 population was used to produce 6 bulks, 3 resistant and 2 susceptible. The RNA was pooled in bulks accordingly. Adapted from (Ramirez-Gonzalez et al., 2015c)

Table 3.1: Arrangement and number of sequenced base pairs per sample.

| Library | name | Bar code | Lane | Reads ($\times 10^8$ bp) |
|---------|-----------------------|----------|------|---------------------------|
| LIB1715 | Bulk R1 | ATCACG | 1 | 0.77 |
| LIB1716 | Bulk R2 | TAGCTT | 1 | 1.20 |
| LIB1717 | Bulk R3 | ACTTGA | 2 | 0.96 |
| LIB1718 | Bulk S1 | GGCTAC | 2 | 1.64 |
| LIB1719 | Bulk S2 | CGTACG | 2 | 1.49 |
| LIB1720 | Bulk S3 | GTGGCC | 1 | 1.88 |
| LIB1721 | AvocetS | N/A | 3 | 4.13 |
| LIB1722 | AvocetS + <i>Yr15</i> | N/A | 4 | 3.99 |

different bulks: 3 resistant and, 3 susceptible (Figure 3.4c).

3.2 Sequencing and mapping

RNA-Seq was used to avoid sequencing the non-coding regions and reduce the search space. The sequencing of the bulks and the parents were done on a single Illumina Hi-Seq2000 each. The bulks were multiplexed and sequenced on a third of a lane each, as shown on Table 3.1. To ensure that the quality of the sequencing was good, **fastqc-0.10** (Babraham Bioinformatics, 2012) was run with its default parameters in each one of the fastq files. The GC content was around 52% in all the samples (Appendx B.2), which is expected as the sample should be of coding regions, and for wheat the reported GC content in genes is around 55%. The quality of the reads is fairly consistent, in general dropping after the base 80 across the samples (Appendix B.1).

When the analysis was started, the draft genome and the corresponding annotation were not yet released, hence gene models were used. All the samples were aligned to the Unigenes v60 (56,954 genes) and the gene models from UCW (Krasileva et al., 2013) using **BWA 0.5.9** (Li and Durbin, 2009). The alignment provided showed that a few genes were overexpressed, however we still have 22,107 and 36,808 genes, on the Unigenes and the UCW gene set respectively, with a coverage greater than 20x in the progenitor with *Yr15*. Both gene sets performed similarly in terms of the percentage of genes with reads and percentage of aligned reads. For AVS and *Yr15*, the percentage of genes with a coverage of at least 20x is 45% and 39% respectively across both references (Figure 3.5a). Since each individual bulk has a lower coverage, the susceptible

Table 3.2: Number of genes with a coverage over 20x, 10x and at least one read (>0x).

| Coverage | Reference | Bulks | | | | | Bulk mixes | | | | Progenitors | | |
|----------|-------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| | | R1 | R2 | R3 | S1 | S2 | S3 | R1+R2 | S1+S2 | R1+R2+R3 | S1+S2+S3 | Yr15 | AVS |
| 20x | UCW | 16,434 17% | 27,871 30% | 27,223 29% | 32,287 34% | 28,669 30% | 34,898 37% | 33,968 36% | 41,019 44% | 40,985 44% | 47,507 50% | 36,808 39% | 42,248 45% |
| | UniGene v60 | 9,643 17% | 16,182 28% | 15,222 27% | 19,549 34% | 17,397 31% | 20,567 36% | 20,219 36% | 25,270 44% | 24,598 43% | 29,052 51% | 22,107 39% | 25,842 45% |
| 10x | UCW | 27,371 29% | 38,282 41% | 37,777 40% | 42,658 45% | 38,999 41% | 44,610 47% | 43,266 46% | 49,473 53% | 49,182 52% | 54,781 58% | 46,356 49% | 50,760 54% |
| | UniGene v60 | 16,201 28% | 22,948 40% | 22,130 39% | 26,200 46% | 24,130 42% | 26,914 47% | 26,318 46% | 30,579 54% | 29,857 52% | 33,557 59% | 28,044 49% | 31,095 55% |
| >0x | UCW | 68,302 73% | 72,484 77% | 72,957 77% | 74,694 79% | 73,290 78% | 75,201 80% | 74,397 79% | 77,093 82% | 76,715 81% | 78,796 84% | 76,275 81% | 77,080 82% |
| | UniGene v60 | 40,717 71% | 42,489 75% | 42,595 75% | 43,625 77% | 43,059 76% | 43,748 77% | 43,393 76% | 44,655 78% | 44,364 78% | 45,392 80% | 43,732 77% | 44,596 78% |

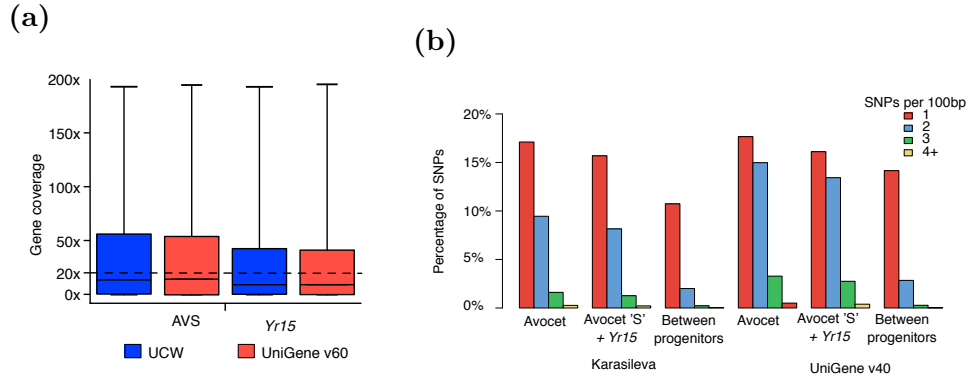


Figure 3.5: (a) Box plot distribution of the gene coverage of the parent reads (AVS and *Yr15*) across the UCW (blue) and the UniGene (red) gene models. The dashed line represents the 209 minimum coverage required for SNP calling. The full line represents the average coverage across all gene models. (b) Percentage of genes exhibiting SNPs across references. The number of SNPs between the parent reads and the corresponding references was calculated (per 100 bp, rounded). The between-parents category corresponds to putative SNPs when comparing the consensus sequence between AVS and *Yr15*. Adapted from Ramirez-Gonzalez et al. (2015c)

and resistant reads were merged *in silico* as: (i) susceptible bulks 1 with 2 (S1 + S2) and resistant bulks 1 with 2 (R1 + R2) and (ii) all the susceptible (S1 + S2 + S3) and resistant bulks (R1 + R2 + R3). The merged samples increased the percentage of genes with coverage over 20x to 44% and 50% in the resistant and susceptible bulks (Table 3.2), which is close to the coverage from the progenitors.

3.3 SNP Calling

The SNP calling was done on positions with a coverage of at least 20x on the progenitor lines against the gene reference. The AVS progenitor had roughly 3% more genes with polymorphisms than *Yr15*, consistent with the difference in coverage, suggesting that with a higher coverage we could recover more SNPs from *Yr15*. The UniGenes have a higher number of SNPs because the UCW gene models have a higher number of monomorphic genes when compared to the UniGenes. (Figure 3.5b; Table 3.3). The difference in the number of relative monomorphic SNPs between references can be explained by the fact that the UniGenes have homoeologues can be represented as a single sequence, as opposed to the UCW set which are homoeologue-specific, improving the mapping to the

Table 3.3: Count of SNPs per 100 bp on genes with at least 20x coverage.

| SNPs per 100bp | UCW | | | UniGene v60 | | |
|----------------------|-----------------|---------------------|------------------------|-----------------|---------------------|------------------------|
| | AVS | AVS+ <i>Yr15</i> | Between progenitors | AVS | AVS+ <i>Yr15</i> | Between progenitors |
| 0 | 67,389 71.6% | 70,338 74.7% | 81,921 87.0% | 36,210 63.6% | 38,339 67.3% | 47,097 82.7% |
| 1 | 16,111 17.1% | 14,770 15.7% | 10,107 10.7% | 10,058 17.7% | 9,175 16.1% | 8,061 14.2% |
| 2 | 8,904 9.5% | 7,676 8.2% | 1,893 2.0% | 8,529 15.0% | 7,648 13.4% | 1,621 2.9% |
| 3 | 1,517 1.6% | 1,192 1.3% | 215 0.2% | 1,870 3.3% | 1,568 2.8% | 59 0.3% |
| 4+ | 253 0.3% | 198 0.2% | 38 0.0% | 287 0.5% | 224 0.4% | 16 0.0% |

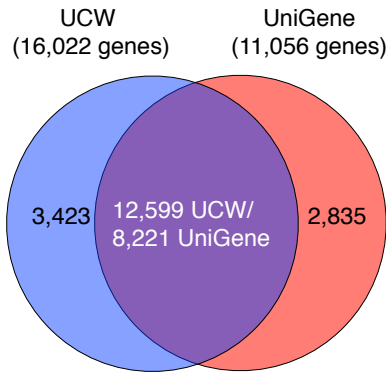


Figure 3.6: Gene models with putative SNPs in common between the UCW and UniGenes reference. The intersection represents the genes that are common in both sets. Adapted from Ramirez-Gonzalez et al. (2015c)

correct homoeologue in the genes from the UCW set over the UniGenes.

Both gene sets were done from varieties different to AVS and are likely to be incomplete, hence we set a low threshold of at least 20% of the observed nucleotides on any position to call an SNP. To represent cases where more than one consensus base is called we use International Union of Pure and Applied Chemistry (IUPAC) codes (Cornish-Bowden (1985); Section 1.5.1; Figure 3.2). To focus the analysis on informative SNPs, the common varietal SNPs and variations between homoeologues were removed by finding the cases when the consensus call on both progenitors is the same. The SNPs that are unique to a single parental were examined in detail. There are 66,426 putative SNPs across 16,022 (17%) UCW genes and 52,262 SNPs on 11,056 UniGenes (19.4%; Figure 3.6).

The high number of genes with SNPs was unexpected as a BC6 NIL used for an F_2 mapping population expects to have $< 1\%$ of the genetic background segregating. The both sets of gene models were aligned

Table 3.4: Number of genes with SNPs assigned to the wheat chromosome arm CSS scaffolds (Mayer et al., 2014) using the best hit from BLAT (Kent, 2002)

| Wheat Chromosome Arm | UCW (94,177) | UniGene v60 (56,954) | Total (151,131) |
|----------------------------|-----------------|----------------------|------------------|
| 1AL | 3,251 (3.45%) | 1,404 (2.47%) | 4,655 (3.08%) |
| 1AS | 1,366 (1.45%) | 560 (0.98%) | 1,926 (1.27%) |
| 1BL | 2,610 (2.77%) | 1,280 (2.25%) | 3,890 (2.57%) |
| 1BS | 1,487 (1.58%) | 693 (1.22%) | 2,180 (1.44%) |
| 1DL | 997 (1.06%) | 1,057 (1.86%) | 2,054 (1.36%) |
| 1DS | 753 (0.80%) | 687 (1.21%) | 1,440 (0.95%) |
| 2AL | 3,491 (3.71%) | 1,460 (2.56%) | 4,951 (3.28%) |
| 2AS | 2,305 (2.45%) | 974 (1.71%) | 3,279 (2.17%) |
| 2BL | 3,658 (3.88%) | 1,546 (2.71%) | 5,204 (3.44%) |
| 2BS | 2,790 (2.96%) | 1,139 (2.00%) | 3,929 (2.60%) |
| 2DL | 1,098 (1.17%) | 1,069 (1.88%) | 2,167 (1.43%) |
| 2DS | 796 (0.85%) | 833 (1.46%) | 1,629 (1.08%) |
| 3AL | 2,135 (2.27%) | 978 (1.72%) | 3,113 (2.06%) |
| 3AS | 1,543 (1.64%) | 718 (1.26%) | 2,261 (1.50%) |
| 3B | 6,559 (6.96%) | 2,839 (4.98%) | 9,398 (6.22%) |
| 3DL | 915 (0.97%) | 938 (1.65%) | 1,853 (1.23%) |
| 3DS | 412 (0.44%) | 450 (0.79%) | 862 (0.57%) |
| 4AL | 3,393 (3.60%) | 1,335 (2.34%) | 4,728 (3.13%) |
| 4AS | 2,011 (2.14%) | 817 (1.43%) | 2,828 (1.87%) |
| 4BL | 2,119 (2.25%) | 898 (1.58%) | 3,017 (2.00%) |
| 4BS | 1,946 (2.07%) | 892 (1.57%) | 2,838 (1.88%) |
| 4DL | 1,069 (1.14%) | 945 (1.66%) | 2,014 (1.33%) |
| 4DS | 800 (0.85%) | 699 (1.23%) | 1,499 (0.99%) |
| 5AL | 2,640 (2.80%) | 1,132 (1.99%) | 3,772 (2.50%) |
| 5AS | 963 (1.02%) | 407 (0.71%) | 1,370 (0.91%) |
| 5BL | 5,324 (5.65%) | 1,943 (3.41%) | 7,267 (4.81%) |
| 5BS | 1,360 (1.44%) | 591 (1.04%) | 1,951 (1.29%) |
| 5DL | 2,067 (2.19%) | 1,688 (2.96%) | 3,755 (2.48%) |
| 5DS | 620 (0.66%) | 614 (1.08%) | 1,234 (0.82%) |
| 6AL | 2,397 (2.55%) | 896 (1.57%) | 3,293 (2.18%) |
| 6AS | 2,285 (2.43%) | 936 (1.64%) | 3,221 (2.13%) |
| 6BL | 1,564 (1.66%) | 820 (1.44%) | 2,384 (1.58%) |
| 6BS | 1,308 (1.39%) | 731 (1.28%) | 2,039 (1.35%) |
| 6DL | 1,399 (1.49%) | 1,050 (1.84%) | 2,449 (1.62%) |
| 6DS | 870 (0.92%) | 680 (1.19%) | 1,550 (1.03%) |
| 7AL | 1,918 (2.04%) | 849 (1.49%) | 2,767 (1.83%) |
| 7AS | 1,717 (1.82%) | 764 (1.34%) | 2,481 (1.64%) |
| 7BL | 1,592 (1.69%) | 776 (1.36%) | 2,368 (1.57%) |
| 7BS | 1,239 (1.32%) | 713 (1.25%) | 1,952 (1.29%) |
| 7DL | 2,040 (2.17%) | 1,301 (2.28%) | 3,341 (2.21%) |
| 7DS | 1,224 (1.30%) | 1,016 (1.78%) | 2,240 (1.48%) |
| Assigned | 80,031 (84.98%) | 41,118 (72.20%) | 121,149 (80.16%) |

Table 3.5: Total number of SNPs scored in parents, individual bulks and in silico merged bulks.

| Gene set | $\frac{R1}{S1}$ | $\frac{R2}{S2}$ | $\frac{R3}{S3}$ | $\frac{R1+R2}{S1+S2}$ | $\frac{R1+R2+R3}{S1+S2+S3}$ | SNPs in parents |
|-------------|------------------|------------------|------------------|-----------------------|-----------------------------|-----------------|
| UCW | 16,269 24.49% | 29,703 44.72% | 31,891 48.01% | 44,224 66.58% | 64,522 97.13% | 66,426 |
| UniGene v60 | 15,261 29.20% | 25,143 48.11% | 24,548 46.97% | 35,698 68.31% | 49,738 95.17% | 52,262 |

with BLAT (Kent, 2002) to the Chinese Spring Chromosome arm survey sequence (CSS; Mayer et al. 2014); the alignment resulted on 80,031 (85.0%) UCW gene models and 41,118 (72.2%) UniGenes assigned to a chromosome arm (Table 3.4). The SNPs found in the mapped genes are evenly distributed across all the chromosomes (Figure 3.10a), suggesting that the Avocet S (JIC, UK) used as parent in the F_2 is different to the Avocet S used for the *Yr15* NIL development (University of Sydney, Australia).

To confirm that the Avocet S seed stocks from JIC are distinct to the stocks in Sydney DNA from both stocks was procured and compared with the iSelect 90k wheat SNP chip. Between two independent Avocet S seeds from JIC only 58 out of 71,972 (0.08%) valid assays were polymorphic. Nonetheless, there are over 5,000 (> 7.5%) assays with polymorphisms between JIC-Avocet S and Avocet S from Sydney. The difference was not expected originally, but considering that the Avocet S seeds are coming from different stocks and the fact that in both countries commercial varieties with the same name had been released, it is not surprising.

3.4 Bulk Frequency Ratios

The objective was to find the SNPs enriched on each bulk and hence linked to the phenotype, variations from *Yr15* to resistance and from AVS to susceptibility in the segregating population. Across individual bulks, it was possible to score between 15,261 (24.5%) to 31,891 (48.0%) SNPs across both reference sets. On the *in silico* mixes over 95% of SNPs were scored (Table 3.5), suggesting that the coverage of individual bulks is not enough to score all the SNPs. The scoring was done with the Bulk Frequency Ratio (Trick et al. 2012; Figure 3.2; Section 3.8.2), which has a value that increases as the *Yr15* allele is observed more times relatively

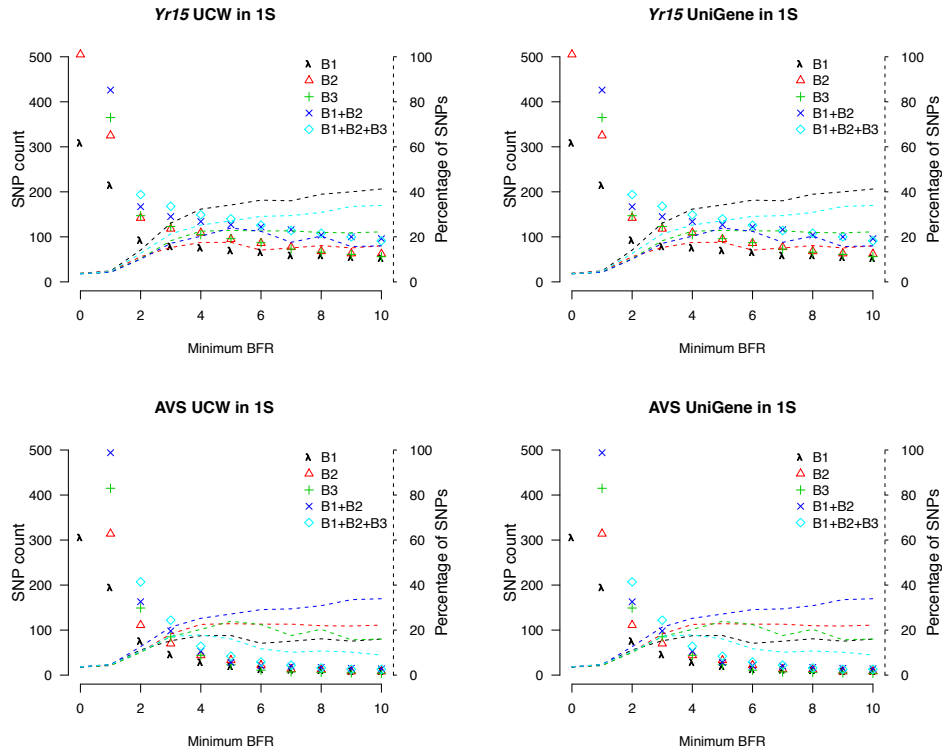


Figure 3.7: Effect of BFR threshold on the number of SNPs across the short arm of chromosome group 1. Figure previously published in Ramirez-Gonzalez et al. (2015c).

to the AVS allele.

