



Faculty 09: Agricultural Sciences, Nutritional Sciences, and Environmental Management

Pangenome Analysis of Haplotypes Associated with Root Dry Mass in Wheat

Thesis of the Master of Science (MSc.) in Agrobiotechnology

Submitted by:

José Antonio Montero Tena

Matriculation No. 7062197

Department of Plant Breeding

Supervisors:

Dr. Christian Obermeier

Prof. Dr. Rod Snowdon

Giessen, September 2021

TABLE OF CONTENTS

| | |
|--|-----------|
| 1. General Introduction..... | 7 |
| 1.1. Wheat. Relevance, evolution and genomics..... | 7 |
| 1.2. Increasing global wheat yield is required to meet future demand..... | 8 |
| 1.3. Plant breeding..... | 9 |
| 1.4. Role of single-nucleotide polymorphisms in wheat breeding and limitations..... | 10 |
| 1.5. Potential impact of sequence-based haploblocks in wheat breeding..... | 11 |
| 1.6. Aim of the master thesis..... | 12 |
| 2. Pangenome Analysis of Haplotypes Associated with Root Dry Mass in Wheat..... | 13 |
| 2.1. Introduction..... | 13 |
| 2.2. Materials and Methods..... | 16 |
| 2.2.1. BLAST calling of SNP-based RDM haplotypes among wheat assemblies..... | 16 |
| 2.2.2. Identification of sequence-based haploblock predictions with crop-haplotypes.com..... | 17 |
| 2.2.3. Validation and redefinition of haploblocks..... | 19 |
| 2.2.4. Statistical analysis of the differences in alignment coverage..... | 21 |
| 2.2.5. Search for KASP marker | 22 |
| 2.2.6. Identification of SNPs through high-stringency variant calling from Illumina paired-end short reads | 23 |

| | |
|--|-----------|
| 2.3. Results and Discussion..... | 26 |
| 2.3.1. Two wheat assemblies contain a SNP-haplotype that increases root dry mass..... | 26 |
| 2.3.2. A sequence-based haplotype matches with the SNP-based RD ^a -h2 according to crop-haplotypes.com..... | 27 |
| 2.3.3. Using scaffolds as query in NUCmer pairwise alignments reduces alignment coverage..... | 28 |
| 2.3.4. Checking alignment properties is essential for accurate haploblock prediction..... | 30 |
| 2.3.5. Low alignment coverage and discontinuities hamper validation and redefinition of the haploblock RD ^a | 33 |
| 2.3.6. SNP markers from an array could target the potentially larger haploblock..... | 35 |
| 2.3.7. Genome sequences of RD ^b -h3-carrying varieties were identified and provided new SNP markers..... | 37 |
| 3. General Discussion..... | 38 |
| 4. Summary..... | 39 |
| 5. References..... | 39 |
| 6. Appendix..... | 46 |
| 7. Declaration of Authorship..... | 79 |

LIST OF FIGURES

| | |
|--|-----------|
| Figure 1: Hybridization events involved in the evolution of bread wheat, Triticum aestivum..... | 7 |
| Figure 2: Evolution of wheat yield (hg/ha) between 1961 and 2019..... | 8 |
| Figure 3: Haploblock predictions on crop-haplotypes.com..... | 28 |
| Figure 4: Boxplots showing differences in the percentage of alignment coverage of the reference chromosome..... | 30 |
| Figure 5: Effects of changing bin size on haploblock predictions on crop-haplotypes.com in the region between 655 and 665 Mbp of Lancer's chromosome 5B..... | 32 |
| Figure 6: Potential wrong assignment of the Norin's chromosome 5B region between 537 and 538 Mbp to a 1-Mbp-long haploblock at 1-Mbp bin size on crop-haplotypes.com..... | 33 |
| Figure 7: Comparison between the four approaches for haploblock redefinition..... | 35 |

LIST OF TABLES

| | |
|---|-----------|
| Table 1: Information about the SNP markers in RDM haploblocks discovered by Voss-Fels et al. (2017)..... | 17 |
| Table 2: SNP-haplotype pattern for the RDM haploblocks among the 15 wheat genome assemblies..... | 26 |
| Table 3: Candidate codominant SNP markers to target the redefined potential RDMA haploblock..... | 36 |

LIST OF SUPPLEMENTARY FILES

| | |
|---|-----------|
| Supplementary file 1: Output of the R script ‘Haplotype-based pangenome analysis in wheat’ showing the analysis of the haploblock RD^a in Lancer and Paragon as example..... | 47 |
| Supplementary file 2: List of potential SNP markers to target the haplotype RD^b-h3, associated with higher root dry mass..... | 70 |

1. GENERAL INTRODUCTION

1.1. Wheat. Relevance, evolution and genomics

Wheat is a staple crop that feeds over 40% of the world's population, providing 20% of the daily protein and food calories and whose production requires around 218 million hectares (Giraldo et al., 2019). Wheat originated in the fertile crescent around 7000 years ago and was one of the first cereals to be domesticated and has been the main food of the major civilizations of Europe, West Asia and North Africa. Modern bread wheat (*Triticum aestivum* L.) is an allohexaploid species ($2n = 6x = 42$ chromosomes, genomic composition AABBDD) that evolved through two polyploidization events between wild relative species (Fig. 1). Hybridization has been an important source of diversity in wheat (Warschefsky et al., 2014).

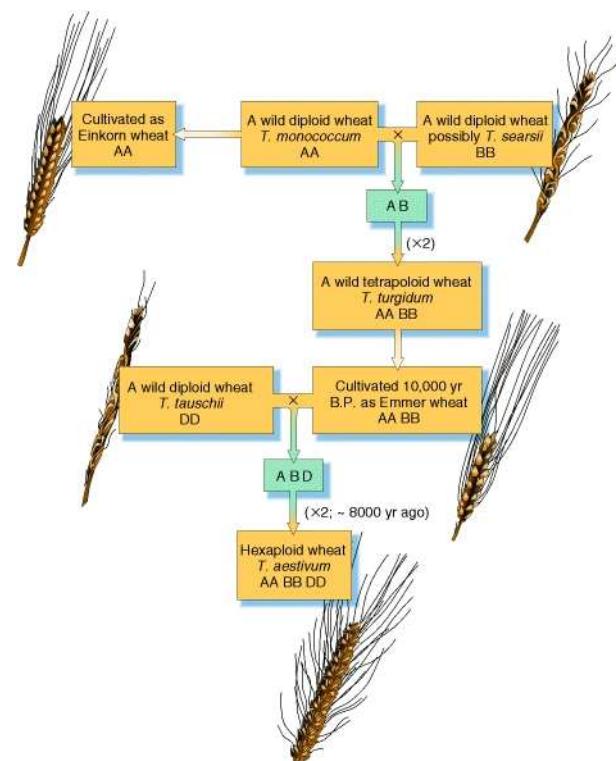


Figure 1: Hybridization events involved in the evolution of bread wheat, *Triticum aestivum*. The 'X2' refers to the doubling of the chromosome complement which gives rise in fertile hybrids. Source: https://www.cerealsdb.uk.net/cerealgenomics/WheatBP/Documents/DOC_Evolution.php

1.2. Increasing global wheat yield is required to meet future demand

Germany is one of the top European wheat producers, mainly bread wheat. In 2019, the last year of record, yield in Germany was nearly 7.4 tons of wheat per hectare, which doubled the average yield across the world (Fig. 2).