When increasing the minimum BFR threshold, enrichment of SNPs was observed in the short arm of the group 1 chromosomes (1S). Without taking in account the BFR, 3.6% of the SNPs are located in the 1S group, similar to the number of SNPs located in other groups 3.4. However, when increasing the threshold (between $BFR > 5$ and $BFR > 7$) the relative number of SNPs in group 1S increases. After $BFR > 7$ the gains in relative enrichment only improves marginally, but the number of called SNPs is reduced (Table 3.6; Figure 3.7). For that reason, SNPs with a $BFR > 6$ were selected for further validation. The method described by Trick et al. (2012) was extended by including cases where there is a complete lack of coverage in one of the samples ($BFR = \infty$), which is an ideal case where the linkage between the SNP and the phenotype is perfect. A total of 1,582 SNPs across 1,173 genes had a $BFR > 6$.

Table 3.6: SNPs in chromosome group 1S vs total number of SNPs with a minimum BFR from 0 to 10. AVS: SNPs coming from Avocet S. Yr15: SNPs coming from Avocet + Yr15.

| Min BFR | Gene Set | R1/S1 Yr15 | R1/S1 AVS | R2/S2 Yr15 | R2/S2 AVS | R3/S3 Yr15 | R3/S3 AVS | S1+2/ R1+2 Yr15 | S1+2/ R1+2 AVS | S1+S2+S3/ R1+R2+R3 Yr15 | S1+S2+S3/ R1+R2+R3 AVS |
|---------|-------------|----------------------|----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-------------------------------|------------------------------|
| 0 | UCW | 308/8,049 (3.83%) | 305/8,220 (3.71%) | 505/14,121 (3.58%) | 556/15,582 (3.57%) | 532/14,875 (3.58%) | 623/17,016 (3.66%) | 670/18,760 (3.57%) | 885/25,464 (3.48%) | 860/24,026 (3.58%) | 1,505/40,496 (3.72%) |
| | UniGene v60 | 307/7,823 (3.92%) | 299/7,438 (4.02%) | 428/12,409 (3.45%) | 421/12,734 (3.31%) | 427/12,050 (3.54%) | 415/12,498 (3.32%) | 536/15,672 (3.42%) | 595/20,026 (2.97%) | 712/19,358 (3.68%) | 901/30,380 (2.97%) |
| 1 | UCW | 214/4,415 (4.85%) | 194/4,108 (4.72%) | 325/7,603 (4.26%) | 314/7,374 (4.26%) | 365/7,920 (4.61%) | 415/8,850 (4.69%) | 426/10,122 (4.21%) | 494/12,185 (4.05%) | 539/13,037 (4.13%) | 842/19,466 (4.33%) |
| | UniGene v60 | 207/4,474 (4.63%) | 194/3,630 (5.34%) | 269/6,649 (4.05%) | 269/6,193 (4.34%) | 279/6,511 (4.29%) | 272/6,436 (4.23%) | 329/8,704 (3.78%) | 369/9,343 (3.95%) | 446/10,860 (4.11%) | 541/14,226 (3.80%) |
| 2 | UCW | 92/651 (14.13%) | 75/671 (11.18%) | 142/1,372 (10.35%) | 111/1,101 (10.08%) | 147/1,162 (12.65%) | 149/1,411 (10.56%) | 167/1,324 (12.61%) | 163/1,478 (11.03%) | 194/1,370 (14.16%) | 207/1,765 (11.73%) |
| | UniGene v60 | 77/568 (13.56%) | 58/527 (11.01%) | 101/1,017 (9.93%) | 81/720 (11.25%) | 105/775 (13.55%) | 84/867 (9.69%) | 122/991 (12.31%) | 116/973 (11.92%) | 132/1,030 (14.08%) | 132/1,210 (10.91%) |
| 3 | UCW | 78/299 (26.09%) | 45/295 (15.25%) | 118/646 (18.27%) | 70/409 (17.11%) | 123/577 (21.32%) | 85/494 (17.21%) | 145/673 (21.55%) | 98/563 (17.41%) | 168/768 (21.88%) | 122/665 (18.35%) |
| | UniGene v60 | 65/254 (25.59%) | 26/186 (13.98%) | 87/499 (17.43%) | 54/294 (18.37%) | 93/379 (24.54%) | 48/315 (15.24%) | 107/525 (20.38%) | 66/379 (17.41%) | 133/617 (21.56%) | 78/489 (15.95%) |
| 4 | UCW | 75/232 (32.33%) | 28/160 (17.50%) | 109/484 (22.52%) | 44/217 (20.28%) | 105/416 (25.24%) | 44/246 (17.89%) | 134/539 (24.86%) | 53/277 (19.13%) | 149/640 (23.28%) | 64/323 (19.81%) |
| | UniGene v60 | 63/192 (32.81%) | 17/104 (16.35%) | 83/390 (21.28%) | 29/155 (18.71%) | 82/288 (28.47%) | 29/173 (16.76%) | 104/431 (24.13%) | 40/214 (18.69%) | 127/519 (24.47%) | 29/266 (10.90%) |
| 5 | UCW | 69/202 (34.16%) | 19/108 (17.59%) | 95/416 (22.84%) | 33/138 (23.91%) | 96/354 (27.12%) | 23/143 (16.08%) | 127/477 (26.62%) | 28/175 (16.00%) | 140/580 (24.14%) | 42/222 (18.92%) |
| | UniGene v60 | 58/163 (35.58%) | 11/70 (15.71%) | 76/337 (22.55%) | 14/102 (13.73%) | 70/228 (30.70%) | 20/112 (17.86%) | 100/389 (25.71%) | 23/146 (15.75%) | 118/469 (25.16%) | 21/178 (11.80%) |
| 6 | UCW | 65/179 (36.31%) | 12/85 (14.12%) | 86/380 (22.63%) | 22/98 (22.45%) | 87/299 (29.10%) | 11/94 (11.70%) | 122/429 (28.44%) | 21/130 (16.15%) | 126/514 (24.51%) | 29/165 (17.58%) |
| | UniGene v60 | 57/151 (37.75%) | 7/48 (14.58%) | 73/300 (24.33%) | 6/71 (8.45%) | 65/191 (34.03%) | 13/84 (15.48%) | 98/358 (27.37%) | 20/122 (16.39%) | 115/439 (26.20%) | 16/143 (11.19%) |
| 7 | UCW | 58/161 (36.02%) | 11/73 (15.07%) | 77/340 (22.65%) | 13/74 (17.57%) | 73/248 (29.44%) | 7/69 (10.14%) | 116/393 (29.52%) | 20/111 (16.02%) | 114/468 (24.36%) | 22/143 (15.38%) |
| | UniGene v60 | 56/132 (42.42%) | 4/37 (10.81%) | 68/273 (24.91%) | 5/58 (8.62%) | 60/171 (35.09%) | 9/64 (14.06%) | 94/334 (28.14%) | 18/103 (17.48%) | 113/412 (27.43%) | 16/124 (12.90%) |
| 8 | UCW | 58/149 (38.93%) | 10/62 (16.13%) | 68/310 (21.94%) | 12/59 (20.34%) | 66/214 (30.84%) | 6/56 (10.71%) | 104/359 (28.97%) | 17/102 (16.67%) | 108/429 (25.17%) | 16/119 (13.45%) |
| | UniGene v60 | 55/126 (43.65%) | 3/33 (9.09%) | 64/255 (25.10%) | 5/50 (10.00%) | 55/150 (36.67%) | 9/55 (16.36%) | 91/313 (29.07%) | 14/89 (15.73%) | 105/376 (27.93%) | 15/108 (13.89%) |
| 9 | UCW | 54/135 (40.00%) | 8/53 (15.09%) | 63/289 (21.80%) | 8/51 (15.69%) | 61/182 (33.52%) | 5/49 (10.20%) | 100/331 (30.21%) | 15/91 (16.48%) | 100/387 (25.84%) | 13/106 (12.26%) |
| | UniGene v60 | 53/117 (45.30%) | 1/30 (3.33%) | 62/244 (25.41%) | 5/46 (10.87%) | 50/136 (36.76%) | 9/48 (18.75%) | 88/291 (36.76%) | 13/83 (15.66%) | 97/345 (28.12%) | 12/99 (12.12%) |
| 10 | UCW | 52/126 (41.27%) | 8/50 (16.00%) | 62/279 (22.22%) | 8/50 (16.00%) | 56/165 (33.94%) | 4/45 (8.89%) | 96/309 (33.94%) | 14/82 (17.07%) | 91/355 (25.63%) | 13/100 (13.00%) |
| | UniGene v60 | 50/105 (47.62%) | 1/28 (3.57%) | 60/226 (26.55%) | 5/39 (12.82%) | 43/119 (36.13%) | 7/45 (15.56%) | 85/272 (31.25%) | 13/82 (15.85%) | 92/318 (28.93%) | 12/97 (12.37%) |

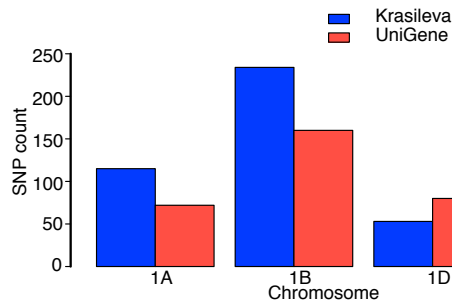


Figure 3.8: Location of SNPs with $BFR > 6$ according to the best alignment of the UniGene (red) and UCW (blue) gene models to the flow-sorted group 1 chromosomes from the Chinese Spring Survey sequence (CSS) (Mayer et al., 2014). Figure adapted from Ramirez-Gonzalez et al. (2015c).

3.5 *In silico* mapping

From the mapped SNPs with a $BFR > 6$, 872 of 1470 ($\sim 60\%$) were assigned to the chromosomes in group 1 of hexaploid wheat, being the only group with more than 4% of the SNPs assigned to it (Table 3.7). From the group 1, the B genome contained the higher proportion of SNPs mapped (54%), having 255 (54%) and 214 (46%) assigned to the long and short arms respectively (Figure 3.8). This results are expected since previous studies have located *Yr15* near the centromere in the short arm of chromosome 1B and, the *Yr15* introgression contains regions from the long and short arm from *T. diccocooides* (Murphy et al., 2009; Peng et al., 2000; Sun et al., 1997).

The CSS assembly was used as a common reference between the reference genes and the SNPs 40,266 SNP markers published at the time when this analysis was done (Wang et al., 2014) to locate the SNPs with a $BFR > 6$ (including $BFR = \infty$) in a genomic position (Figures 3.9, 3.10). From the 1,582 SNPs across 1,173 genes, only 678 SNPs (43%, 474 genes) were successfully located in the genetic map. Since the CSS assembly is quite fragmented, the low percentage of located SNPs can be because not all candidate SNPs had a corresponding scaffold that has at least one of the 40,266 markers in the genetic map. Even if the number of located SNPs was not enough to give a position for over 50% of the SNPs from the parental line, the resolution of the genetic position SNPs that were assigned improved over just having the chromosome arm information from the CSS assembly. The mapping position further confirmed an enrichment of SNPs near the centromere of chromosome 1B with 325 out of 678 SNPs. Furthermore, 311 of those were located within an interval of 30cM (Figures 3.10b, 3.9a).

Studies in diploid organisms using QTL-Seq (Takagi et al., 2013a) or other NGS-enable genetic approaches (James et al., 2013) have shown

Table 3.7: SNP and genes with BFR i 6 mapping to each of the chromosomes from the CSS assemblies. The chromosome assignment on the "Genetically mapped" column correspond to the map published in Wang et al. (2014)).

| Reference Chromosome | CSS assemblies | | | | Genetically mapped | | | |
|-------------------------|-----------------|--------|-------------|--------|--------------------|--------|-------------|--------|
| | UCW gene models | | UniGene v60 | | UCW gene models | | UniGene v60 | |
| | SNPs | Genes | SNPs | Genes | SNPs | Genes | SNPs | Genes |
| 1AL | 113 | 13.15% | 79 | 12.29% | 78 | 10.79% | 50 | 9.43% |
| 1AS | 26 | 3.03% | 21 | 3.27% | 20 | 2.77% | 17 | 3.21% |
| 1BL | 157 | 18.28% | 110 | 17.11% | 98 | 13.55% | 64 | 12.08% |
| 1BS | 120 | 13.97% | 74 | 11.51% | 94 | 13.00% | 44 | 8.30% |
| 1DL | 30 | 3.49% | 21 | 3.27% | 58 | 8.02% | 47 | 8.87% |
| 1DS | 40 | 4.66% | 25 | 3.89% | 38 | 5.26% | 24 | 4.53% |
| 2AL | 22 | 2.56% | 20 | 3.11% | 14 | 1.94% | 12 | 2.26% |
| 2AS | 11 | 1.28% | 11 | 1.71% | 10 | 1.38% | 7 | 1.32% |
| 2BL | 17 | 1.98% | 15 | 2.33% | 18 | 2.49% | 17 | 3.21% |
| 2BS | 11 | 1.28% | 10 | 1.56% | 12 | 1.66% | 7 | 1.32% |
| 2DL | 2 | 0.23% | 2 | 0.31% | 15 | 2.07% | 10 | 1.89% |
| 2DS | 0 | 0.00% | 0 | 0.00% | 5 | 0.69% | 3 | 0.57% |
| 3AL | 7 | 0.81% | 7 | 1.09% | 2 | 0.28% | 2 | 0.38% |
| 3AS | 1 | 0.12% | 1 | 0.16% | 4 | 0.55% | 4 | 0.75% |
| 3B | 31 | 3.61% | 26 | 4.04% | 28 | 3.87% | 24 | 4.53% |
| 3BL | 0 | 0.00% | 0 | 0.00% | 0 | 0.00% | 0 | 0.00% |
| 3BS | 0 | 0.00% | 0 | 0.00% | 0 | 0.00% | 0 | 0.00% |
| 3DL | 7 | 0.81% | 6 | 0.93% | 2 | 0.28% | 2 | 0.38% |
| 3DS | 1 | 0.12% | 1 | 0.16% | 2 | 0.28% | 2 | 0.38% |
| 4AL | 18 | 2.10% | 15 | 2.33% | 6 | 0.83% | 6 | 1.13% |
| 4AS | 5 | 0.58% | 5 | 0.78% | 6 | 0.83% | 5 | 0.94% |
| 4BL | 11 | 1.28% | 10 | 1.56% | 6 | 0.83% | 6 | 1.13% |
| 4BS | 6 | 0.70% | 5 | 0.78% | 13 | 1.80% | 10 | 1.89% |
| 4DL | 4 | 0.47% | 4 | 0.62% | 5 | 0.69% | 5 | 0.94% |
| 4DS | 2 | 0.23% | 2 | 0.31% | 5 | 0.69% | 4 | 0.75% |
| 5AL | 7 | 0.81% | 5 | 0.78% | 3 | 0.41% | 3 | 0.57% |
| 5AS | 1 | 0.12% | 1 | 0.16% | 2 | 0.28% | 2 | 0.38% |
| 5BL | 31 | 3.61% | 28 | 4.35% | 14 | 1.94% | 14 | 2.64% |
| 5BS | 7 | 0.81% | 5 | 0.78% | 6 | 0.83% | 5 | 0.94% |
| 5DL | 8 | 0.93% | 7 | 1.09% | 15 | 2.07% | 14 | 2.64% |
| 5DS | 4 | 0.47% | 3 | 0.47% | 6 | 0.83% | 5 | 0.94% |
| 6AL | 22 | 2.56% | 17 | 2.64% | 9 | 1.24% | 7 | 1.32% |
| 6AS | 8 | 0.93% | 8 | 1.24% | 11 | 1.52% | 10 | 1.89% |
| 6BL | 7 | 0.81% | 6 | 0.93% | 3 | 0.41% | 2 | 0.38% |
| 6BS | 7 | 0.81% | 5 | 0.78% | 2 | 0.28% | 2 | 0.38% |
| 6DL | 11 | 1.28% | 10 | 1.56% | 7 | 0.97% | 7 | 1.32% |
| 6DS | 5 | 0.58% | 3 | 0.47% | 2 | 0.28% | 2 | 0.38% |
| 7AL | 9 | 1.05% | 8 | 1.24% | 7 | 0.97% | 6 | 1.13% |
| 7AS | 5 | 0.58% | 5 | 0.78% | 8 | 1.11% | 7 | 1.32% |
| 7BL | 10 | 1.16% | 10 | 1.56% | 4 | 0.55% | 4 | 0.75% |
| 7BS | 3 | 0.35% | 3 | 0.47% | 4 | 0.55% | 4 | 0.75% |
| 7DL | 15 | 1.75% | 10 | 1.56% | 12 | 1.66% | 12 | 2.26% |
| 7DS | 8 | 0.93% | 4 | 0.62% | 6 | 0.83% | 6 | 1.13% |
| Unmapped | 49 | 5.70% | 35 | 5.44% | 63 | 8.71% | 46 | 8.68% |
| Mapped | 810 | 94.30% | 608 | 94.56% | 660 | 91.29% | 484 | 91.32% |
| | | | | | 460 | 53.55% | 358 | 55.68% |
| | | | | | 399 | 46.45% | 285 | 44.32% |
| | | | | | 444 | 61.41% | 341 | 64.34% |
| | | | | | 279 | 38.59% | 189 | 35.66% |

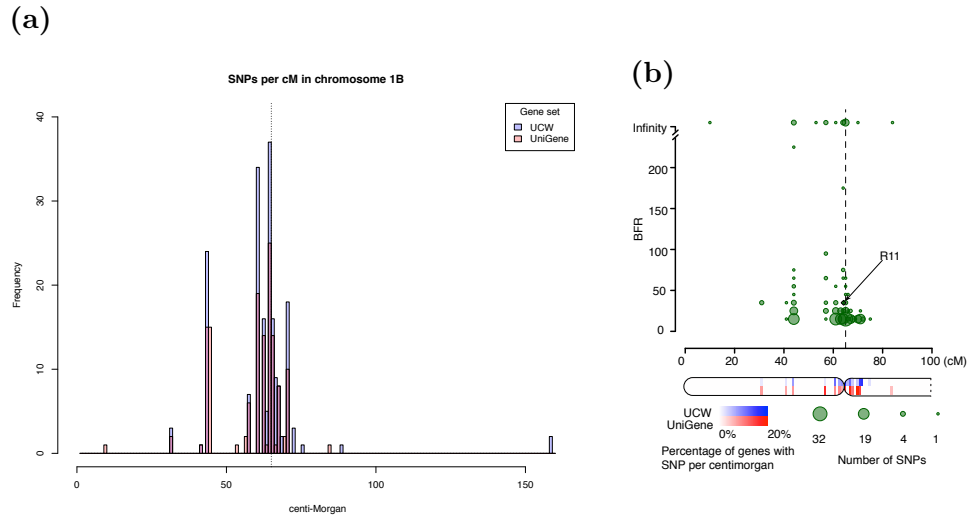


Figure 3.9: (a) Number of SNPs with $BFR > 6$ per cM in chromosome 1B. (b) BFRs of mapped genes with SNPs on chromosome 1B. The area of the circle represents the number of SNPs clustered by location (windows size: 10 cM) and BFR (window size: 5cM). R11 is the only marker near the *Yr15* locus that had a corresponding position in the genetic map. The percentage of genes with SNPs per cM is also illustrated based on UCW (blue) and UniGene (red) gene models. The centromere is imputed by the centre of a window of 10 cM where the short arm switches to the long in the genetic map. BFRs correspond to those from the mixed in silico bulk S1 + S2 + S3/R1 + R2 + R3. Adapted from (Ramirez-Gonzalez et al., 2015c).

smooth curves with a defined peak in the region linked to the studied trait. In practice, we only observe clusters of SNPs with enriched BFRs near the centromere of chromosome 1B (Figures 3.9a, 3.10b).

The location of the clusters with an enrichment of SNPs near the centromere is not expected on a random selection of genes, as the gene density increases with the distance to the centromere (Akhunov et al., 2003). This suggests that the experiment was successful on finding SNPs linked to *Yr15*. There are several factor that prevent a clear peak; like the biases induced by the differential expression, the fragmented reference sequence with scaffolds that are not long enough to go across genetic positions. Since there are several SNPs with a high BFR and the genetic map is not enough to locate a single region linked to *Yr15*, multiple criteria was needed to prioritise SNPs that were more likely to yield on successful genetic markers.

3.6 Assay selection

Three independent criteria were use to prioritize the SNPs for marker development and validation:

High BFR. SNPs with a $BFR > 6$ in at least two independent bulk replicates or in either of the *in silico* mixes were selected to ensure consistency and recover SNPs with a low coverage on a particular bulk.

Group 1S. SNPs that are in CSS scaffolds in the short arm of chromosome group 1 were selected. This is to be consistent with the *in silico* genetic map and with previous studies (Murphy et al., 2009; Peng et al., 2000; Sun et al., 1997).

Yr15 parent. The SNPs should originate from the *Yr15* parent to ensure that the SNP is coming from the *T. diccocooides* introgression and not from a SNP in the AVS genetic background, who would be less useful in breeding programs with a different background.

Only SNPs meeting the three criteria were selected for further analysis.

With the multiple criteria the number of genes with a putative SNP went down from $> 27,000$ to just 175; 77 and 98 from the UniGene and

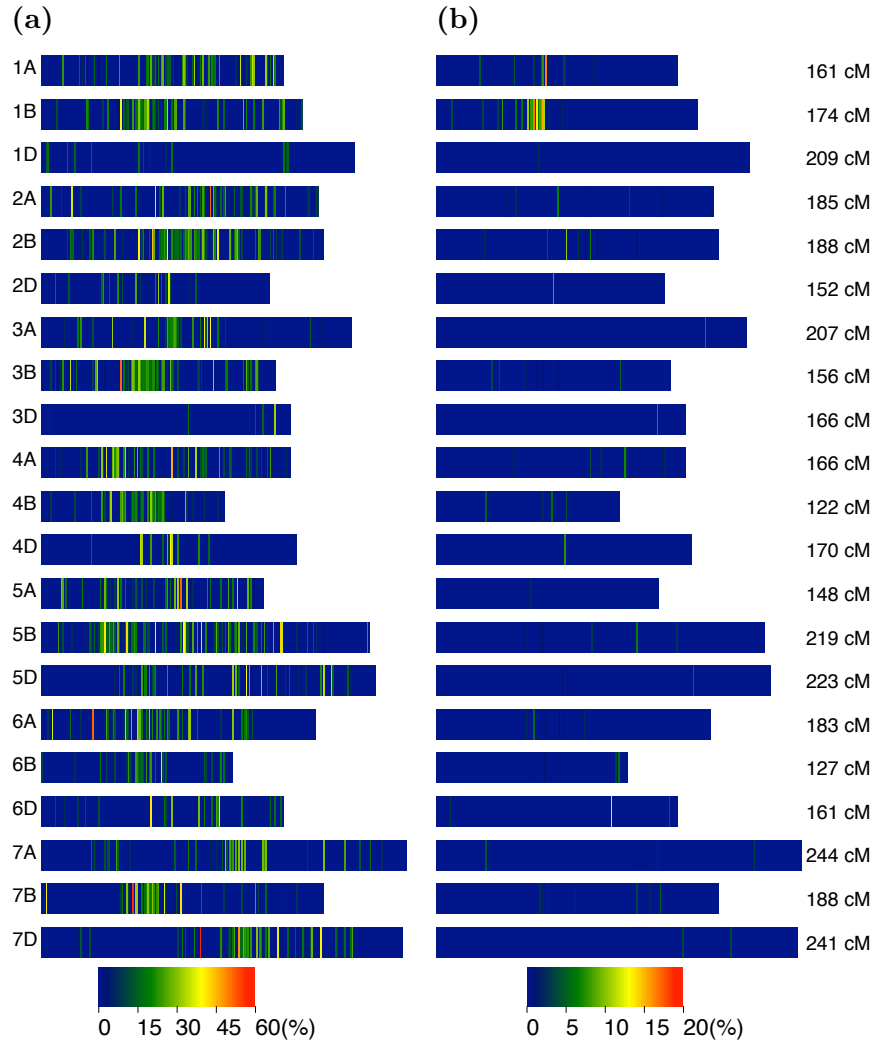


Figure 3.10: Genetic location of genes with SNPs between AVS and Yr15. The colour scale indicates the percentage of genes with SNPs per centi-Morgan (cM) across the 21 wheat chromosomes. The location of the genes was determined by the best alignment to the CSS scaffolds, and the location of these was determined by their position on a genetic map (Wang et al., 2014) (a). All the SNPs between progenitors. Note the lack of enrichment across any individual chromosome. (b) SNPs with BFR > 6. Note the enrichment in Chromosome 1B

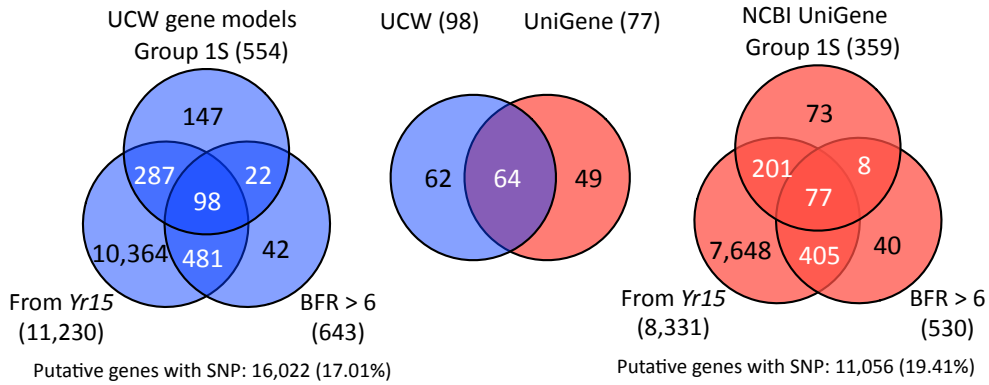


Figure 3.11: Selection criteria for marker design. Venn diagrams based on the three selection criteria (SNP in the short arm of chromosome group 1; SNP has a $BFR > 6$; and SNP is from the *Yr15* parent) for the UCW (blue) and UniGene (red) gene models. The centre diagram shows the intersection between common genes matching all three criteria across both data sets. Note that the numbers are not directly additive as in cases, multiple models from one reference set will relate to a single gene model in the other values. Published in (Ramirez-Gonzalez et al., 2015c)

UCW gene sets respectively. The selected genes from both references were aligned between references, as they come from independent sources an overlap in the selection between them is expected and, as expected, around half of the genes between gene sets overlap (Figure 3.11). The 50 SNPs with the highest BFRs, out of the 175 genes, were selected for validation, 15 of them were redundant between references, resulting on 35 SNPs to validate.

The separate bulks and the *in silico* mixes were evaluated in detail to understand the behaviour and value of having multiple bulks. The initial expectation was that as the number of SNPs with $BFR = \infty$ should drop in the mixes, as the improved coverage should reduce the instances were the absence of an allele is because of the lack of coverage on a particular sample. However, the opposite happened, the additional coverage in the *in silico* mixes recovered SNPs in genes with a low expression at the time of the sampling (Figure 3.12). Some SNPs were present across all the samples, however the value of the BFR changed depending on the sample(marker R5). On some cases a SNP are missing in an individual bulk, but present in the rest of them and in the mixes (marker R8). The main reason affecting the scoring is the coverage in the sample for each particular gene, hence an strategy with a consistent coverage would be preferred for this kind of analysis. Previous studies have shown that a

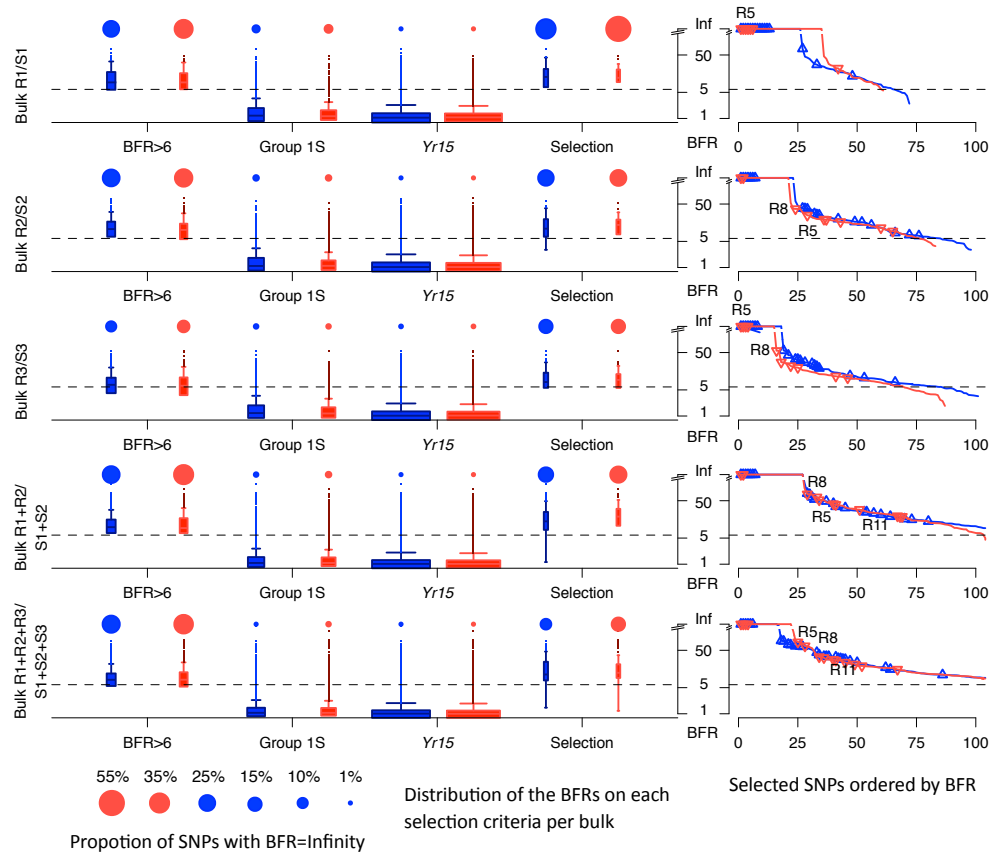


Figure 3.12: Bulk frequency ratio (BFR) of selected SNPs across the individual bulks and in silico mixes (UCW, red; UniGene, blue). The dotted line represents the BFR threshold of 6 (logarithmic scale). Left: Distribution of the BFRs for each selection criteria and the selected SNPs for validation. The circles on the top of each plot represent the percentage of SNPs with $BFR = \infty$. The Selection may include SNPs with $BFR < 6$ when the same SNP has a higher score on the complementing reference (ie. $BFR > 6$ on UCW, but $BFR < 6$ on UniGenes). Right: The BFR values of selected SNPs were sorted in descending order across the different bulks and according to their origin. Validated SNPs are indicated by open triangles, and SNPs corresponding to markers R5, R8 and R11 are labelled across different bulks and mixes. Note that some SNPs are below the threshold in a specific bulk as they meet the BFR criteria across others.

Table 3.8: Number of genes (and SNPs) with a unique hit ($> 99\%$ sequence identity) to a single wheat survey sequence scaffold.

| Chromosome 1 | | All SNPs | | BFR >6 | | % BFR >6 | |
|--------------|------------|----------|--------|----------|--------|------------|--------|
| | | SNP | Genes | SNP | Genes | SNPs | Genes |
| UCW | Unique | 5,283 | 1,245 | 311 | 214 | 5.89% | 17.19% |
| | Total | 8,086 | 1,954 | 486 | 330 | 6.01% | 16.89% |
| | Percentage | 65.34% | 63.72% | 63.99% | 64.85% | | |
| UniGene | Unique | 3,687 | 745 | 213 | 139 | 5.78% | 18.66% |
| | Total | 6,422 | 1,318 | 386 | 246 | 6.01% | 18.66% |
| | Percentage | 57.41% | 56.53% | 49.17% | 56.07% | | |
| UCW + | Unique | 8,970 | 1,990 | 524 | 353 | 5.84% | 17.74% |
| | Total | 14,508 | 3,272 | 872 | 576 | 6.01% | 17.60% |
| | Percentage | 61.83% | 60.82% | 60.09% | 61.28% | | |

| All SNPs | | All SNPs | | BFR >6 | | % BFR >6 | |
|----------|------------|----------|--------|----------|--------|------------|-------|
| | | SNP | Genes | SNP | Genes | SNPs | Genes |
| UCW | Unique | 39,247 | 9,585 | 481 | 368 | 1.23% | 3.84% |
| | Total | 66,426 | 16,022 | 859 | 643 | 1.29% | 4.01% |
| | Percentage | 59.08% | 59.82% | 56.00% | 57.23% | | |
| UniGene | Unique | 27,292 | 5,698 | 344 | 252 | 1.26% | 4.42% |
| | Total | 52,262 | 11,056 | 723 | 530 | 1.38% | 4.79% |
| | Percentage | 52.22% | 51.54% | 47.58% | 47.55% | | |
| UCW + | Unique | 66,539 | 15,283 | 825 | 620 | 1.24% | 4.06% |
| | Total | 118,688 | 27,078 | 1,582 | 1,173 | 1.33% | 4.33% |
| | Percentage | 56.06% | 56.44% | 52.15% | 52.86% | | |

coverage of $< 5x$ is enough to call for SNPs in model organisms with a high-quality reference (Schneeberger and Weigel, 2011). However, the results on this study are in line with other studies using populations for SNP calling (Abe et al., 2012; Takagi et al., 2013a). The non-uniform distribution of the coverage in RNA-Seq experiments affects the number of reads that can be used to call for SNPs, specially on genes with a low expression level (Mortazavi et al., 2008).

Around 60% of the gene models, across both references, had a unique hit with $> 99\%$ sequence identity to a single CSS scaffold (Table 3.8). This is likely because there is no unique homoeologue in the gene models, leading to reads, from two different homoeologues, mapping to the same region. To reduce the number of spurious SNPs we used IUAPC ambiguity codes (Section 1.5.1, Cornish-Bowden (1985)) when two different alleles were observed. This had as side effect that in order to keep only high confidence SNPs we required a higher coverage ($> 20x$). On the original study introducing the BFR in tetraploid wheat, the authors show that increasing the coverage, from $8x$ to $16x$, reduces the putative

SNPs by 60%, but the validated SNPs increases from 57% to 83% (Trick et al., 2012). Hence, a compromise between increasing the minimum coverage at the cost of reducing the SNP candidates has to be reached in line with the objectives and available resources for a particular study.

3.7 Genetic map

The three versions of the genetic map: With a subset of the F_2 population

3.8 Bioinformatic methods

3.8.1 Alignment reads to gene models

The raw output from the Illumina HiSeq 2000 was processed with Casava v1.8 (Illumina, 2011). Lanes 1 and 2, containing multiplexed bulks (Table 3.1) was demultiplexed with a tolerance of 1 mismatch in the barcode. Lanes 3 and 4 contained the parental sequences without a barcode. The FastQ files were left compressed and in chunks of 40,000, as the default for the BCL conversion pipeline from Casava to allow parallel processing in a cluster environment. The quality of the sequencing lanes was assessed with FastQC v0.10.1 (Babraham Bioinformatics, 2012). The RNA-Seq reads were aligned with BWA 0.5.9 (Li and Durbin, 2009) to the wheat UniGene database v60 (Pontius et al., 2002) and to the UCW gene models (Krasileva et al., 2013), including the *T. turgidum* and complementary ORFs (MAS Wheat, 2013). The alignments were sorted and stored as single BAM files to have random access (Li et al., 2009).

3.8.2 Bulk Frequency Ratios

Listing 3.1: Method to find best BLAT alignment

```
def self.each_best_hit(text = '')
  emptyHit = Bio::Blat::Report::Hit.new
  emptyHit.score = 0
  best_aln = Hash.new(emptyHit)
  self.each_hit(text) do |hit|
    current_score = hit.score
    if current_score > best_aln[current_name].score
```

```

        best_aln[current_name] = hit
    end
end
best_aln.each_value { |val| yield val }
end

```

Listing 3.2: Extensions to Bio::Blat::Report::Hit to get the percentage of coverage

```

class Bio::Blat::Report::Hit
  def covered
    match + mismatch
  end
  def query_percentage_covered
    covered * 100.0 / query_len.to_f
  end
  def target_percentage_covered
    covered * 100.0 / target_len.to_f
  end
end
end

```

3.8.3 Alignment between gene models

3.9 Discussion

Remarks on how this technique can be used to do fine-mapping and that if I were to start the project now I would use exome capture or Ren-Seq.

The references have changed since we started

When the study started the study the Chinese Spring Chromosome arm survey sequence (CSS) was about to be released and, despite being fragmented, it served as a valuable resource to assign the SNPs to a chromosome arm. Currently the TGACv1 assembly is available in ensembl REF and, the NRGene assembly, which promises to be on the level of pseudomolecules, is about to be released REF.

There are new annotations, now we don't necessarily need to use unigenes anymore.

Importance of genotyping everything used in an experiment

In Silico mapping wasn't enough to locate all the SNPs, but we ciykd have found one of the SNPs.

Importance of uniform coverage across samples, new developments that could help.

The recombination in the population is not enough to make a finer map, even if we had found more markers.

Mention other people using a similar strategy since this was published.

We can use different techniques now (exome capture, ren-seq)

The markers are now used by our collaborators.

Chapter 4

Gene expression (expVIP)

4.1 Expression experiments (Introduction)

Describe the list of previously published expression experiments and how they can potentially be used as a framework for new experiments.

4.2 Database design

Description of how the database was designed and the flexibility given by having the factors and units as variables

4.3 Analysis pipeline

Implementation of the pipeline, from running kallisto to load the data in the database

4.4 Graphical interface

How the expression can be displayed filtered, and sorted

4.5 Conclusions

The use of previously published studies is a valuable resource. Also, mention that despite the fact that there are several expression/gene browsers, none of them allow comparisons between species and don't consider polyploids.

Chapter 5

Conclusions and final remarks

This section wraps up by showing the relationship and importance of a comprehensive approach to data analysis, from the field, genetics, molecular biology and genomics. I will also remark how the technology and the resources have changed in the last 4 years. As at the references used at beginning where superseded during the PhD.