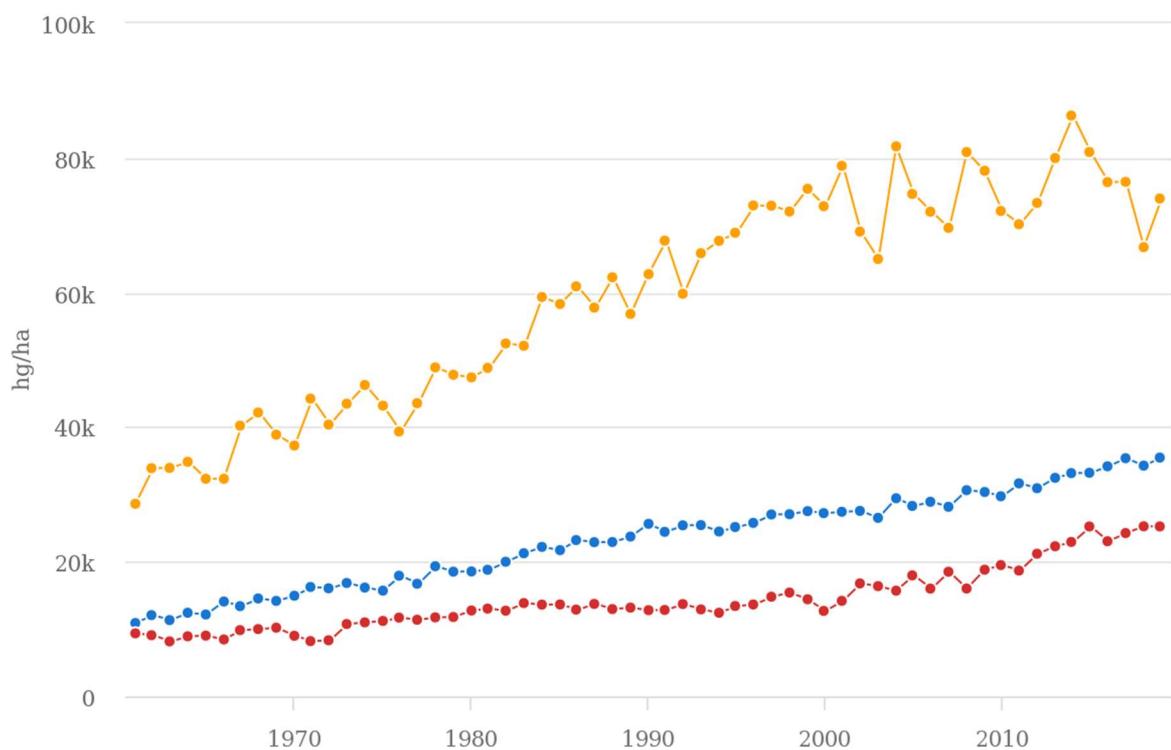


Figure 2: Evolution of wheat yield (hg/ha) between 1961 and 2019. Data is shown for the world average (blue), least developed countries (red) and Germany (orange). Source: FAOSTATS (2021).

The rapidly growing demand for food, feed and fuel due to the world population expansion, which is predicted to surpass nine billion by 2050, will require further increase in crop yields (ca. 2% per annum). However, wheat yield seems to have begun stagnating since the mid 1990s, especially in the biggest wheat-producing countries as Germany, France, United Kingdom or United States (Brisson et al., 2010). On the other hand, the world wheat yields follow a positive trend, although this seems to be due to the growing wheat yields in least developed countries and yet their annual growth is not bigger than 1% (FAOSTATS). According to Reynolds et al. (2012), current genetic gains are not enough to meet the food demands by 2050. Yield increases need to be increased in the future mainly through genetic yield improvement for temperate

regions in Europe with high wheat yields as Germany (Ewert et al., 2005). The "Green Revolution" in the 1960s set an example of the importance of genetic improvement in wheat production, with the introgression of dwarfing genes into modern high-yielding cultivars. Semi-dwarf wheat is easier to harvest and less prone to lodge in fertile and humid conditions and the fast-growing high yields achieved in the following decades were acknowledged to this trait (Casebow et al., 2016).

1.3. Plant breeding

Plant breeding has been used for thousands of years to stack beneficial plant traits and has a great value to human nutrition, feeding of farm animals or fibre production. The realization that many traits are controlled by genes led to the evolution from traditional plant breeding, aimed towards the selection and cross of plants with superior traits, to modern plant breeding, which incorporates molecular biology techniques to select specific genes involved in desired traits (Hartung and Schiemann, 2014). This is the purpose of marker-assisted selection (MAS), which is based on the use of DNA markers in plant breeding.

DNA markers are DNA sequences with known location in a genome. DNA markers can be linked by proximity with locus, e.g. specific position of a gene or DNA sequence on a chromosome, which can at the same time be associated with specific traits. In this case, we say that there is a marker-trait association (MTA). Collard and Mackill (2008) established the requirements that DNA markers should meet for their implementation in MAS: reliability; quantity and quality, technical procedure for marker assay, level of polymorphism and cost. The reliability of DNA markers depends on how close their location is to the target loci, since DNA markers are often not causal but simply associated with the trait of interest (Brinton et al., 2020). Also, higher polymorphism is a positive characteristic of DNA markers, since multiallelic markers can discriminate between more genotypes.