Appendix A

Supplemental tables

Table A.1: Count of KASP assays designed for the 40,267 SNP markers located in the genetic map from Wang et al. (2014). 4,228 assays did not align to the target chromosome. Not designed: Primer3 could not find viable primers flanking the SNP.

| | Homoeologous variant | Varietal SNP | Percentage |
|---------------|-------------------------|-----------------|------------|
| Non-specific | 1,765 | 5,857 | 21.15% |
| Semi-specific | 7,942 | 6,907 | 41.20% |
| Specific | 6,813 | 5,957 | 35.43% |
| Not designed | 242 | 556 | 2.21% |
| Total | 16,762 | 19,277 | 36,039 |

A.1 PolyMarker supplemental tables.

Table A.2: PolyMarker used to genotype PST

| Assay | Contig | Position | X | Y | Cluster I isolates | | Cluster II isolates | | Cluster III isolates | | Cluster IV isolates | |
|-------|--------------|----------|---|---|--------------------|--------|---------------------|--------|----------------------|-------|---------------------|-------|
| | | | | | 13/26 | 13/123 | CL1 | T-13/3 | 13/09 | 13/23 | 13/182 | 13/36 |
| 1 | PST130.14470 | 268 | C | T | X:Y | X:X | X:X | X:X | X:X | X:X | X:X | X:X |
| 2 | PST130.8160 | 11876 | C | T | Y:Y | X:Y | X:Y | X:Y | X:Y | X:Y | X:Y | X:Y |
| 3 | PST130.14628 | 1712 | A | C | X:Y | - | X:X | X:X | X:X | X:X | X:X | X:X |
| 4 | PST130.14898 | 503 | G | A | X:X | X:X | X:Y | X:Y | X:Y | - | X:Y | X:Y |
| 5 | PST130.28344 | 2372 | A | G | Y:Y | Y:Y | X:Y | X:Y | Y:Y | Y:Y | Y:Y | Y:Y |
| 6 | PST130.7634 | 3463 | A | C | Y:Y | Y:Y | X:Y | X:Y | Y:Y | Y:Y | Y:Y | Y:Y |
| 7 | PST130.7629 | 11699 | G | A | Y:Y | Y:Y | X:Y | X:Y | Y:Y | Y:Y | Y:Y | Y:Y |
| 8 | PST130.10943 | 2979 | C | T | X:Y | X:Y | X:Y | X:Y | X:X | X:X | X:Y | X:Y |
| 9 | PST130.10126 | 6216 | G | T | Y:Y | Y:Y | X:X | X:X | X:X | X:X | Y:Y | Y:Y |
| 10 | PST130.22010 | 172 | C | T | Y:Y | Y:Y | Y:Y | Y:Y | X:Y | X:Y | X:Y | X:Y |
| 11 | PST130.16961 | 1098 | C | T | X:X | X:X | X:Y | X:Y | Y:Y | Y:Y | X:Y | X:Y |
| 12 | PST130.6915 | 2710 | A | T | Y:Y | Y:Y | Y:Y | Y:Y | X:Y | X:Y | Y:Y | Y:Y |
| 13 | PST130.12479 | 1428 | C | T | X:X | X:X | Y:Y | Y:Y | X:X | X:X | Y:Y | X:X |
| 14 | PST130.7634 | 3883 | C | G | X:X | X:X | X:Y | X:Y | X:X | X:X | X:Y | X:X |
| 15 | PST130.14470 | 456 | T | C | Y:Y | Y:Y | X:Y | X:Y | Y:Y | Y:Y | Y:Y | Y:Y |

Table A.3: Validation of homozygous deletions on line Cadenza0423.

| Marker | Deletion | chr | cM | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | C | C | C | C | Result |
|--|----------|-----|--------|---|---|---|---|---|---|---|---|---|----|----|----|---|---|---|---|--------------|
| 5BS_2297308_Cadenza0423_12664_C12664T | - | 5B | 4.551 | X | X | - | X | X | X | X | X | X | X | - | X | Y | Y | Y | Y | HOM Mutation |
| 5BL_10812849_Cadenza0423_5664_G5664T | - | 5B | 38.769 | X | X | - | X | X | X | X | X | X | X | - | X | Y | Y | Y | Y | HOM Mutation |
| 5BL_10825062_Cadenza0423_7917_G7917A | - | 5B | 38.769 | X | X | - | X | X | X | X | X | X | X | - | X | Y | Y | Y | Y | HOM Mutation |
| IWGSC_CSS_5BL_scaff_10847976:27068-27231 | + | 5B | 38.769 | X | X | - | X | X | X | X | X | X | X | - | X | H | H | H | H | Hom Deletion |
| IWGSC_CSS_5BL_scaff_10847976:28118-28674 | + | 5B | 38.769 | X | X | - | X | X | X | X | X | X | X | - | X | H | H | H | H | Hom Deletion |
| IWGSC_CSS_5BL_scaff_10865441:15863-15946 | + | 5B | 38.769 | X | X | - | X | X | X | X | X | X | X | - | X | H | H | H | H | Hom Deletion |
| 5BL_10837222_Cadenza0423_4616_G4616A | - | 5B | 39.905 | X | X | - | X | X | X | X | X | X | X | - | X | Y | Y | Y | Y | HOM Mutation |
| 5BL_10891320_Cadenza0423_18847_C18847T | - | 5B | 45.594 | Y | Y | - | Y | H | X | X | Y | H | Y | - | H | Y | Y | Y | Y | HET Mutation |

Table A.4: Validation of mutations on M_4 on Cadenza

| IWGSC contig | Line | Pos | WT | Mut | Predicted | M_4 | Primer 1 (Cadenza) | Primer 2 (mutant) | Common Primer |
|------------------------------|-------------|-------|----|-----|-----------|-------|---------------------------|---------------------------|---------------------------|
| IWGSC_CSS_3B_scaff_10445294 | Cadenza1772 | 6019 | C | T | het | het | caggatAgtGggactgtcaaaG | caggatAgtGggactgtcaaaA | ggagacGGctGtggacatT |
| IWGSC_CSS_3DL_scaff_6955403 | Cadenza1772 | 2418 | C | T | het* | hom | tcagCggattgtcgggatG | tcagCggattgtcgggatA | tgteCatgaaTcttgtccacG |
| IWGSC_CSS_4AL_scaff_7106846 | Cadenza1772 | 11277 | G | A | hom | hom | tgggatccatgcctacactG | tgggatccatgcctacactA | gatggTgatttgcgcctA |
| IWGSC_CSS_4AS_scaff_5991335 | Cadenza1772 | 15710 | G | A | hom | hom | ctggccctgcgctgctaC | ctggccctgcgctgctaT | gtggaaGttcagaaggaccaG |
| IWGSC_CSS_4BS_scaff_4956646 | Cadenza1772 | 252 | G | A | het* | hom | gcaggttgacttcccgaG | gcaggttgacttcccgaA | tGaggtacgaGcTaaagAaagC |
| IWGSC_CSS_4DS_scaff_1715962 | Cadenza1772 | 1225 | G | A | hom | hom | cagctgtggTatctcaactgG | cagctgtggTatctcaactgA | CcCtGaaACACcGtttggAT |
| IWGSC_CSS_5AL_scaff_2763407 | Cadenza1772 | 2119 | G | A | hom | hom | gcgacGaacctcgagatctG | gcgacGaacctcgagatctA | gaTggcaAtcgtCgtgcA |
| IWGSC_CSS_5AS_scaff_1548786 | Cadenza1772 | 12625 | C | T | het | het | AtaggcacattgctagactgaG | AtaggcacattgctagactgaA | ggattgggtgttcacgC |
| IWGSC_CSS_5BL_scaff_10849226 | Cadenza1772 | 2289 | C | T | het* | hom | cctgacatcattgttcacgatC | cctgacatcattgttcacgatT | cactccgaggtgtccatgaT |
| IWGSC_CSS_5BS_scaff_2270737 | Cadenza1772 | 2262 | G | A | hom | — | attcCTgtgttggtggCaaatgaG | attcCTgtgttggtggCaaatgaA | taaGcaciaAccctccagctgG |
| IWGSC_CSS_1AL_scaff_3022915 | Cadenza1661 | 891 | C | T | hom | hom | ccacagtgcgactcctattgaCG | ccacagtgcgactcctattgaCA | atgtctgattcGtcGtagtcC |
| IWGSC_CSS_1AS_scaff_3297240 | Cadenza1661 | 1970 | C | T | het | het | catcccgccGtttctcC | catcccgccGtttctcT | gtctcgccgatgaagagcT |
| IWGSC_CSS_1BL_scaff_3828996 | Cadenza1661 | 1340 | G | A | hom | hom | agccggatgttagtgttaacC | agccggatgttagtgttaacT | agcagcttgTcgcgttaaC |
| IWGSC_CSS_1DS_scaff_1884529 | Cadenza1661 | 10575 | G | A | hom | hom | aCagatacaAttgtcatgcaggC | aCagatacaAttgtcatgcaggT | acctgggTTgtccaatacttC |
| IWGSC_CSS_2AL_scaff_6318370 | Cadenza1661 | 19142 | C | T | het | — | cgtggcCgaatCtcGacG | cgtggcCgaatCtcGacA | ttcttggggagccgggC |
| IWGSC_CSS_2AS_scaff_5213460 | Cadenza1661 | 1358 | G | A | hom | hom | gtcacgaaCccgctcagA | gtcacgaaCccgctcagA | aggaaagagagaaaagaGcG |
| IWGSC_CSS_2BS_scaff_5179331 | Cadenza1661 | 5604 | G | A | het | het | actctcgtcaagaactgatacaG | actctcgtcaagaactgatacaA | gcaGagaatgttcttgaacT |
| IWGSC_CSS_2DS_scaff_5341235 | Cadenza1661 | 4673 | G | A | het | het | ggtaggagatctcggagctG | ggtaggagatctcggagctA | gcgcggtcgtacaggttG |
| IWGSC_CSS_3AL_scaff_4250995 | Cadenza1661 | 7046 | G | A | hom | hom | cCaagaaacgggtgggtccaG | cCaagaaacgggtgggtccaA | ctgcagctgtccatcatcgT |
| IWGSC_CSS_3B_scaff_10404421 | Cadenza1661 | 4303 | G | A | het | het | ccttcgtcgaCaggacctG | ccttcgtcgaCaggacctA | GCcagtaactCacAtgtctC |
| IWGSC_CSS_5DL_scaff_2390496 | Cadenza1538 | 2125 | C | T | hom | het | gcagttttatcctcagtagtcttgG | gcagttttatcctcagtagtcttgA | ttctgagaaTgtaagtgtcGatG |
| IWGSC_CSS_6AL_scaff_5753680 | Cadenza1538 | 3920 | C | T | hom | hom | tgctccaaatttgagcacaaTaaC | tgctccaaatttgagcacaaTaaT | aaatgcaaggggtaagtttttG |
| IWGSC_CSS_6AS_scaff_4425792 | Cadenza1538 | 4307 | G | A | hom | het | agatgcttgtCggGccaG | agatgcttgtCggGccaA | gctgaagcaacgcgatcaaT |
| IWGSC_CSS_6BS_scaff_3003630 | Cadenza1538 | 6933 | C | T | het | het | ggcagtaagtgtgtgctgagC | ggcagtaagtgtgtgctgagT | tTgaCttctggttgggtggcA |
| IWGSC_CSS_6DL_scaff_3246988 | Cadenza1538 | 9186 | G | A | het | het | gctaaagaagagcttgagagaattC | gctaaagaagagcttgagagaattT | aattttctgaagagaggtgtgtatG |
| IWGSC_CSS_7AL_scaff_4480114 | Cadenza1538 | 3446 | C | T | het | — | gatatctcccacacggcgG | gatatctcccacacggcgA | tgagccactcttcgagtttT |
| IWGSC_CSS_7AS_scaff_4193541 | Cadenza1538 | 8359 | C | T | hom | het | agcaattctttggctatcaattagC | agcaattctttggctatcaattagT | tcactGTcttaactctactgctG |
| IWGSC_CSS_7BL_scaff_6721572 | Cadenza1538 | 9223 | C | T | het | het | gctCaggaggagagacaagaaG | gctCaggaggagagacaagaaA | tgctatgaagaattccgacctC |
| IWGSC_CSS_7BS_scaff_3152545 | Cadenza1538 | 3960 | G | A | hom | — | tcagcaaaatcacctgcCgC | tcagcaaaatcacctgcCgT | gCtgccccatcatcgtttaT |
| IWGSC_CSS_7DS_scaff_3963838 | Cadenza1538 | 2913 | G | A | het | het | tCgttgcaagcCttTtgtgT | tCgttgcaagcCttTtgtgT | agaGttaTcaageTactgtcacA |
| IWGSC_CSS_1AL_scaff_3903380 | Cadenza1469 | 6193 | G | A | hom | hom | ctcttcAgagatgaacgcgA | ctcttcAgagatgaacgcgA | tcGtGagatGTgggttGTTA |
| IWGSC_CSS_1AS_scaff_3287728 | Cadenza1469 | 3817 | C | T | het* | hom | ccgaccaAttcactaacccG | ccgaccaAttcactaacccA | accctctttcccAgacatgaT |
| IWGSC_CSS_1BL_scaff_3815304 | Cadenza1469 | 513 | G | A | hom | hom | aacatttgctTaCcaaaacGC | aacatttgctTaCcaaaacGT | acacagcaagttataatgCAAAGC |
| IWGSC_CSS_1DL_scaff_2266648 | Cadenza1469 | 5926 | C | T | het | het | caacatgagacacacaccttC | caacatgagacacacaccttT | gtcaacgcgtgaggattgtC |
| IWGSC_CSS_1DS_scaff_1906671 | Cadenza1469 | 3697 | C | T | hom | hom | tggTGTgtagacacttggcgA | tggTGTgtagacacttggcgA | catggcgaccaccAcctG |
| IWGSC_CSS_2AL_scaff_6337088 | Cadenza1469 | 7334 | G | A | het* | hom | acaatgccAagttgacaggttG | acaatgccAagttgacaggttA | gggagtggtgttCagaacaT |

| IWGSC contig | Line | Pos | WT | Mut | Predicted | M_4 | Primer 1 (Cadenza) | Primer 2 (mutant) | Common Primer |
|------------------------------|-------------|-------|----|-----|-----------|-------|---------------------------|---------------------------|--------------------------|
| IWGSC_CSS_2BL_scaff.7972799 | Cadenza1469 | 8995 | C | T | het | hom | gTgCtcctcGgcaccccttC | gTgCtcctcGgcaccccttT | gatccgGgcaaacacgTG |
| IWGSC_CSS_2DL_scaff.9832343 | Cadenza1469 | 3262 | G | A | het | het | TtgtctaAcagcacCGcagG | TtgtctaAcagcacCGcagA | agatctcggtcagcctttcT |
| IWGSC_CSS_2DS_scaff.5327939 | Cadenza1469 | 3889 | G | A | het | het | ttttTgccttatgtgactctagtaC | ttttTgccttatgtgactctagtaT | gaggccatcacagatagcG |
| IWGSC_CSS_3B_scaff.10395219 | Cadenza1469 | 1292 | G | A | hom | — | agggtccttgctgcttgctgG | agggtccttgctgcttgctgA | cctctctgggggctttataC |
| IWGSC_CSS_3B_scaff.10592217 | Cadenza0580 | 2994 | C | T | het | — | acagcagtatcaagccctcC | acagcagtatcaagccctcT | tgatactgttgTggCggagG |
| IWGSC_CSS_3DS_scaff.2596771 | Cadenza0580 | 1037 | G | A | het | het | tggttatgCAcaggataatCagG | tggttatgCAcaggataatCagA | tggcaaatgtgatgtcattaggT |
| IWGSC_CSS_4AL_scaff.7093953 | Cadenza0580 | 9881 | C | T | hom | hom | GacaggaagccggtaacaC | GacaggaagccggtaacaT | ctccAgcaggcatgggaT |
| IWGSC_CSS_4BL_scaff.7037448 | Cadenza0580 | 1837 | C | T | hom | hom | CgttgaaaaGctgcaagaacttaaC | CgttgaaaaGctgcaagaacttaaT | cagttcttccTtCaGagcagataT |
| IWGSC_CSS_4BS_scaff.4929479 | Cadenza0580 | 10668 | G | A | hom | — | tggattttcccgactgttC | tggattttcccgactgttT | gtaaacaggcatttcaagagtcA |
| IWGSC_CSS_4DL_scaff.14359838 | Cadenza0580 | 1408 | G | A | hom | — | gCtcAttcagggatTGTcCtaTatG | gCtcAttcagggatTGTcCtaTatA | tgaCagaacagttggctacatC |
| IWGSC_CSS_4DS_scaff.2276484 | Cadenza0580 | 8034 | G | A | hom | hom | gccgtggttgatggAgaG | gccgtggttgatggAgaA | cgtccagattactgatacttgcA |
| IWGSC_CSS_5AL_scaff.2756579 | Cadenza0580 | 5278 | G | A | het | het | tgaatggatttttctgccgttC | tgaatggatttttctgccgttT | ggAAatCCTATgCagaAgAaaCTG |
| IWGSC_CSS_5BL_scaff.10787208 | Cadenza0580 | 10627 | G | A | het | — | gcctctcacatcgcgagaC | gcctctcacatcgcgagaT | acgatgtcAggtggGcgT |
| IWGSC_CSS_5BS_scaff.2282179 | Cadenza0580 | 5267 | G | A | het | — | tgatgggctacgacgtgC | tgatgggctacgacgtgT | tcggcgcccttgaaAtcC |
| IWGSC_CSS_5DL_scaff.4498073 | Cadenza0423 | 4937 | C | T | hom | hom | gcaccctctggttggtcatC | gcaccctctggttggtcatT | tgacagcaAagcagccG |
| IWGSC_CSS_5DS_scaff.2738970 | Cadenza0423 | 2319 | C | T | het | — | cgtgaggtgggtgatttG | cgtgaggtgggtgatttT | tggaaactagtacactgcagtTC |
| IWGSC_CSS_6AL_scaff.5757109 | Cadenza0423 | 2788 | G | A | hom | hom | caggaGcctggcaataaaGG | caggaGcctggcaataaaGA | ctttcGagtcctcttagtttcG |
| IWGSC_CSS_6AS_scaff.4387871 | Cadenza0423 | 2543 | G | A | hom | hom | gcatgctaacaggcgaaaaG | gcatgctaacaggcgaaaaA | ctcatgctcctgatcttaaggtC |
| IWGSC_CSS_6BL_scaff.4271391 | Cadenza0423 | 4660 | C | T | hom | hom | tacgtgcatgatgtggtagtctgaC | tacgtgcatgatgtggtagtctgaT | gtttgaaagtgcacagatgTaccA |
| IWGSC_CSS_6DS_scaff.1880206 | Cadenza0423 | 9159 | G | A | het | het | ctgCgaaggctccacaaG | ctgCgaaggctccacaaA | ggatgagaagtttgcattgctC |
| IWGSC_CSS_7AS_scaff.4227506 | Cadenza0423 | 952 | G | A | het | — | ccatgtgtttccaatgttagagC | ccatgtgtttccaatgttagagT | tgccctagctggtatgcT |
| IWGSC_CSS_7BL_scaff.6681782 | Cadenza0423 | 1486 | C | T | hom | hom | agtaagCGtgacagcaatggG | agtaagCGtgacagcaatggA | AtgtctTtgGtggaagtacatcA |
| IWGSC_CSS_7BS_scaff.3160328 | Cadenza0423 | 7801 | C | T | het | het | tgttaaatGatacagCctgcagC | tgttaaatGatacagCctgcagT | tggaaatgggtCgttgtttT |
| IWGSC_CSS_7DS_scaff.407428 | Cadenza0423 | 2051 | G | A | het | het | gtcGCgccatcctgacaG | gtcGCgccatcctgacaA | actcatcAggtcagcccaA |
| IWGSC_CSS_3AL_scaff.442479 | Cadenza0364 | 3198 | C | T | het | het | gagtcaTtaagttggtaagattggC | gagtcaTtaagttggtaagattggT | GCaGaTaaCaacaggatcacG |
| IWGSC_CSS_3AL_scaff.4447942 | Cadenza0364 | 11917 | G | A | het | het | gtcataaaagattgctcctgtgaaG | gtcataaaagattgctcctgtgaaA | ctcGgatgtgggaggaagA |
| IWGSC_CSS_3AS_scaff.1557483 | Cadenza0364 | 2547 | C | T | het | het | aaagtcacatcatgcttaccataaG | aaagtcacatcatgcttaccataaA | cgaataccaacgcctcatcA |
| IWGSC_CSS_3AS_scaff.2648747 | Cadenza0364 | 2688 | G | A | het | het | tggAagcAcaaggggccC | tggAagcAcaaggggccT | GccgccgatggagactcG |
| IWGSC_CSS_3AS_scaff.3304956 | Cadenza0364 | 1017 | G | A | het | het | gtcccttgacacagctttG | gtcccttgacacagctttA | cctgctggactacaactcaaT |
| IWGSC_CSS_3AS_scaff.3321091 | Cadenza0364 | 4585 | C | T | het | het | caagaatgATgctgatgttggaG | caagaatgATgctgatgttggaA | acatgctgaatgccgaatC |
| IWGSC_CSS_3AS_scaff.3371333 | Cadenza0364 | 538 | G | A | het | het | gggaaaCgAgAcgagcgG | gggaaaCgAgAcgagcgA | ccgtgcttctctaccctT |
| IWGSC_CSS_3AS_scaff.3371815 | Cadenza0364 | 1061 | C | T | het | het | atccccacggcacagagG | atccccacggcacagagA | aAttggcccttggtgattcC |
| IWGSC_CSS_3AS_scaff.3440912 | Cadenza0364 | 4498 | G | A | het | het | ccgtaaaaactttctgtgcttgC | ccgtaaaaactttctgtgcttgT | atActgacaaactacatgatgtgC |
| IWGSC_CSS_3B_scaff.10343586 | Cadenza0364 | 2242 | G | A | het | — | gggttcTgTcctctcttccactG | gggttcTgTcctctcttccactA | tgtgttgaaccgcgaagcA |
| IWGSC_CSS_3AL_scaff.442479 | Cadenza0364 | 3198 | C | T | het | het | gagtcaTtaagttggtaagattggC | gagtcaTtaagttggtaagattggT | GCaGaTaaCaacaggatcacG |
| IWGSC_CSS_3AL_scaff.4447942 | Cadenza0364 | 11917 | G | A | het | het | gtcataaaagattgctcctgtgaaG | gtcataaaagattgctcctgtgaaA | ctcGgatgtgggaggaagA |
| IWGSC_CSS_3AS_scaff.1557483 | Cadenza0364 | 2547 | C | T | het | het | aaagtcacatcatgcttaccataaG | aaagtcacatcatgcttaccataaA | cgaataccaacgcctcatcA |
| IWGSC_CSS_3AS_scaff.2648747 | Cadenza0364 | 2688 | G | A | het | het | tggAagcAcaaggggccC | tggAagcAcaaggggccT | GccgccgatggagactcG |
| IWGSC_CSS_3AS_scaff.3304956 | Cadenza0364 | 1017 | G | A | het | het | gtcccttgacacagctttG | gtcccttgacacagctttA | cctgctggactacaactcaaT |
| IWGSC_CSS_3AS_scaff.3321091 | Cadenza0364 | 4585 | C | T | het | het | caagaatgATgctgatgttggaG | caagaatgATgctgatgttggaA | acatgctgaatgccgaatC |

| IWGSC contig | Line | Pos | WT | Mut | Predicted | M_4 | Primer 1 (Cadenza) | Primer 2 (mutant) | Common Primer |
|------------------------------|-------------|-------|----|-----|-----------|-------|---------------------------|---------------------------|----------------------------|
| IWGSC_CSS_3AS_scaff_3371333 | Cadenza0364 | 538 | G | A | het | het | gggaaaCgAgAcgagcgG | gggaaaCgAgAcgagcgA | cctgtccttctcaccT |
| IWGSC_CSS_3AS_scaff_3371815 | Cadenza0364 | 1061 | C | T | het | het | atccccacggcacagagG | atccccacggcacagagA | aAttggcccttggtgattC |
| IWGSC_CSS_3AS_scaff_3440912 | Cadenza0364 | 4498 | G | A | het | het | ccgtaaaaactttctgtgcttG | ccgtaaaaactttctgtgcttG | atActgacaaactacatgatgtG |
| IWGSC_CSS_3B_scaff_10343586 | Cadenza0364 | 2242 | G | A | het | — | ggttcTgTcctctcttccactG | ggttcTgTcctctcttccactA | tgtgttgaaccgcaagcA |
| IWGSC_CSS_5DL_scaff_242342 | Cadenza0281 | 2433 | C | T | hom | hom | catggCgacggtGtctG | catggCgacggtGtctA | aAccctcatTTtggCTACTtCT |
| IWGSC_CSS_5DL_scaff_4538822 | Cadenza0281 | 1208 | G | A | hom | — | acgtcagaacaaccgtttgaC | acgtcagaacaaccgtttgaT | ttaaattggttggcgccacC |
| IWGSC_CSS_6AL_scaff_5813297 | Cadenza0281 | 4532 | C | T | hom | — | gggagaggggacgtctcgG | gggagaggggacgtctcgA | ttctctcgcaacgattccG |
| IWGSC_CSS_6AS_scaff_4378990 | Cadenza0281 | 6748 | C | T | hom | hom | cccaggttctgctcttttcC | cccaggttctgctcttttcT | caagtatcaaaaaatgaaggGgT |
| IWGSC_CSS_6BL_scaff_4360781 | Cadenza0281 | 5426 | C | T | het | het | aCtactcaaatggcttGgtgtaG | aCtactcaaatggcttGgtgtaA | tcagtccaacatgTcaagagatT |
| IWGSC_CSS_7AL_scaff_4488310 | Cadenza0281 | 3808 | G | A | hom | hom | gttctctttagtagcagccG | gttctctttagtagcagccA | ggcgctttcttcggcctA |
| IWGSC_CSS_7BL_scaff_6696509 | Cadenza0281 | 9232 | G | A | het | het | gctctaggGgtggcaaAagG | gctctaggGgtggcaaAagA | ggcttGaGgtcGcagtgT |
| IWGSC_CSS_7BS_scaff_3143575 | Cadenza0281 | 1866 | C | T | het | het | agatgttgagaggcgcttC | agatgttgagaggcgcttT | gcttggAtgtgggcaagtT |
| IWGSC_CSS_7DL_scaff_3346250 | Cadenza0281 | 1663 | G | A | het | het | acgtgcagcaacatcctaaC | acgtgcagcaacatcctaaT | TttcccaccaggccaagA |
| IWGSC_CSS_7DS_scaff_3933917 | Cadenza0281 | 1243 | C | T | het | het | tgCtgagcCttTcaccttgC | tgCtgagcCttTcaccttgT | agaggtttggttccatcGG |
| IWGSC_CSS_3B_scaff_10626860 | Cadenza0148 | 7847 | G | A | het | het | gcagctctgggaaggagA | gcagctctgggaaggagA | gttaatgtacCTcctagctcG |
| IWGSC_CSS_3DL_scaff_6915683 | Cadenza0148 | 6904 | C | T | het | het | cgtcaaCctgtgggcaattG | cgtcaaCctgtgggcaattA | tcatgctcataatgTcatagggT |
| IWGSC_CSS_4AS_scaff_5929057 | Cadenza0148 | 4238 | G | A | hom | hom | gcgcaacgtagCacctacC | gcgcaacgtagCacctacT | ttatctggtgaagtgcacaggtCA |
| IWGSC_CSS_4AS_scaff_5950625 | Cadenza0148 | 10590 | C | T | het | het | agaTattCaaaTcggtggAttggC | agaTattCaaaTcggtggAttggT | cctgCtccctcacgtcC |
| IWGSC_CSS_4AS_scaff_5967119 | Cadenza0148 | 11626 | C | T | hom | hom | cgtGgacaccccgagctG | cgtGgacaccccgagctA | gacgacgcactgcacgaC |
| IWGSC_CSS_4DL_scaff_14455742 | Cadenza0148 | 1946 | C | T | hom | hom | gCctgaggagatcgcgC | gCctgaggagatcgcgT | aaccgGtAaCTGtGgGcA |
| IWGSC_CSS_4DS_scaff_2318993 | Cadenza0148 | 4000 | C | T | hom | hom | tccagtttgacacagattgaatggG | tccagtttgacacagattgaatggA | tgagaTtctgtttctttcacAttG |
| IWGSC_CSS_5AL_scaff_2750707 | Cadenza0148 | 4603 | G | A | het | het | ccttgggtgtagccatttcaagTaG | ccttgggtgtagccatttcaagTaA | ccaggaTgcAgtgcaatatttcaagG |
| IWGSC_CSS_5BL_scaff_10794137 | Cadenza0148 | 9235 | C | T | hom | hom | gaagctgcttctcgcttG | gaagctgcttctcgcttA | agtatcccttccatataagcagtG |
| IWGSC_CSS_5BS_scaff_1646558 | Cadenza0148 | 2916 | C | T | het | het | gccGtacactcacctAtccttG | gccGtacactcacctAtccttA | gcaaTgtccacttAtcatcccT |
| IWGSC_CSS_1AL_scaff_3883106 | Cadenza0110 | 27536 | C | T | het | het | accttccatcactggctgG | accttccatcactggctgA | gtgaagaacaacaggttgaagC |
| IWGSC_CSS_1BL_scaff_3812829 | Cadenza0110 | 10770 | G | A | het* | hom | ccccactccattccagA | ccccactccattccagA | gGatgtgttctgtgctggaA |
| IWGSC_CSS_1DL_scaff_2266648 | Cadenza0110 | 6156 | G | A | het | het | actgcgtggttatgggacC | actgcgtggttatgggacT | ccccatcactgaacacaacA |
| IWGSC_CSS_1DS_scaff_1889435 | Cadenza0110 | 8826 | C | T | hom | hom | aaccatgaattactcggacagG | aaccatgaattactcggacagA | gcctgaagaattgtatcaaaacaG |
| IWGSC_CSS_2AS_scaff_5268634 | Cadenza0110 | 4636 | G | A | het | het | gatccatgtgattggcatgttG | gatccatgtgattggcatgttA | TgctgtTggatagcagttacT |
| IWGSC_CSS_2BL_scaff_7965110 | Cadenza0110 | 15801 | C | T | hom | hom | cattgaagcAtacacAattgcAtaC | cattgaagcAtacacAattgcAtaT | gccagagatccagataaggTttA |
| IWGSC_CSS_2DL_scaff_9852812 | Cadenza0110 | 13788 | G | A | hom | hom | atttttgtatggtctcaatcttcG | atttttgtatggtctcaatcttcT | gaacgtTcattctgtactgtcT |
| IWGSC_CSS_2DS_scaff_5371379 | Cadenza0110 | 2166 | C | T | hom | hom | agacacaaaactagtGatgCG | agacacaaaactagtGatgCT | gctgctgagaatgttTtgtatttG |
| IWGSC_CSS_3AL_scaff_4384278 | Cadenza0110 | 1276 | C | T | het | het | agcTgaactgccccTgtaG | agcTgaactgccccTgtaA | agggaacctCgGtgatgaA |
| IWGSC_CSS_3AS_scaff_3340122 | Cadenza0110 | 1467 | C | T | hom | hom | attcctAgtgttgcggaacatG | attcctAgtgttgcggaacatA | gagaagactagaaagttttcAgcaT |
| IWGSC_CSS_5DL_scaff_4554222 | Cadenza2103 | 6528 | C | T | het* | hom | gctgccctacaaagaaacaaattG | gctgccctacaaagaaacaaattA | aTcccaactatCGaTttgtcataC |
| IWGSC_CSS_6AL_scaff_5833640 | Cadenza2103 | 7346 | C | T | hom | hom | aagaaaagccacaatggtttctC | aagaaaagccacaatggtttctT | aCTctgTcagtgtttcccagC |
| IWGSC_CSS_6AS_scaff_4429974 | Cadenza2103 | 3867 | G | A | hom | hom | GagatgaAttatttgagcatgtggC | GagatgaAttatttgagcatgtggT | ggttccgggtgcataagT |
| IWGSC_CSS_6DL_scaff_3307626 | Cadenza2103 | 4970 | C | T | hom | hom | tgcagatgttgcctgtgtaG | tgcagatgttgcctgtgtaA | tgtagaaggtgattttgtactGtC |
| IWGSC_CSS_6DS_scaff_2059604 | Cadenza2103 | 5224 | G | A | het | — | gctcaatgcatgcTgagtgG | gctcaatgcatgcTgagtgA | tgtcaagtattattttcgtctG |
| IWGSC_CSS_7AL_scaff_4552322 | Cadenza2103 | 1412 | C | T | het | het | gcaaaggcTgatactccaacaG | gcaaaggcTgatactccaacaA | ggcAAGccAgtataaaagtaaGC |

| IWGSC contig | Line | Pos | WT | Mut | Predicted | M_4 | Primer 1 (Cadenza) | Primer 2 (mutant) | Common Primer |
|------------------------------|-------------|-------|----|-----|-----------|-------|----------------------------|----------------------------|---------------------------|
| IWGSC_CSS_7BS_scaff_3147455 | Cadenza2103 | 4607 | G | A | het | — | gcaccttaggatgtgagTtatgC | gcaccttaggatgtgagTtatgT | gcatgtagggtttatttgactgttA |
| IWGSC_CSS_7DL_scaff_3382467 | Cadenza2103 | 3473 | C | T | hom | — | GGTtctgCaGTTTCATAActcatC | GGTtctgCaGTTTCATAActcatT | attgaatcaactgatacGaaGactC |
| IWGSC_CSS_3B_scaff_10457010 | Cadenza0277 | 10599 | G | A | het | het | aaccttggccgcagaacaC | aaccttggccgcagaacaT | actggtcgcacgagaggG |
| IWGSC_CSS_3B_scaff_10593852 | Cadenza0277 | 10124 | C | T | het | het | tgacaggggacgctatacaG | tgacaggggacgctatacaA | gtctaaCTtACattAcccatcagC |
| IWGSC_CSS_3DS_scaff_2583390 | Cadenza0277 | 663 | G | A | hom | hom | actgcactcatacaatActtCtgC | actgcactcatacaatActtCtgT | tcCacctggacagcaagtG |
| IWGSC_CSS_4AL_scaff_7093953 | Cadenza0277 | 10004 | C | T | hom | hom | ccttgattcaatggaTtgTtttgG | ccttgattcaatggaTtgTtttgA | ttcccaaaTaaaaaggaagagC |
| IWGSC_CSS_4AL_scaff_7176064 | Cadenza0277 | 6220 | C | T | het | het | gtgccgtaTtcCgcctgG | gtgccgtaTtcCgcctgA | atgttcgaggggatgggG |
| IWGSC_CSS_4DL_scaff_14122349 | Cadenza0277 | 1010 | C | T | hom | hom | gtcgctgctgCttgtgaG | gtcgctgctgCttgtgaA | ggaacaggcccaaggagG |
| IWGSC_CSS_5AL_scaff_2736916 | Cadenza0277 | 4296 | G | A | het | het | aagaactATgAaaGtaacacacgaC | aagaactATgAaaGtaacacacgaT | ttcGcTttTaagGcAttCtcG |
| IWGSC_CSS_5BL_scaff_10883744 | Cadenza0277 | 2080 | C | T | hom | hom | gcctctttCtgttTagcctcaG | gcctctttCtgttTagcctcaA | cgacaaggtctggtatTgcA |
| IWGSC_CSS_1AL_scaff_3932013 | Cadenza0548 | 11765 | C | T | hom | hom | accgccaaCccaagacaG | accgccaaCccaagacaA | cccatTAcccTgcAacG |
| IWGSC_CSS_1BS_scaff_3417505 | Cadenza0548 | 373 | C | T | het | het | gtggtgaggaGGgtgGaG | gtggtgaggaGGgtgGaA | tggtcgcGccagttgttA |
| IWGSC_CSS_2AS_scaff_5305619 | Cadenza0548 | 2786 | C | T | hom | hom | atacagatgccttAAGtggTtC | atacagatgccttAAGtggTtT | ggaagacaAtGctccagtaC |
| IWGSC_CSS_2AS_scaff_5306489 | Cadenza0548 | 46953 | T | G | het | wt | aggttccatgtccatagaagGT | aggttccatgtccatagaagGG | aggctaTAGactcctgtACAgT |
| IWGSC_CSS_2BL_scaff_7984123 | Cadenza0548 | 11660 | G | A | het | het | cattgtggcatagtaatcagtacaG | cattgtggcatagtaatcagtacaA | aatacattgaggaatacaagccC |
| IWGSC_CSS_2DL_scaff_9907477 | Cadenza0548 | 1363 | C | T | hom | hom | tgcttccctttgccagaaC | tgcttccctttgccagaaT | ggcaaacctgatgtggcatC |
| IWGSC_CSS_2DS_scaff_5330886 | Cadenza0548 | 5449 | G | A | hom | hom | gcattgtccattataactgaacCgtG | gcattgtccattataactgaacCgtA | catgtctcttctctggacC |
| IWGSC_CSS_3AL_scaff_4449951 | Cadenza0548 | 633 | C | T | het | het | tccaaacctaacagtcataactaG | tccaaacctaacagtcataactaA | gtctgcagTGCaatgtgC |
| IWGSC_CSS_3B_scaff_10479889 | Cadenza0097 | 3339 | C | T | hom | — | ttgTttctGgagaagatgcCG | ttgTttctGgagaagatgcCA | ggtgtcattcaAcGgcA |
| IWGSC_CSS_3B_scaff_10562262 | Cadenza0097 | 7819 | C | T | het | het | agaggggtgctatccatAttgG | agaggggtgctatccatAttgA | agcgatgccaaaggcttcC |
| IWGSC_CSS_4AL_scaff_7040796 | Cadenza0097 | 10772 | G | A | hom | hom | acacaacattgccaccagaG | acacaacattgccaccagaA | CAatCgattgtctgtTctcC |
| IWGSC_CSS_4AL_scaff_7063488 | Cadenza0097 | 6360 | C | T | het | het | gcctctcacCttAattgaaagctgC | gcctctcacCttAattgaaagctgT | aggcagtgagatgtgaaagtT |
| IWGSC_CSS_4AL_scaff_7091701 | Cadenza0097 | 5050 | G | A | het | het | catgagcatctgggaggaaaatG | catgagcatctgggaggaaaatA | agcaagggaAtaatgaacggaaA |
| IWGSC_CSS_4DS_scaff_1845841 | Cadenza0097 | 7110 | G | A | hom | hom | aatgTAGctccccatacCgG | aatgTAGctccccatacCgA | actgaacTgcaatcgtTtatggA |
| IWGSC_CSS_5AL_scaff_2767581 | Cadenza0097 | 3737 | G | A | het | het | gagaggtcctcactAtcggC | gagaggtcctcactAtcggT | cgTcatcacaatatattgtcggG |
| IWGSC_CSS_5BL_scaff_10784643 | Cadenza0097 | 1568 | C | T | hom | hom | agaaaTAcatggatggatggaCG | agaaaTAcatggatggatggaCA | catctcCCttcaCgGaaaG |
| IWGSC_CSS_1AL_scaff_3952258 | Cadenza2092 | 8107 | C | T | het | — | tgagtagagaaattgacagtgtgG | tgagtagagaaattgacagtgtgA | tgccaccattgacatgagaG |
| IWGSC_CSS_1BL_scaff_3858008 | Cadenza2092 | 10278 | G | A | hom | hom | tttgagcaggcaggatcgC | tttgagcaggcaggatcgT | actcaggcctatacActattC |
| IWGSC_CSS_1DL_scaff_2265172 | Cadenza2092 | 9094 | C | T | hom | hom | tgcaTGTcatttgttcttatcagC | tgcaTGTcatttgttcttatcagT | agtgtccaacttccGttcatC |
| IWGSC_CSS_2AL_scaff_6435867 | Cadenza2092 | 16201 | G | A | hom | hom | tttctgTactttaacgtcaattgaC | tttctgTactttaacgtcaattgaT | gtgaggatgatgagtgaaagC |
| IWGSC_CSS_2AL_scaff_6439430 | Cadenza2092 | 25101 | C | T | het | — | caagaaagggCagCtCagC | caagaaagggCagCtCagT | tcGttAcTctttcActggtgaA |
| IWGSC_CSS_2DL_scaff_9760848 | Cadenza2092 | 4733 | C | T | het | het | gcaccatgggtctcaggtaC | gcaccatgggtctcaggtaT | tcagtcagtttGCTCgtTCTG |
| IWGSC_CSS_3AL_scaff_4407012 | Cadenza2092 | 2785 | C | T | hom | hom | acatatAgtgttctcatccaccatC | acatatAgtgttctcatccaccatT | acctctcatgttaaataggtttgT |
| IWGSC_CSS_3AS_scaff_3441108 | Cadenza2092 | 541 | G | A | het | het | GtgatgaccttgagacGgaC | GtgatgaccttgagacGgaA | aggcaTgacaaCgcgcaA |
| IWGSC_CSS_3B_scaff_10449827 | Cadenza1551 | 4779 | G | A | hom | hom | ggcaaggtcaagaaacGgtC | ggcaaggtcaagaaacGgtT | aCagaGtgggttagaggcaG |
| IWGSC_CSS_3B_scaff_10550638 | Cadenza1551 | 3250 | C | T | het | het | ctccttcacttgttgcggC | ctccttcacttgttgcggT | gcaacATtTgatactgcaagG |
| IWGSC_CSS_3DL_scaff_6945816 | Cadenza1551 | 589 | C | T | hom | hom | agcatctcacctgcaaCaataC | agcatctcacctgcaaCaataT | TgtgccCTctgaAtattttcaTG |
| IWGSC_CSS_3DL_scaff_6954177 | Cadenza1551 | 3508 | C | T | het | het | tgtagcatcacattaaacttctcG | tgtagcatcacattaaacttctcA | gcttggtataaacCttacgacA |
| IWGSC_CSS_4AS_scaff_5938272 | Cadenza1551 | 19080 | G | A | hom | hom | agAcCccgAtcgccatgG | agAcCccgAtcgccatgA | GggAgatAcaggtaaaActcTtcG |
| IWGSC_CSS_4AS_scaff_5977594 | Cadenza1551 | 11092 | C | T | het | het | gccttgattcggaacaacaaaC | gccttgattcggaacaacaaaT | gcgtctctcagtcctgcA |