1.4. Role of single-nucleotide polymorphisms in wheat breeding and limitations

A single nucleotide polymorphism (SNP) is a single position in DNA that varies between individuals. SNPs are usually biallelic, e.g. one of two different nucleotides can be found in the SNP, each one representing a different allele. SNPs can be DNA markers when they show linkage disequilibrium with traits, e.g. when the allele distribution of a SNP is not random. Currently, SNPs markers represent the most common DNA marker because they occur abundantly in plant genomes and are easy to detect (Qian et al., 2017). SNPs can be discovered by pan-genome diversity analysis based on assemblies of wheat thanks to the availability of a chromosome-level reference assembly, Chinese Spring (IWGSC, 2014), plus the recent publication of further 9 reference-quality chromosome-level assemblies and 5 scaffold-level assemblies by the 10+ Wheat Genome Project (Walkowiak et al., 2020). These breakthroughs were enabled by the next-generation sequencing (NGS) technologies, despite the high complexity of wheat genome due to its enormous (~17,000 Mbp) repetitive (80% consistent of repetitive sequences) genome (IWGSC, 2014). NGS also allows SNP discovery without reference genome through genotyping-by-sequencing (GBS) methods. Poland et al. (2012) or Lin et al. (2015) provide examples of SNP discovery by GBS. The genetic diversity that SNPs provide can be exploited by high-density SNP genotyping arrays that allow high throughput analysis of large wheat populations. Wheat counts with 9k, 35 k, 90 k and 800 k SNPs platforms (Voss-Fels and Snowdon, 2015) that can be used in genome-wide association studies (GWAS) to link SNP markers with complex genetic traits controlled by quantitative trait loci (QTL). However, once marker-trait associations are established, SNP technologies can also be implemented with only one or a few markers to select lines carrying the diagnostic markers as in competitive allele specific PCR (KASP), which is a very common cost-effective technique in MAS (He, Holme and Anthony, 2014).

Nevertheless, SNP markers also have limitations. First, each SNP can only discriminate between two genotypes due to their biallelic nature. To overcome this issue SNPs tend to be clustered into haploblocks based on linkage disequilibrium (Qian et al., 2017). Bernardo (2010) defined haplotypes as two or more SNP alleles that tend to be inherited as a unit. Multiallelic SNP combinations are therefore capable of calling more than two genotypes and the resulting haploblocks are effective units of recombination by breeders (Walkowiak et al., 2020). Secondly, the non-causality of the SNPs identified in QTL analysis by GWAS towards the traits they are in LD with (Platten et al., 2019) can lead to the selection of false positive or negative individuals during the breeding process (Walkowiak et al., 2020). Also, SNP arrays lack of resolution to detect underrepresented rare alleles (Voss-Fels et Snowdon, 2016). Finally, SNP arrays are not exempt of genotyping errors. To avoid the difficulties related with these errors, maximum level of diversity (1-3%) are implemented to account for genotyping arrays, grouping together highly similar haploblocks (Brinton et al., 2020).

1.5. Potential impact of sequence-based haploblocks in wheat breeding

Despite SNPs being the traditional method to define haploblocks, these can also be defined by different methods to overcome the limitations of SNPs. The abovementioned new wheat assemblies allowed Brinton et al. (2020) to define wheat haploblocks based on full DNA sequence comparison between assemblies. First, data was generated by pairwise alignments between wheat assemblies, followed by alignment filtering by length to get rid of repetitive non-syntenic transposons. Next, alignments were clustered into bins by position to finally call haploblocks in those bins whose median percentage of identity was ≥ 99.99 , to account for sequencing errors. Bins had a fixed size, which provided different resolutions to haploblock detection. This sequence-based approach to call haploblocks has been implemented more often on monoploid or diploid genomes as humans (Rizzi et al., 2019; Matsumoto and Kiryu, 2013; Barrett et al., 2005) or some human pathogens (Hong et al., 2019; Boose et al., 2011; Zheng et

al., 2004), especially by targeting specific loci by Sanger sequencing. Nevertheless, the rapid advancement of wheat pangenomics explained by the implementation of NGS technologies enabled the development of this new method. Brinton et al. (2020) displayed the sequence-based wheat haploblocks across the 15 assemblies on the website crop-haplotypes.com. This website is available for free and provides the opportunity to understand genetic diversity across the wheat pangenome. One potential application would be the redefinition of SNP-based haploblocks discovered by SNP arrays into sequence-based haploblocks. This would facilitate the exploration of conserved genomic regions between assemblies to develop new more precise SNP markers to discriminate between more genotypes with higher sensibility than is usually possible with the markers obtained from GWAs.