| IWGSC contig | Line | Pos | WT | Mut | Predicted | M_4 | Primer 1 (Cadenza) | Primer 2 (mutant) | Common Primer |
|------------------------------|-------------|-------|----|-----|-----------|-------|---------------------------|---------------------------|---------------------------|
| IWGSC_CSS_5AL_scaff_2671035 | Cadenza1551 | 5859 | C | T | het | het | cggatgatattTttagacttcgacG | cggatgatattTttagacttcgacG | ggcagttcagcGacccatT |
| IWGSC_CSS_5BL_scaff_10889480 | Cadenza1551 | 2530 | G | A | hom | hom | gagcttaactcgagatggaG | gagcttaactcgagatggaA | tccatgCAacGccttgG |
| IWGSC_CSS_3B_scaff_10528396 | Cadenza2088 | 8059 | G | A | hom | — | ctttccgctccgtaagcaataG | ctttccgctccgtaagcaataA | gtgcaactgttcaggcctgA |
| IWGSC_CSS_3B_scaff_10637573 | Cadenza2088 | 16815 | G | A | het | het | agcaagcttaccGgtctgC | agcaagcttaccGgtctgT | cgagcAactacgagcagctT |
| IWGSC_CSS_4AL_scaff_7086469 | Cadenza2088 | 6697 | G | A | het | het | gccgtctacttcaacgcG | gccgtctacttcaacgcA | ccaGaggcttgTGCattttT |
| IWGSC_CSS_4AL_scaff_7126302 | Cadenza2088 | 3627 | G | A | hom | hom | gttcaaaaacaagtggtAatttgC | gttcaaaaacaagtggtAatttgT | cacaaggatatgaagcTctctagA |
| IWGSC_CSS_4BL_scaff_7041808 | Cadenza2088 | 10234 | G | A | hom | hom | tcaatggatgagggtgcttC | tcaatggatgagggtgcttT | ccatagcagcatcagccacA |
| IWGSC_CSS_5AL_scaff_2794167 | Cadenza2088 | 13162 | G | A | het | — | agtattcaggacaagcatCttCaG | agtattcaggacaagcatCttCaA | caatgaacctctcgaagaaGaG |
| IWGSC_CSS_5BL_scaff_10889232 | Cadenza2088 | 3885 | G | A | het | het | cTcaaccacaatgggcaAatC | cTcaaccacaatgggcaAatT | tccttcatcaatcatcaattgtgG |
| IWGSC_CSS_5BS_scaff_2267405 | Cadenza2088 | 11113 | C | T | hom | hom | ctttgatgatcctaggcctctTG | ctttgatgatcctaggcctctTA | tgatttggTctggtAgagtttGA |
| IWGSC_CSS_3B_scaff_10475354 | Cadenza1409 | 2203 | G | A | hom | hom | agCGaacaagagGtcaaacG | agCGaacaagagGtcaaacA | ctgaaacacaCtagaCAattAccG |
| IWGSC_CSS_3B_scaff_10674115 | Cadenza1409 | 4555 | C | T | het | het | gcttcagtgcatgccttcaG | gcttcagtgcatgccttcaA | cttcacaccGagataatGtattG |
| IWGSC_CSS_4AL_scaff_7153568 | Cadenza1409 | 13073 | C | T | hom | hom | tccgaccgAtcaaccttgG | tccgaccgAtcaaccttgA | gaccggaactcctcggcC |
| IWGSC_CSS_4DL_scaff_14314966 | Cadenza1409 | 2010 | G | A | het | hom | gtaggccccctctCAGgG | gtaggccccctctCAGgA | cggcgTcaCaAgttgCcT |
| IWGSC_CSS_4DS_scaff_2324074 | Cadenza1409 | 7606 | G | A | het | het | tGcatgaaaatgtgtGcaGaG | tGcatgaaaatgtgtGcaGaA | gggtaAgttcAaaactGaaagtgaG |
| IWGSC_CSS_5AS_scaff_1517889 | Cadenza1409 | 3561 | G | A | het | het | tctcgacatcttcccggtgaC | tctcgacatcttcccggtgaT | gtgccttgaacattgcttattA |
| IWGSC_CSS_5AS_scaff_1523866 | Cadenza1409 | 8054 | G | A | hom | — | ggatgatctaccgcaGgaC | ggatgatctaccgcaGgaT | tctcgagCcTctctcA |
| IWGSC_CSS_5BL_scaff_10917655 | Cadenza1409 | 19073 | G | A | hom | hom | caaatgacatgcaaaagaagttgC | caaatgacatgcaaaagaagttgT | cgcttcatcactacaAaatatgtcT |
| IWGSC_CSS_1AL_scaff_3886649 | Cadenza1599 | 5204 | C | T | het | het | tgatgcccaaccacaatGcC | tgatgcccaaccacaatGcT | ggactgactgtgacatatttaG |
| IWGSC_CSS_1BL_scaff_3810267 | Cadenza1599 | 6634 | C | T | hom | hom | ccCaggaaatgagcacctC | ccCaggaaatgagcacctT | cgcaggcgaagatgtgaTtG |
| IWGSC_CSS_1DL_scaff_2291677 | Cadenza1599 | 12856 | C | T | hom | hom | GgtagacaagtcgccgaG | GgtagacaagtcgccgaA | cctcctccttcaacGCcG |
| IWGSC_CSS_2AL_scaff_6354492 | Cadenza1599 | 7566 | G | A | het | het | gGagaatgcaCAgtAacTtctgG | gGagaatgcaCAgtAacTtctgA | ttccgaagaaccacaTccTG |
| IWGSC_CSS_2AS_scaff_5282937 | Cadenza1599 | 9736 | G | A | het | het | gctgtagattttatagctgctatgC | gctgtagattttatagctgctatgT | cacCagaattgttCactgatttTC |
| IWGSC_CSS_2BL_scaff_7952427 | Cadenza1599 | 19249 | G | A | hom | hom | cgTccctCcttagcacgaC | cgTccctCcttagcacgaT | aTcactcattagcgcgAG |
| IWGSC_CSS_2DL_scaff_9897981 | Cadenza1599 | 5627 | C | T | het | het | cttgggtgctTgattgcttactC | cttgggtgctTgattgcttactT | gTttgctCtctctgactTtgtG |
| IWGSC_CSS_3AL_scaff_4446105 | Cadenza1599 | 1765 | G | A | hom | — | aaatgctttcctaCcgctagtG | aaatgctttcctaCcgctagtA | ttctAgaggcaatagctTatatgcT |

Table A.5: Validation of mutations on M_4 on Kronos

| IWGSC contig | Line | Pos | WT | Mut | Predicted | M_4 | Primer 1 (Kronos) | Primer 2 (mutant) | Common Primer |
|-----------------------------|------------|------|----|-----|-----------|-------|--------------------------|--------------------------|--------------------------|
| IWGSC_CSS_1AS_scaff_3284790 | Kronos3085 | 7449 | G | A | Het | Het | ccacaccttgagcctcgC | ccacaccttgagcctcgT | gtgattttgccaggggagA |
| IWGSC_CSS_1BL_scaff_3897513 | Kronos3085 | 1515 | C | T | Het | Het | gcttcactGggtcctgC | gcttcactGggtcctgT | acAaggactgcttcagaGaC |
| IWGSC_CSS_2AL_scaff_6434745 | Kronos3085 | 3424 | C | T | Het | Het | cctcGgttttgcaaatttctatgC | cctcGgttttgcaaatttctatgT | gGCaaTggcataacaacagatA |
| IWGSC_CSS_3AS_scaff_3408995 | Kronos3085 | 732 | C | T | Het | Het | aggccatttcgaattccgC | aggccatttcgaattccgT | ggTgttaTccagAacctgagTG |
| IWGSC_CSS_3B_scaff_10708748 | Kronos3085 | 2675 | G | A | Het | Het | gttgcatgcttcacccagG | gttgcatgcttcacccagA | gtaacaactctgagttcgtagcaC |
| IWGSC_CSS_4AL_scaff_7132733 | Kronos3085 | 1799 | C | T | Hom | Hom | caccctgtagtgaccctC | caccctgtagtgaccctT | aCcGcctaGaaagaaagcttC |

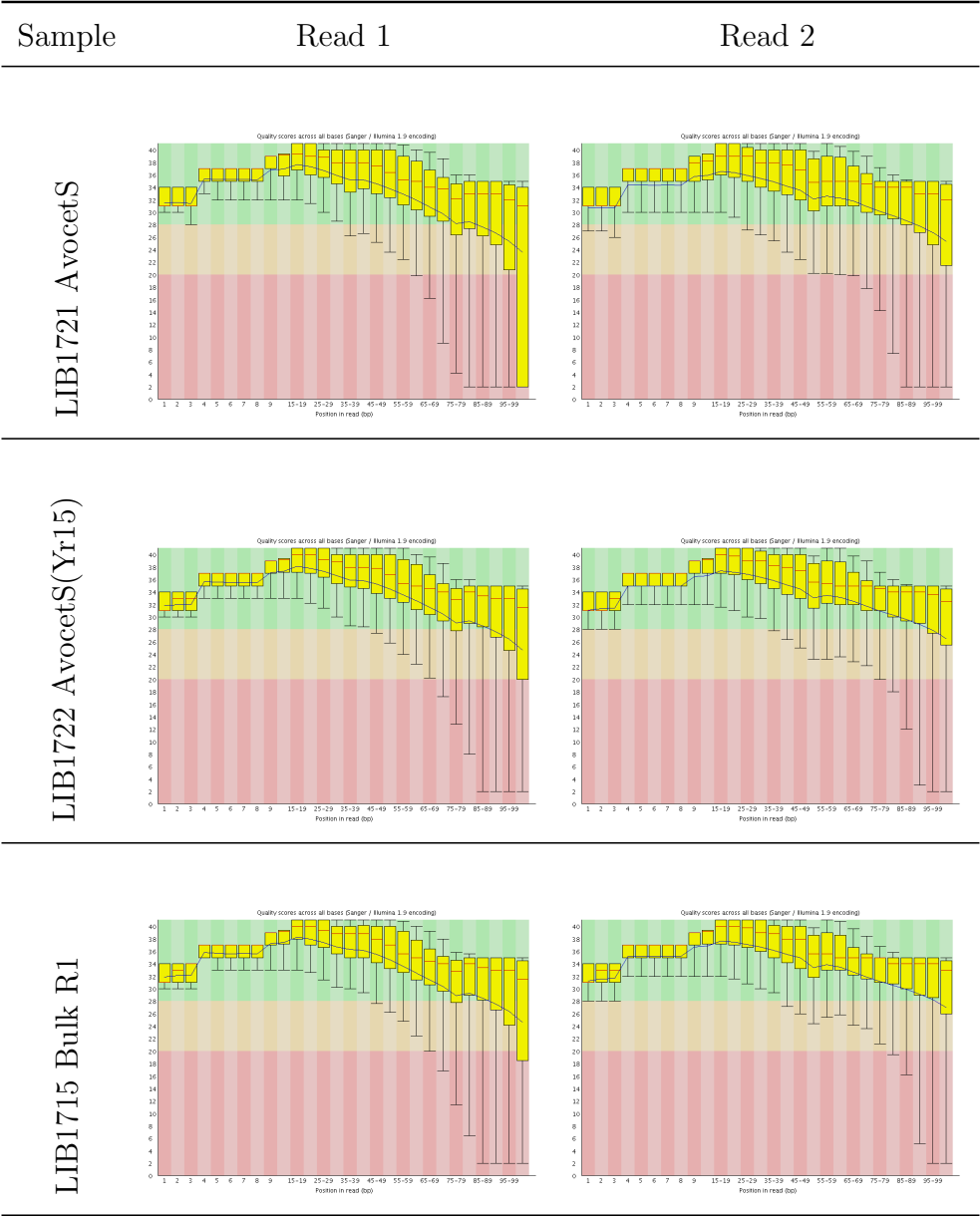
| IWGSC contig | Line | Pos | WT | Mut | Predicted | M_4 | Primer 1 (Kronos) | Primer 2 (mutant) | Common Primer |
|------------------------------|------------|-------|----|-----|-----------|-------|--------------------------|--------------------------|--------------------------|
| IWGSC.CSS_5AS_scaff.1534693 | Kronos3085 | 4605 | C | T | Het | Het | cagcttctggccctcAtC | cagcttctggccctcAtT | gtaCctcagAGtcaTgagAG |
| IWGSC.CSS_6AS_scaff.4361911 | Kronos3085 | 8857 | G | A | Het | Het | tcacgaaagacgacttcaacctcC | tcacgaaagacgacttcaacctcT | catgagggtgctgcatctccatcA |
| IWGSC.CSS_6BS_scaff.3008326 | Kronos3085 | 1528 | G | A | Het | Het | ccatgttgactggtggtgC | ccatgttgactggtggtgT | ggaagcatggCaagtgcA |
| IWGSC.CSS_7AS_scaff.4214385 | Kronos3085 | 27835 | C | T | Hom | Hom | cgtaccttcggtgggaaagG | cgtaccttcggtgggaaagA | ctcttggtcagctgataaagacT |
| IWGSC.CSS_1AL_scaff.3929964 | Kronos3191 | 1336 | C | T | Het | Het | tttcggccataacctgacatC | tttcggccataacctgacatT | attgcctcagcttcttgcaG |
| IWGSC.CSS_1BL_scaff.3899789 | Kronos3191 | 7925 | C | T | Het | Het | actctcacTggcagcagC | actctcacTggcagcagT | caacgtggtgcccatcGtA |
| IWGSC.CSS_2AL_scaff.6426728 | Kronos3191 | 1481 | G | A | Hom | Hom | gaaActgccgcagctCgC | gaaActgccgcagctCgT | ccaGcaGctcgtgagaaA |
| IWGSC.CSS_2BL_scaff.7960273 | Kronos3191 | 690 | C | T | Hom | Hom | gccattcatccttaggcgC | gccattcatccttaggcgT | acatgcaattgctgatgactG |
| IWGSC.CSS_3AS_scaff.3286603 | Kronos3191 | 2975 | G | A | Het* | Hom | ccgtgtggtttgttggtG | ccgtgtggtttgttggtA | gaaaggaacgtgTcaTgcaG |
| IWGSC.CSS_5AL_scaff.2694249 | Kronos3191 | 2399 | C | T | Het | Het | gccttcagatagagccGC | gccttcagatagagccGT | cgccacatcgacattctcG |
| IWGSC.CSS_5BL_scaff.10923577 | Kronos3191 | 3713 | C | T | Het | Het | gtggattgcctgagcttgC | gtggattgcctgagcttgT | tgggtggcctcttgggaC |
| IWGSC.CSS_6AL_scaff.5823017 | Kronos3191 | 13225 | C | T | Hom | Hom | ccctttcagcctctggaG | ccctttcagcctctggaA | ttcgagaaggcccatcgA |
| IWGSC.CSS_6BS_scaff.2955394 | Kronos3191 | 1622 | C | T | Het* | Hom | gtggagatgaaggtctagcaaG | gtggagatgaaggtctagcaaA | gatactcTgcaatgggtgT |
| IWGSC.CSS_7BL_scaff.6739382 | Kronos3191 | 12261 | G | A | Hom | Hom | gagacaagctttgaattgctcC | gagacaagctttgaattgctcT | CgagtgaacctTcatttcccG |
| IWGSC.CSS_1AS_scaff.3276389 | Kronos3288 | 9720 | C | T | Hom | Hom | aCcaGcaggaccAatgtctC | aCcaGcaggaccAatgtctT | atgatgcaacctcagccaT |
| IWGSC.CSS_2AL_scaff.6367515 | Kronos3288 | 6976 | G | A | Het | Het | caggtcgagTgtctccgG | caggtcgagTgtctccgA | ggggtgatCtggaaaggG |
| IWGSC.CSS_2AL_scaff.6422019 | Kronos3288 | 4523 | G | A | Het | Het | cgtaggtccctgcatagG | cgtaggtccctgcatagA | acgcAcgctaagccgtaC |
| IWGSC.CSS_3AL_scaff.4284850 | Kronos3288 | 7901 | C | T | Hom | Hom | tgctttggacaacatcgG | tgctttggacaacatcgA | tgtcAgtcatgcagaccaG |
| IWGSC.CSS_4AS_scaff.5962359 | Kronos3288 | 13049 | G | A | Het | Hom | ccatcaagaagtacgagttcgaC | ccatcaagaagtacgagttcgaT | accatgccagcttgcA |
| IWGSC.CSS_6AL_scaff.5778773 | Kronos3288 | 6853 | G | A | Het | Het | gagtgaccttcccgtcttC | gagtgaccttcccgtcttT | ggagaacagctactcggcT |
| IWGSC.CSS_6AS_scaff.4392100 | Kronos3288 | 3434 | C | T | Het | Het | atggaagcacaggtgaccG | atggaagcacaggtgaccA | ggAagcgaaagtgaacaaA |
| IWGSC.CSS_7BL_scaff.6744240 | Kronos3288 | 9772 | G | A | Het | Het | agctgttcttctcctacttcaaG | agctgttcttctcctacttcaaA | caggtcgttcttgagctcC |
| IWGSC.CSS_1AL_scaff.3887185 | Kronos3413 | 9708 | C | T | Hom | Hom | gcacgcctttatcgaggtaaaG | gcacgcctttatcgaggtaaaA | AgaacagcagagcgcaA |
| IWGSC.CSS_2BS_scaff.3381362 | Kronos3413 | 5160 | C | T | Het* | Hom | caacttctgggctgtagtG | caacttctgggctgtagtA | tgAgaattctgacGcaaaagaC |
| IWGSC.CSS_3AS_scaff.3296605 | Kronos3413 | 6154 | G | A | Het | Het | ctggtcacgggctctagC | ctggtcacgggctctagT | cagcactgagacatggaC |
| IWGSC.CSS_3BL_scaff.10693516 | Kronos3413 | 12632 | C | T | Het | Het | ctaggcttgacaaaacaggC | ctaggcttgacaaaacaggT | agcttgcatctatgggcatT |
| IWGSC.CSS_5AS_scaff.1547699 | Kronos3413 | 2686 | G | A | Het | Het | gCtacaaccttcaccaatcgC | gCtacaaccttcaccaatcgT | gacggctttgaagtgtcatC |
| IWGSC.CSS_5BL_scaff.10856077 | Kronos3413 | 5853 | G | A | Het | Het | agagcttcaccccatgctC | agagcttcaccccatgctT | acgCacatttAatagctgaagC |
| IWGSC.CSS_6AL_scaff.5750718 | Kronos3413 | 11046 | G | A | Hom | Hom | cacgcTtcccgaacttcttataG | cacgcTtcccgaacttcttataA | AgacgatgtgatcaggattcaG |
| IWGSC.CSS_7AL_scaff.4433177 | Kronos3413 | 3511 | C | T | Het | Het | GaTgctccGtcaggctgG | GaTgctccGtcaggctgA | cactactggacaagctcttgG |
| IWGSC.CSS_7BL_scaff.6742567 | Kronos3413 | 667 | C | T | Het | Het | gttgcttgctggcgagaC | gttgcttgctggcgagaT | cattttgcacctgtgtcTG |
| IWGSC.CSS_1AL_scaff.3976389 | Kronos3935 | 10941 | C | T | Hom | Hom | ggtagagagatcggCgatG | ggtagagagatcggCgatA | cagtcactcatagagaggtcaG |
| IWGSC.CSS_1BL_scaff.3873362 | Kronos3935 | 1392 | G | A | Het | Het | cagatctgaagcctaGcacatG | cagatctgaagcctaGcacatA | actaccagaatcagcacaacaaAC |
| IWGSC.CSS_2BL_scaff.7882382 | Kronos3935 | 2721 | C | T | Het | Het | gcaagctaagatgtaccgtagC | gcaagctaagatgtaccgtagT | gccacagttaggagaaagactT |
| IWGSC.CSS_3AL_scaff.4242376 | Kronos3935 | 2410 | C | T | Het | Het | agaacccaaaacccgTacttaG | agaacccaaaacccgTacttaA | gtagGgtCcatCTaaagcttG |
| IWGSC.CSS_3B_scaff.10485067 | Kronos3935 | 3349 | C | T | Hom | Hom | gcttgagcaactactccaatG | gcttgagcaactactccaatA | gcaatttctttaTccgcagT |
| IWGSC.CSS_4AS_scaff.5984153 | Kronos3935 | 6006 | G | A | Het | Het | agCaggctctggccaagttG | agCaggctctggccaagttA | cgaaTGatgaGtaggcgcT |
| IWGSC.CSS_4BL_scaff.7019402 | Kronos3935 | 9081 | C | T | Het | Het | tgcaatcatgtagtgcgtgG | tgcaatcatgtagtgcgtgA | agcatgatccctagaaCCataC |
| IWGSC.CSS_5BL_scaff.10842786 | Kronos3935 | 3304 | G | A | Het | Het | tggttcccGaaagcctgaaC | tggttcccGaaagcctgaaT | cgcatactgaaacaTGagcAC |
| IWGSC.CSS_6BS_scaff.3045205 | Kronos3935 | 2293 | G | A | Het | Het | aaggaccaagcccaactctcG | aaggaccaagcccaactctcA | agtgatcaagcccaatgtgcA |