1.6. Aim of the master thesis

This master thesis focuses on the analysis of SNP-based haploblocks involved in higher root dry mass (Voss Fels et al., 2017) with bioinformatic methods, as the sequence-based method developed by Brinton to call haplotypes in wheat or variant calling from Illumina short reads. Root-related trait breeding is often discriminated in favor of shoot-related traits due to the difficulty of measuring root phenotypic indicators, despite of roots being essential for absorbing nutrients and maintaining physical stability in wheat (Voss Fels et al., 2017). Selection for shoot-related traits could have eroded below-ground traits meanwhile having higher root dry mass could support the adaptation of wheat to challenging environments (Voss-Fels, Snowdon and Hickey (2018)). At the time of the discovery of these SNP markers by GWAS, pan-genomic diversity analysis was not possible since only the reference genome Chinese Spring (IWGSC, 2014) was available. The bioinformatic methods applied for this study have in common the use of pangenomic resources in wheat to analyze genomic diversity.

2. PANGENOME ANALYSIS OF HAPLOTYPES ASSOCIATED WITH ROOT DRY MASS

2.1. INTRODUCTION

Haplotypes are combinations of conserved allele sequences that are located in genome regions known as haploblocks or haplotype blocks, which have not been disrupted significantly by recombination events. Therefore, haplotypes have diverged from each other by random mutations and each haploblock contains at least a few haplotypes (Gabriel et al., 2002). Since gene combinations tend not to be inherited independently but in genetic linkage as haplotype blocks, haplotypes are important targets for trait selection and haploblocks can facilitate marker development for applied breeding (Walkowiak et al., 2020).

Haploblocks have been defined in the past by SNP markers because these markers occur frequently in the genome and are easy to detect (Qian et al., 2017). Bernardo (2010) describes a haplotype as “two or more SNP alleles that tend to be inherited as a unit.” However, the method used for defining haplotype blocks changes between studies and the term ‘haploblock’ might seem ambiguous (Wall and Pritchard, 2003). The most common method is to assign polymorphic sites together into haplotype blocks based on linkage disequilibrium between markers showing little or no evidence of historical recombination and a joint inheritance in the same chromosomal block (Gabriel). Hereafter, we will refer to haplotypes and haploblocks defined by this method as SNP-haplotypes and SNP-haploblocks to avoid confusion with those defined by other methods.

One example of the discovery of SNP-haploblocks in allohexaploid wheat (*Triticum aestivum*) was provided by Voss-Fels et al. (2017), who defined two adjacent haploblocks associated with root dry mass (RDM) in chromosome 5B, Hap-5B-RDMa and Hap-5B-RDMb (hereafter shortened to RDMA and RDMb). These blocks were defined upon pairwise linkage

disequilibrium between 9 and 6 SNP markers, respectively from the 90k SNP Illumina Infinium wheat genotyping array in a genome-wide association study (GWAS) with 215 homozygous wheat accessions. After the application of locus-specific KASP markers designed from the original SNP chips by Makhoul et al. (2020), 2 and 4 haplotypes were identified among the tested population for either RDMa (h1 and h2) and RDMb (h1, h2, h3 and h8). A small group of wheat accessions carrying the combination of favourable haplotype variants RDMa-h2 and RDMb-h3 had significantly higher root dry mass than the rest of the accessions (Voss-Fels et al., 2017).

SNP-haploblocks are a popular target in marker-assisted selection (MAS) to improve the efficiency and speed of the breeding process in many crops. MAS is based on the correlation between the marker and the trait, but this correlation is not direct, since most SNP markers are not causal for phenotypes but are rather linked with the quantitative trait loci (QTL), e.g. genetic regions associated with specific phenotypic traits, which can result in selection of false-positive and false-negative individuals and hampers the use of SNP markers in breeding programs (Platten et al., 2019).

Other methods to define haploblocks are those based on sequence comparisons and have been exploited mainly in research about human genomics (Rizzi et al., 2019; Matsumoto and Kiryu, 2013; Barrett et al., 2005) and population genetics studies of pathogens (Hong et al., 2019; Boose et al., 2011; Zheng et al., 2004). However, these methods are difficult to apply in *Triticum aestivum* because of its large genome size, high sequence similarity between homoeologous chromosomes and abundance of repetitive elements (IWGSC, 2014).

The release of new assemblies helped Brinton et al. (2020) to detect haploblocks across wheat genome by sequence comparison between assemblies. They reported the shared haplotypes on the website crop-haplotypes.com. Brinton's method defined haplotype blocks as physical

regions with ≥99.99% sequence identity (across a fixed bin size of 5-, 2.5- or 1-Mbp) between pairwise comparisons of wheat lines in the 10+ Wheat Genome Project, which included 15 assemblies: 10 chromosome-level assemblies (CLA), including the IWGSC Chinese Spring reference genome plus reference-quality pseudomolecules (ArinaLrFor, Jagger, Julius, Lancer, Landmark, Mace, Norin61, Stanley, SY-Mattis) and 5 scaffold-level assemblies (SLA; Cadenza, Claire, Paragon, Robigus and Weebill). These 14 additional non-reference genomes were published by Walkowiak et al. to expand the reference genome assembly. Brinton and colleagues intended to develop a systematic haplotype-led approach to identify genetic diversity and enable its exploitation. The advantage of this method to the traditional SNP-based methods are increased precision in haplotype definition, since SNP markers tend not to be causal for traits and can cause false positives, and increased sensibility, being capable of differentiating ‘identical-by-state’ haplotypes from ‘near-identical’ sequences. Brinton and colleagues provided evidence of the importance of this sensibility for plant breeding that can be relevant in some cases.

The aim of this master thesis was to analyze the SNP-haploblocks RD^a and RD^b using the new sequence-based method of haplotype detection developed by Brinton et al. (2020). It was considered that the new wheat pangenome could help to validate the haplotypes and redefine the haplotype blocks to a more precise location for increased correlation marker-trait. Additionally, this could lead to the discovery of further SNP markers to discriminate between varieties carrying the RD^a haplotypes and those carrying nearly-identical alleles. Results suggest that Brinton et al. reported inaccurate haploblock predictions mainly due to lower alignment coverage in comparisons using scaffolds as query. Considering the necessity of good-quality alignment properties in haploblock predictions, a new bioinformatic pipeline was developed to include additional information about alignment properties in haploblock predictions obtained by sequence comparisons. This tool was applied in RD^a-h2 to redefine

the SNP-haploblock to a new sequence-based haploblock in a more precise way, which could help increasing the correlation marker-trait. The new pipeline is user-friendly and can be applied to any input file containing alignments.

2.2. MATERIALS AND METHODS

2.2.1. BLAST calling of SNP-based RDM haplotypes among wheat assemblies

BLASTN, version 2.6.0+ (Zhang et al., 2000) was used with each of the 101 bp flanking sequences that Makhoul et al. (2020) used to develop KASP markers as query (Table 1) and the 15 publicly available wheat genome assemblies described in introduction as database in order to extract the SNP allele sequences from these blocks in each of the genotypes. These databases are available at the BLAST tool from EnsemblPlants, on the recently-released BLAST server IPK, Galaxy Blast Suite (<https://galaxy-web.ipk-gatersleben.de>), and can be downloaded at <https://webblast.ipk-gatersleben.de/downloads/wheat/pseudomolecules/>. Default parameters were applied.

For chromosome-level assemblies, hits were selected if their BLAST location was in chromosome 5B or in chromosome 7B in case of ArinaLrFor and SY Mattis, since these lines suffered a translocation in this chromosome arm between these two chromosomes (Miura et al., 1992). If there were more than one hit in these chromosomes, only the hit with the highest score was retained. For scaffold-level assemblies, the BLAST location of the hits were checked by blasting the context region against the Chinese Spring reference genome.

Table 1: Information about the SNP markers in RDM haploblocks discovered by Voss-Fels et al. (2017). Markers were considered significantly associated with RDM if their -log(p-value) was over the arbitrary threshold 4 and markers with the symbol – in this column were added to construct the block RDMy, where only one marker, Tdurum_con48959_1172, was found in LD with RDM. Notice that the marker Kukri_c46570_214 presents another homologous target upstream in the IWGSC chromosome 5B.