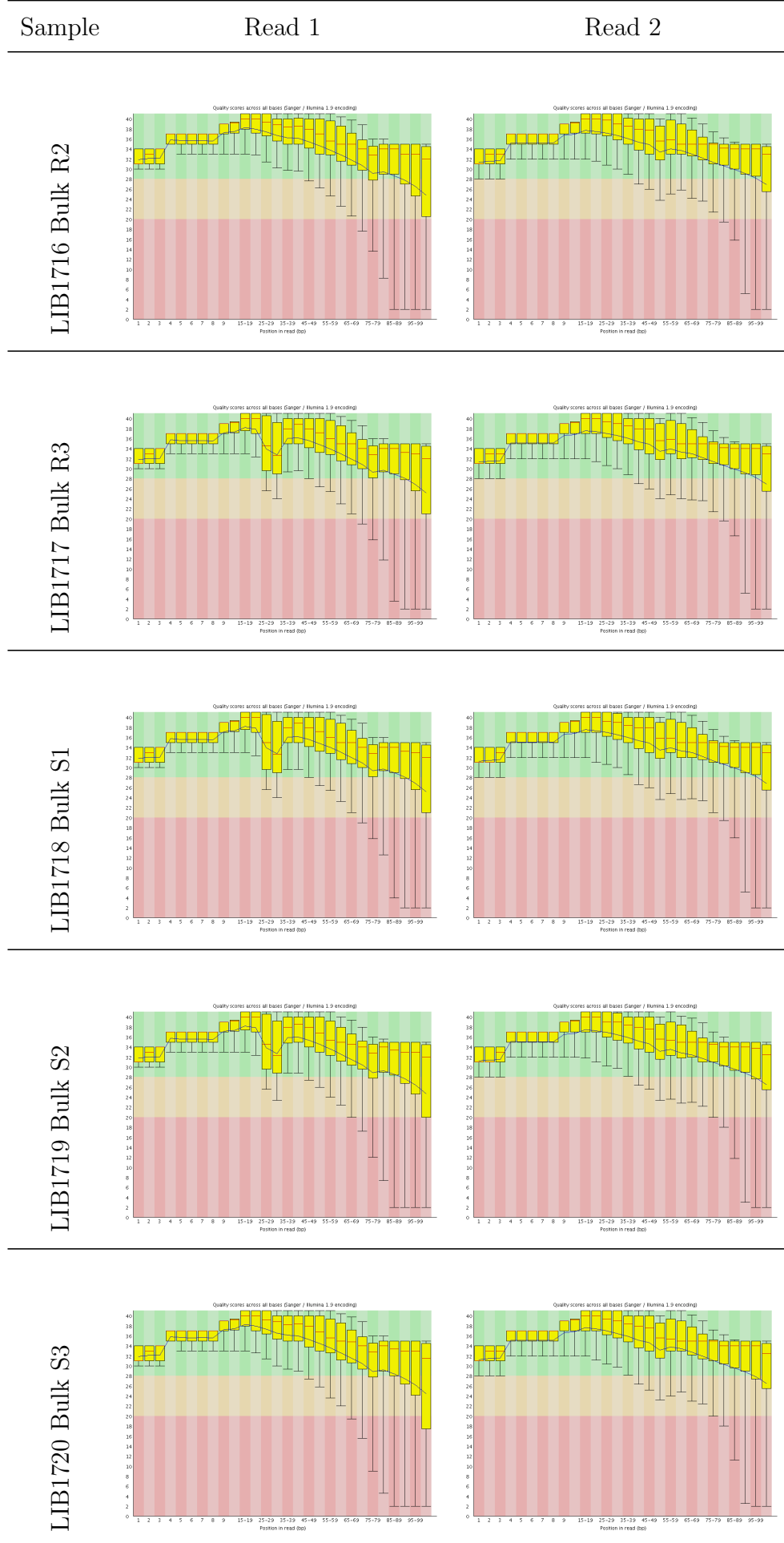
| IWGSC contig | Line | Pos | WT | Mut | Predicted | M_4 | Primer 1 (Kronos) | Primer 2 (mutant) | Common Primer |
|------------------------------|------------|-------|----|-----|-----------|-------|-------------------------|-------------------------|--------------------------|
| IWGSC.CSS_7AL_scaff.4555249 | Kronos3935 | 4487 | C | T | Het | Het | cAgtgctcgagatggcgC | cAgtgctcgagatggcgT | cCttgcaacctctctgatT |
| IWGSC.CSS_1BL_scaff.3918498 | Kronos4240 | 6096 | G | A | Het | Het | ttgcatgccccagaagaG | ttgcatgccccagaagaA | tgggcgaactggtaatgtgG |
| IWGSC.CSS_2BS_scaff.5131713 | Kronos4240 | 5900 | G | A | Het | Het | cctttatcgaggaaagagacacC | cctttatcgaggaaagagacacT | caccattgtagggttcctTttC |
| IWGSC.CSS_5AL_scaff.2769540 | Kronos4240 | 9626 | C | T | Het | Het | tgCagtgtgggaacggaG | tgCagtgtgggaacggaA | catgagtGagatcttctgcT |
| IWGSC.CSS_5BL_scaff.10871091 | Kronos4240 | 7062 | G | A | Het | Het | gccaaaggAaccataacctgC | gccaaaggAaccataacctgT | GgactcttggcAaccggA |
| IWGSC.CSS_6AL_scaff.5800333 | Kronos4240 | 2360 | G | A | Het | Het | cgacaggattgtgagCgC | cgacaggattgtgagCgT | tcagatgctgcaagattcatcT |
| IWGSC.CSS_7BL_scaff.6716931 | Kronos4240 | 2613 | G | A | Het | Het | gGtgGgtattTgcttgggtgaG | gGtgGgtattTgcttgggtgaA | tgGtggactcgacaGtGtA |
| IWGSC.CSS_2BL_scaff.8029221 | Kronos4346 | 2860 | G | A | Het | Het | tgcttccgctcttgctcC | tgcttccgctcttgctcT | atTtgcATCgAtcgggcC |
| IWGSC.CSS_3B_scaff.10460714 | Kronos4346 | 14359 | C | T | Hom | Hom | ctaccttgccatgcgacatG | ctaccttgccatgcgacatA | agcaccccgactctttgacG |
| IWGSC.CSS_4AS_scaff.5989735 | Kronos4346 | 6404 | G | A | Hom | Hom | acgcatgctaacatcagcG | acgcatgctaacatcagcT | actcaagataccaCcgcacG |
| IWGSC.CSS_5BL_scaff.7648030 | Kronos4346 | 6893 | C | T | Het | Het | taccttttctactggcagG | taccttttctactggcagA | ttttcagaggaaacacaggtatcA |
| IWGSC.CSS_6AL_scaff.5755840 | Kronos4346 | 778 | C | T | Het | Het | atcgagtaagctgtcacCgC | atcgagtaagctgtcacCgT | acctgcatgtcaCatccaC |
| IWGSC.CSS_6BS_scaff.2972151 | Kronos4346 | 7876 | G | A | Hom | Hom | gcagcaatgtcActgtttgG | gcagcaatgtcActgtttgA | gcttggactgggcattttatG |
| IWGSC.CSS_7AL_scaff.4542983 | Kronos4346 | 18700 | G | A | Het | Het | gcagggctAccggatacC | gcagggctAccggatacT | catctgccGgttaaacatgC |
| IWGSC.CSS_7BS_scaff.3098098 | Kronos4346 | 5183 | C | T | Het | Het | gCgatatggtacttgcaatgaG | gCgatatggtacttgcaatgaA | ttacattgcttataG'TtgCcgG |
| IWGSC.CSS_1AS_scaff.3259804 | Kronos4485 | 219 | C | T | Het | Het | gtcggcacacccttgC | gtcggcacacccttgT | gcttctttaaggaggcgA |
| IWGSC.CSS_2AL_scaff.6315418 | Kronos4485 | 10490 | G | A | Hom | Hom | gccccctctcaaCcttctcagC | gccccctctcaaCcttctcagT | ttcagacgtCGaggaaatttcC |
| IWGSC.CSS_2BS_scaff.5181092 | Kronos4485 | 3742 | G | A | Het | Het | TggccagcacacctgcaG | TggccagcacacctgcaA | tggacgatgagTgatggAaaT |
| IWGSC.CSS_3B_scaff.10425015 | Kronos4485 | 2372 | C | T | Het | Het | gctactgaagttggCtcGG | gctactgaagttggCtcGA | cttcacatccttgggggTtC |
| IWGSC.CSS_3B_scaff.10775915 | Kronos4485 | 4701 | C | T | Het | Het | ccaagggctgcagagagG | ccaagggctgcagagagA | agacctcacgatGtcctcC |
| IWGSC.CSS_5AL_scaff.2754304 | Kronos4485 | 2301 | G | A | Het | Het | taaccTgccatcgcccG | taaccTgccatcgcccA | cattgGccagccaTgacT |
| IWGSC.CSS_5BL_scaff.10919959 | Kronos4485 | 1867 | C | T | Hom | Hom | gatgccctttgtggagaagG | gatgccctttgtggagaagA | tcttgttcccgaacatgtcA |
| IWGSC.CSS_7AS_scaff.4245431 | Kronos4485 | 3402 | G | A | Hom | Hom | aaggcgctgtgtttcC | aaggcgctgtgtttcT | agtaagtggAacagctaagatcaT |
| IWGSC.CSS_7BL_scaff.6667357 | Kronos4485 | 641 | C | T | Het | Het | gatcAgctgctcattcgagG | gatcAgctgctcattcgagA | ttccctgtcaattgatgccC |