2.2.2. Identification of sequence-based haploblock predictions with crop- haplotypes.com

The website crop-haplotypes.com provides a graphical visualization of the shared sequence-based haploblock that were called from a combination of whole-chromosome NUCmer and gene-based BLAST pairwise alignments between wheat assemblies by Brinton et al. (2020). Each NUCmer alignment compared one reference against one query genome. However, scaffold-level assemblies could only be queries in the chromosome comparisons since their genomes are not organized in chromosomes so only NUCmer pairwise comparisons between chromosomes (CLA-CLA) or chromosome-scaffold (CLA-SLA) were conducted by Brinton and colleagues. Consequently, gene-based BLAST pairwise alignments were performed

between any two lines to allow for SLA-SLA comparisons and the genes compared with this method needed to be consistent with the expected chromosomes and have only one projection. Delta files were output of nucmer using the option -mum, to avoid repeats by keeping only maximal unique matches between the reference and the query, and alignments of length under 20 Kbp were filtered out with delta-filter (-l 20000) and specifying additionally -r -q, to retain only the longest consistent set of alignments for the reference and query. After the NUCmer pairwise comparisons, delta files were imported to R (R Core Team, 2021) where alignments were clustered into bins of fixed size (5-, 2.5- and 1-Mbp) across chromosomes (Brinton et al, 2020). Median percentage of identity was calculated across bins, which were assigned to haploblocks if their median percentage of identity was \geq 99.99. Consecutive assigned bins were assigned to the same haploblock and two errors (<99.99% identity) in a row within each haploblock were allowed. Regarding the BLAST-based method, alignments containing Ns were discarded and haploblocks were assigned by using a 25-gene sliding window. The 3 alignments with the lowest percentage of identity were also discarded in the windows and the mean percentage of identity was calculated across windows. Only 100% identical windows were assigned to blocks. Both sets of haploblocks defined from NUCmer or BLAST were combined and the longest predictions were selected in haploblocks called by both methods.

The RDMA block position (coordinates in chromosome 5B: Chinese Spring, 663,100,000-664,000,000 base pairs; Lancer, 655,700,000-656,600,000 bp) was searched on the website (crop-haplotypes.com/Wheat/haplotype/5B). All assemblies were compared at the RDMA region by selecting each of the available bin size views ('Haplotypes for 5Mbp', '[...] 2.5Mbp' and '[...] 1-Mbp'). The SNP-based RDMA block position was selected as the main highlighted bin, shown as dark grey boxes. If haploblocks were found in this bin, the haplotypes were analyzed throughout their whole length, meanwhile if no haploblock was found, other haploblocks were checked in the nearest bins. Notice that coordinates might differ between

identical regions in different assemblies due to insertion/deletion events. Finally, tables showing the haploblock coordinates in the lines sharing RDMa-h2 were produced (Fig. 3a).

2.2.3. Validation and redefinition of haploblocks

To assess the accuracy of the haploblock predictions called on crop-haplotypes.com, the scripts and files reported in Brinton *et al.* (2020) that were used for the creation of crop-haplotypes.com were applied. These original scripts are named ‘assign_mummer_blocks_whole_genome.r’ and ‘assign_BLAST_blocks_whole_genome.r’ and can be found at <https://github.com/Uauy-Lab/pangenome-web>. The original scripts were edited to create a new script in R (R Core Team, 2016), called ‘Haplotype-based Pangenome Analysis’ (hereafter referred to as HBPA).

Brinton’s method produced information on sequence-based haploblocks from bins of fixed size across wheat chromosomes but does not consider alignment properties during the haploblock calling process that could be essential to obtain accurate predictions. This is because, when clustering alignment into bins, the median percentage of identity was the only criterium used to assign bins into blocks. As a result, some predictions meeting the cutoff percentage of identity criterium could be based on poor number of alignments and therefore potential false positive or negative calls. HBPA provides adds new opportunities to Brinton’s method. First, HBPA performs a chromosome-scale analysis that checks for alignment properties across the whole target chromosome chromosome. The alignment properties considered in the analysis are (1) percentage of coverage of the reference chromosome by alignments (hereafter referred to as alignment coverage), e.g. the sum of the lengths of all alignments in a pairwise comparison divided by the chromosome length of the reference assembly, (2) average alignment length, (3) number of alignments and (4) average expected number of alignment per bin, e.g. total number of alignments divided by the number of bins of a fixed bin size that fit in the chromosome, supposing homogeneous alignment distribution. Alignment coverage depends is proportional to the next two properties, which makes it the most important indicator. These properties