Appendix B

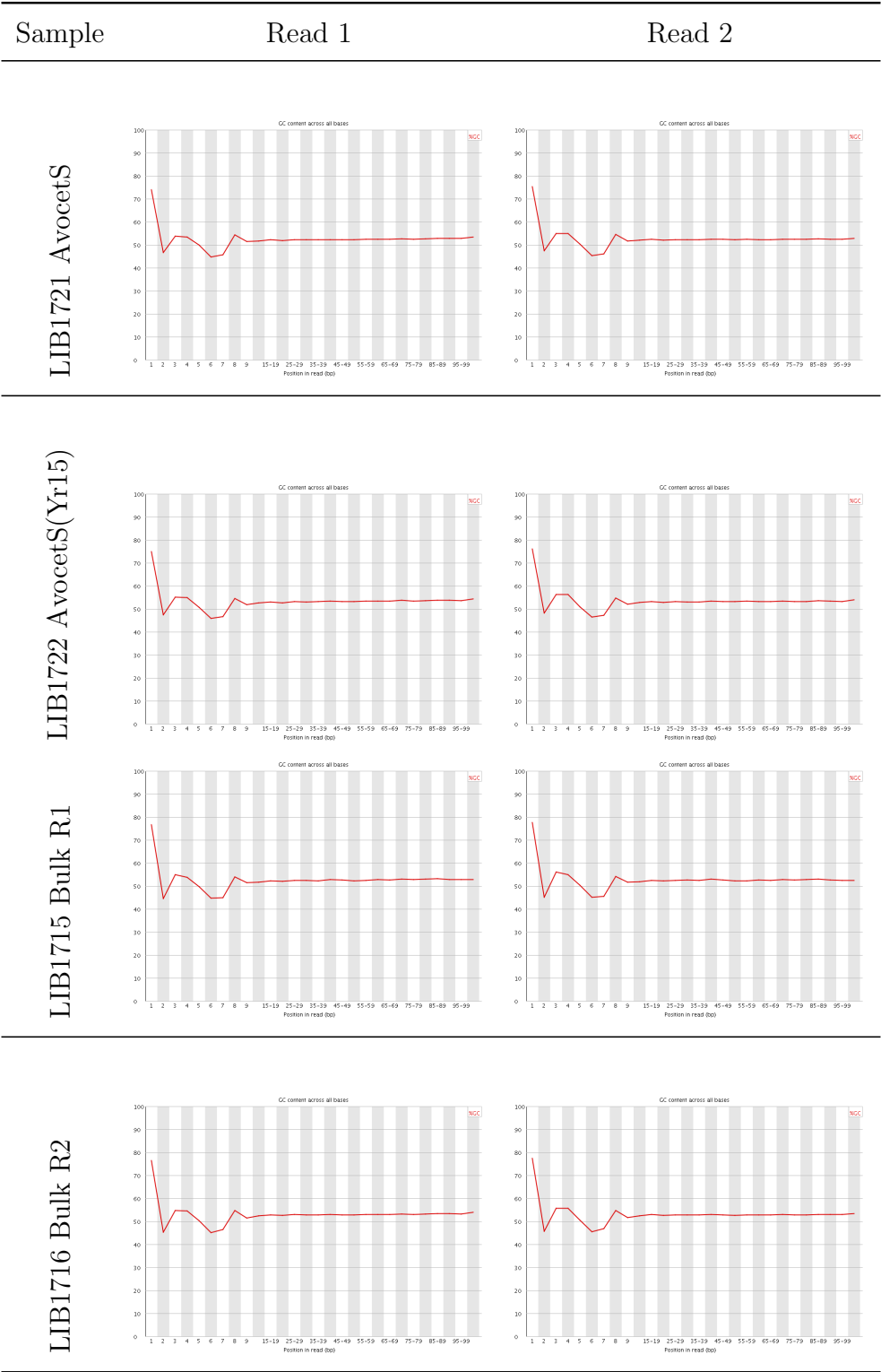
Quality Control

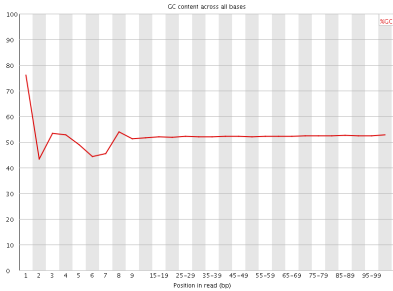
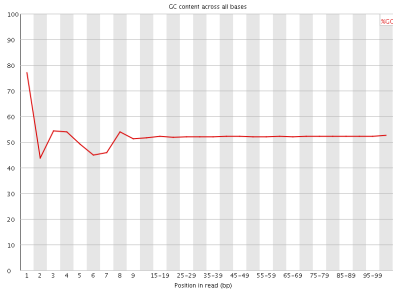
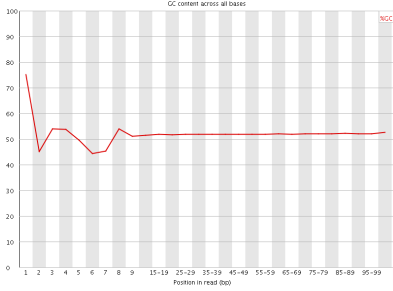
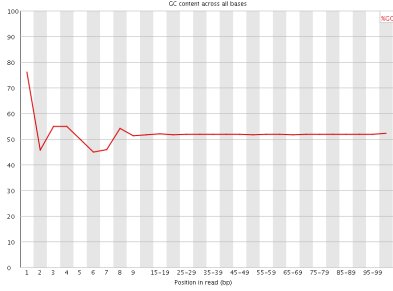
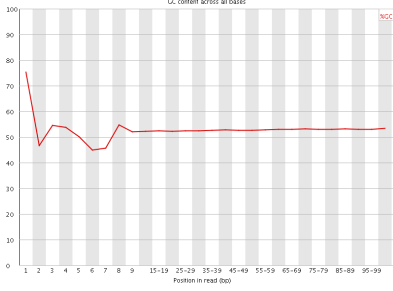
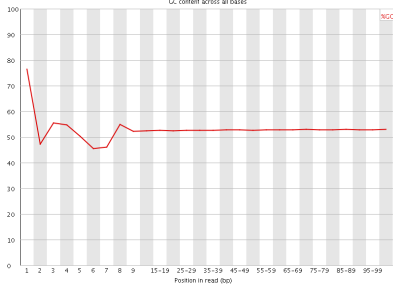
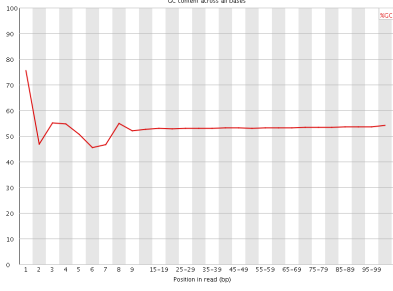
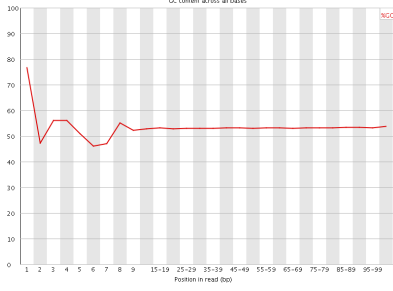
B.1 Sequence read quality





B.2 Sequence GC content



| Sample | Read 1 | Read 2 |
|-----------------|---|--|
| LIB1717 Bulk R3 |  |  |
| | | |
| LIB1718 Bulk S1 |  |  |
| | | |
| LIB1719 Bulk S2 |  |  |
| | | |
| LIB1720 Bulk S3 |  |  |
| | | |

Bibliography

- Abe, A., Kosugi, S., Yoshida, K., et al. Genome sequencing reveals agronomically important loci in rice using MutMap. *Nature biotechnology*, 30(2):174–8, February 2012. ISSN 1546-1696. doi: 10.1038/nbt.2095.
- Akhunov, E. D., Goodyear, A. W., Geng, S., et al. The organization and rate of evolution of wheat genomes are correlated with recombination rates along chromosome arms. *Genome research*, 13(5):753–63, May 2003. ISSN 1088-9051. doi: 10.1101/gr.808603.
- Allen, A. M., Barker, G. L. a., Berry, S. T., et al. Transcript-specific, single-nucleotide polymorphism discovery and linkage analysis in hexaploid bread wheat (*Triticum aestivum* L.). *Plant biotechnology journal*, 9(9):1086–99, December 2011. ISSN 1467-7652. doi: 10.1111/j.1467-7652.2011.00628.x.
- Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–10, October 1990. ISSN 0022-2836. doi: 10.1016/S0022-2836(05)80360-2.
- Babraham Bioinformatics. FastQC A Quality Control tool for High Throughput Sequence Data, 2012. URL <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- Bonnal, R. J. P., Aerts, J., Githinji, G., et al. Biogem: An effective tool-based approach for scaling up open source software development in bioinformatics. *Bioinformatics*, 28(7):1035–1037, April 2012. ISSN 13674803. doi: 10.1093/bioinformatics/bts080.
- Burt, C., Steed, A., Gosman, N., et al. Mapping a type 1 fhb resistance on chromosome 4as of triticum macha and deployment in combination

- with two type 2 resistances. *Theoretical and Applied Genetics*, 128(9): 1725–1738, 2015. ISSN 1432-2242. doi: 10.1007/s00122-015-2542-9.
- Cantu, D., Govindarajulu, M., Kozik, A., et al. Next generation sequencing provides rapid access to the genome of *Puccinia striiformis* f. sp. *tritici*, the causal agent of wheat stripe rust. *PLoS ONE*, 6(8):1–8, 08 2011. doi: 10.1371/journal.pone.0024230.
- Chapman, J. a., Mascher, M., Buluç, A., et al. A whole-genome shotgun approach for assembling and anchoring the hexaploid bread wheat genome. *Genome Biology*, 16(1):1–17, 2015. ISSN 1465-6906. doi: 10.1186/s13059-015-0582-8.
- Cornish-Bowden, A. Nomenclature for incompletely specified bases in nucleic acid sequences: recommendations 1984. *Nucleic acids research*, 13(9):3021–30, May 1985. ISSN 0305-1048.
- Eddy, S. R. What is a hidden Markov model? *Nature biotechnology*, 22 (10):1315–6, October 2004. ISSN 1087-0156. doi: 10.1038/nbt1004-1315.
- Etherington, G. J., Ramirez-Gonzalez, R. H., and MacLean, D. bio-samtools 2: a package for analysis and visualization of sequence and alignment data with SAMtools in Ruby. *Bioinformatics*, pages 1–2, 2015. ISSN 1367-4803. doi: 10.1093/bioinformatics/btv178.
- Gerechter-Amitai, Z. K., van Silfhout, C. H., Grama, A., and Kleitman, F. Yr15 — a new gene for resistance to puccinia striiformis in triticum dicoccoides sel. g-25. *Euphytica*, 43(1):187–190, 1989. ISSN 1573-5060. doi: 10.1007/BF00037912.
- GM, C. Chloroplasts and Other Plastids. In *The Cell: A Molecular Approach*. Sinauer Associates, 2000. URL http://www.ncbi.nlm.nih.gov/books/NBK9905/?redirect-on-error=__HOME__.
- Goto, N., Prins, P., Nakao, M., et al. BioRuby: bioinformatics software for the Ruby programming language. *Bioinformatics (Oxford, England)*, 26(20):2617–9, October 2010. ISSN 1367-4811. doi: 10.1093/bioinformatics/btq475.

Hodges, E., Xuan, Z., Balija, V., et al. Genome-wide in situ exon capture for selective resequencing. *Nat Genet*, 39(12):1522–1527, Dec 2007. ISSN 1061-4036. doi: 10.1038/ng.2007.42.

Hubbard, A., Lewis, C. M., Yoshida, K., et al. Field pathogenomics reveals the emergence of a diverse wheat yellow rust population. *Genome Biology*, 16(1):1–15, 2015. ISSN 1465-6906. doi: 10.1186/s13059-015-0590-8.

Hutchison, C. a. DNA sequencing: bench to bedside and beyond. *Nucleic acids research*, 35(18):6227–37, January 2007. ISSN 1362-4962. doi: 10.1093/nar/gkm688.

Illumina. *CASAVA v1.8.2 User guide*. Revc edition, 2011. URL http://support.illumina.com/sequencing/sequencing_software/casava/documentation

James, G. V., Patel, V., Nordström, K. J., et al. User guide for mapping-by-sequencing in Arabidopsis. *Genome biology*, 14(6):R61, June 2013. ISSN 1465-6914. doi: 10.1186/gb-2013-14-6-r61.

Katoh, K. and Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Molecular biology and evolution*, 30(4):772–80, April 2013. ISSN 1537-1719. doi: 10.1093/molbev/mst010.

Kent, W. J. BLAT—The BLAST-Like Alignment Tool. *Genome Research*, 12(4):656–664, March 2002. ISSN 1088-9051. doi: 10.1101/gr.229202.

King, R., Bird, N., Ramirez-Gonzalez, R., et al. Mutation scanning in wheat by exon capture and next-generation sequencing. *PLoS ONE*, 10(9):1–18, 09 2015. doi: 10.1371/journal.pone.0137549.

Krasileva, K., Vasquez-Gross, H., Howell1, T., et al. Uncovering hidden variation in young polyploid wheat genomes. submitted 2016.

Krasileva, K. V., Buffalo, V., Bailey, P., et al. Separating homeologs by phasing in the tetraploid wheat transcriptome. *Genome biology*, 14(6):R66, June 2013. ISSN 1465-6914. doi: 10.1186/gb-2013-14-6-r66.

Lander, E. S., Linton, L. M., Birren, B., et al. Initial sequencing and analysis of the human genome. *Nature*, 409(6822):860–921, February 2001. ISSN 0028-0836. doi: 10.1038/35057062.

- LGC Genomics. <http://www.lgcgroup.com/services/genotyping/>, 2013. URL <http://www.lgcgroup.com/services/genotyping/>.
- Li, H. and Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics (Oxford, England)*, 25(14):1754–60, July 2009. ISSN 1367-4811. doi: 10.1093/bioinformatics/btp324.
- Li, H., Handsaker, B., Wysoker, A., et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford, England)*, 25(16):2078–9, August 2009. ISSN 1367-4811. doi: 10.1093/bioinformatics/btp352.
- Liu, Y. and Schmidt, B. Long read alignment based on maximal exact match seeds. *Bioinformatics (Oxford, England)*, 28(18):i318–i324, September 2012. ISSN 1367-4811. doi: 10.1093/bioinformatics/bts414.
- Ma, J., Stiller, J., Zheng, Z., et al. A high-throughput pipeline for detecting locus-specific polymorphism in hexaploid wheat (*triticum aestivum* l.). *Plant Methods*, 11(1), aug 2015. doi: 10.1186/s13007-015-0082-6.
- MAS Wheat. Mas wheat transcriptome. supplemental file 17., 2013. URL <http://maswheat.ucdavis.edu/Transcriptome/index.htm>.
- Mayer, K. F. X., Rogers, J., Dole el, J., et al. A chromosome-based draft sequence of the hexaploid bread wheat (*Triticum aestivum*) genome. *Science*, 345(6194):1251788–1251788, July 2014. ISSN 0036-8075. doi: 10.1126/science.1251788.
- Metzker, M. L. Sequencing technologies - the next generation. *Nature reviews. Genetics*, 11(1):31–46, January 2010. ISSN 1471-0064. doi: 10.1038/nrg2626.
- Michelmore, R. W., Paran, I., and Kesseli, R. V. Identification of markers linked to disease-resistance genes by bulked segregant analysis: a rapid method to detect markers in specific genomic regions by using segregating populations. *Proceedings of the National Academy of Sciences*, 88(21):9828–9832, November 1991. ISSN 0027-8424. doi: 10.1073/pnas.88.21.9828.
- Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L., and Wold, B. Mapping and quantifying mammalian transcriptomes by RNA-

- Seq. *Nature methods*, 5(7):621–8, July 2008. ISSN 1548-7105. doi: 10.1038/nmeth.1226.
- Murphy, L. R., Santra, D., Kidwell, K., et al. Linkage Maps of Wheat Stripe Rust Resistance Genes and for Use in Marker-Assisted Selection. *Crop Science*, 49(5):1786, 2009. ISSN 1435-0653. doi: 10.2135/crop-sci2008.10.0621.
- Myllykangas, S., Buenrostro, J., and Ji, H. P. *Bioinformatics for High Throughput Sequencing*. Springer New York, New York, NY, 2012. ISBN 978-1-4614-0781-2. doi: 10.1007/978-1-4614-0782-9. URL <http://www.springerlink.com/index/10.1007/978-1-4614-0782-9>.
- Needleman, S. B. and Wunsch, C. D. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology*, 48(3):443–453, March 1970. ISSN 00222836. doi: 10.1016/0022-2836(70)90057-4.
- Peng, J., Fahima, T., Röder, M., and Huang, Q. High-density molecular map of chromosome region harboring stripe-rust resistance genes YrH52 and Yr15 derived from wild emmer wheat, *Triticum dicoccoides*. *Genetica*, pages 199–210, 2000.
- Pickrell, J. K., Marioni, J. C., Pai, A. A., et al. Understanding mechanisms underlying human gene expression variation with rna sequencing. *Nature*, 464(7289):768–772, Apr 2010. ISSN 0028-0836. doi: 10.1038/nature08872.
- Pontius, J., Wagner, L., and Schuler, G. UniGene: A Unified View of the Transcriptome. In *The NCBI Handbook [Internet]*, chapter 21. National Center for Biotechnology Information (US), October 2002. URL <http://www.ncbi.nlm.nih.gov/books/NBK21083/>.
- Ramirez-Gonzalez, R. H., Uauy, C., and Caccamo, M. PolyMarker: A fast polyploid primer design pipeline. *Bioinformatics*, pages 2–3, 2015a. ISSN 1367-4803. doi: 10.1093/bioinformatics/btv069.
- Ramirez-Gonzalez, R. H., Bonnal, R., Caccamo, M., and Maclean, D. Bio-samtools: Ruby bindings for SAMtools, a library for accessing BAM files containing high-throughput sequence alignments. *Source*

- code for biology and medicine*, 7(1):6, January 2012. ISSN 1751-0473. doi: 10.1186/1751-0473-7-6.
- Ramirez-Gonzalez, R. H., Segovia, V., Bird, N., Caccamo, M., and Uauy, C. *Next Generation Sequencing Enabled Genetics in Hexaploid Wheat*, pages 201–209. Springer Japan, Tokyo, 2015b. ISBN 978-4-431-55675-6. doi: 10.1007/978-4-431-55675-6_22. URL http://dx.doi.org/10.1007/978-4-431-55675-6_22.
- Ramirez-Gonzalez, R. H., Segovia, V., Bird, N., et al. Rna-seq bulked segregant analysis enables the identification of high-resolution genetic markers for breeding in hexaploid wheat. *Plant Biotechnology Journal*, 13(5):613–624, 2015c. ISSN 1467-7652. doi: 10.1111/pbi.12281.
- Randhawa, H. S., Mutti, J. S., Kidwell, K., et al. Rapid and targeted introgression of genes into popular wheat cultivars using marker-assisted background selection. *PloS one*, 4(6):e5752, January 2009. ISSN 1932-6203. doi: 10.1371/journal.pone.0005752.
- Rozen, S. and Skaletsky, H. Primer3 on the WWW for general users and for biologist programmers. *Methods in molecular biology (Clifton, N.J.)*, 132:365–86, January 2000. ISSN 1064-3745.
- Schneeberger, K. and Weigel, D. Fast-forward genetics enabled by new sequencing technologies. *Trends in Plant Science*, 16(5):282–288, 2011.
- Schneeberger, K., Ossowski, S., Lanz, C., et al. Shoremap: simultaneous mapping and mutation identification by deep sequencing. *Nat Meth*, 6(8):550–551, Aug 2009. ISSN 1548-7091. doi: 10.1038/nmeth0809-550.
- Shendure, J. and Ji, H. Next-generation DNA sequencing. *Nature biotechnology*, 26(10):1135–45, October 2008. ISSN 1546-1696. doi: 10.1038/nbt1486.
- Slater, G. S. C. and Birney, E. Automated generation of heuristics for biological sequence comparison. *BMC bioinformatics*, 6:31, January 2005. ISSN 1471-2105. doi: 10.1186/1471-2105-6-31.
- Smith, T. and Waterman, M. Identification of common molecular subsequences. *Journal of Molecular Biology*, 147(1):195–197, March 1981. ISSN 00222836. doi: 10.1016/0022-2836(81)90087-5.

- Sun, G. L., Fahima, T., Korol, A. B., et al. Identification of molecular markers linked to the Yr15 stripe rust resistance gene of wheat originated in wild emmer wheat, *Triticum dicoccoides*. *TAG Theoretical and Applied Genetics*, 95(4):622–628, September 1997. ISSN 0040-5752. doi: 10.1007/s001220050604.
- Takagi, H., Abe, A., Yoshida, K., et al. QTL-seq: rapid mapping of quantitative trait loci in rice by whole genome resequencing of DNA from two bulked populations. *The Plant journal : for cell and molecular biology*, 74(1):174–83, May 2013a. ISSN 1365-313X. doi: 10.1111/tpj.12105.
- Takagi, H., Uemura, A., Yaegashi, H., et al. MutMap-Gap: whole-genome resequencing of mutant F2 progeny bulk combined with de novo assembly of gap regions identifies the rice blast resistance gene Pii. *The New phytologist*, 200(1):276–83, October 2013b. ISSN 1469-8137. doi: 10.1111/nph.12369.
- Tang, H., Zhang, X., Miao, C., et al. Allmaps: robust scaffold ordering based on multiple maps. *Genome Biology*, 16(1):1–15, 2015. ISSN 1465-6906. doi: 10.1186/s13059-014-0573-1.
- Trick, M., Adamski, N., Mugford, S. G., et al. Combining SNP discovery from next-generation sequencing data with bulked segregant analysis (BSA) to fine-map genes in polyploid wheat. *BMC plant biology*, 12(1):14, January 2012. ISSN 1471-2229. doi: 10.1186/1471-2229-12-14.
- Van Ooijen, J. and Jansen, J. *Estimation of recombination frequencies; Genetic Mapping in Experimental Populations*, pages 73–133. Cambridge University Press, Cambridge, 2013. ISBN 978-1-107-0132-16.
- Wang, S., Wong, D., Forrest, K., et al. Characterization of polyploid wheat genomic diversity using a high-density 90 000 single nucleotide polymorphism array. *Plant biotechnology journal*, 12(6):787–796, March 2014. ISSN 1467-7652. doi: 10.1111/pbi.12183.
- Wang, Y., Tiwari, V. K., Rawat, N., et al. GSP: a web-based platform for designing genome-specific primers in polyploids. *Bioinformatics*, page btw134, mar 2016. doi: 10.1093/bioinformatics/btw134.

- Wellings, C. and McIntosh, R. A. Host-pathogen studies of wheat stripe rust in australia. In Slinkard, A., editor, *Proceedings 9th International Wheat Genetics Symposium*, pages 336–338. University of Saskatchewan, Saskatoon, SK, Canada, 1998.
- Wilkinson, P. a., Winfield, M. O., Barker, G. L. a., et al. CerealsDB 2.0: an integrated resource for plant breeders and scientists. *BMC bioinformatics*, 13(1):219, January 2012. ISSN 1471-2105. doi: 10.1186/1471-2105-13-219.