provide an overview of the alignment quality of the pairwise comparison. Secondly, HBPA conducts a ‘target-scale’ analysis to validate haploblock predictions based on their alignment properties in chromosome regions of interest. For this, HBPA produces dotplots using the function ‘plot_aln_and_bins’, which contain graphical and numerical information on the haploblock predictions made in a selected region and its alignment coverage, as can be appreciated in Figure 5 and Figure 6. In this graph, the x axis marks the position of individual alignments, the y axis shows the alignment percentage of identity against the query and dot color indicates alignment length. Graphs are generated for the self-defined bin size and the bins are shown as the area between two consecutive vertical lines, for which the number of alignments is shown. Simultaneously, haploblock predictions from crop-haplotypes.com are highlighted in green and areas of interest, for example SNP-haploblocks, can also be highlighted on the graph. Therefore, this dotplot assists with haploblock validation and redefinition. Finally, HBPA extracts a list of the genes, e.g. projections and sources in IWGSC RefSeq v1.0, contained in haploblocks of interest. HBPA can be found at <https://github.com/MonteroJLU/Haplotype-based-Pangenome-Wheat-Analysis>, where several scripts that allow the replication of this project analyses or the personalized used of this method are hosted. The script containing the analysis of this study is named ‘demonstration_rdma_hbpa.r’. Any person interested in applying a similar analysis with personal data can use the script ‘template_hbpa.r’. However, reading the file ‘ReadMe.txt’ is recommended before downloading any of these resources.

HBPA is structured in three sections: (1) calling haploblocks from NUCmer pairwise comparisons, (2) Extracting a list of genes located within a haploblock and (3) calling haploblocks from BLAST pairwise comparisons. Notice the BLAST-based method was deemed insufficient to fully differentiate between haplotypes, with haplotypes called by NUCmer not being called by BLAST and viceversa.

For the analysis of the RDMA haplotype 2 shared between LongReach Lancer and Paragon, the raw data about alignments was imported into the script from a RDS file containing all filtered delta files from all NUCmer pairwise comparisons in chromosome 5B. This file is available at https://opendata.earlham.ac.uk/wheat/under_license/toronto/Brinton_et.al_2020-05-20-Haplotypes-for-wheat-breeding/nucmer/. The target-scale analysis was conducted on a target region between 655,000,000 and 665,000,000 in Lancer's chromosome 5B where the SNP-haplotype is located. For the gene identification, files were downloaded from https://webblast.ipk-gatersleben.de/downloads/wheat/gene_projection/. These included 'geneid_2_chinese.sourceid.txt', that matches gene model projections in assemblies with their source in the reference genome and 'projectedGenes__Triticum_aestivum_LongReach_Lancer_v1.0.gff', containing gene model projections in Lancer genome. For the third section, 'varieties_all_identities_2000bp.tab.gz' was downloaded from https://opendata.earlham.ac.uk/wheat/under_license/toronto/Brinton_et.al_2020-05-20-Haplotypes-for-wheat-breeding/pairwise_blast/ and 'iwgsc_refseq_v1.2_gene_annotation.zip' from https://urgi.versailles.inrae.fr/download/iwgsc/IWGSC_RefSeq_Annotations/v1.2/. These files contain the haplotype predictions obtained by Brinton and colleagues from gene-based BLAST pairwise alignments and the gene annotation in the reference genome respectively.

2.2.4. Statistical analysis of the differences in alignment coverage

To verify differences in alignment coverage caused by using scaffolds or chromosomes as query assemblies, two data groups were made based on the type of comparison, which is a relevant feature for Brinton's method: comparisons in which both reference and query are chromosome-level assemblies (CLA-CLA) and comparisons in which the reference is a chromosome- and the query is a scaffold-level assembly (CLA-SLA). To prove if there were significant

differences in alignment coverage between chromosomes, data was split into the 21 chromosomes of allohexaploid wheat (7x3). Functions were generated to calculate alignment coverage and can be found in the script ‘functions_hbpa.r’ at the github repository. The script generated for the data generation and statistical analysis can be found at this repository under the name ‘alignment_props_supplementary_hbpa.r’.

The percentage of coverage by alignment of the reference chromosome was calculated for the groups CLA-CLA and CLA-SLA as the alignment coverage mean values from all possible NUCmer pairwise comparison within each group and chromosome: 90 bidirectional pairwise comparisons for CLA-CLA (10 CLA x 10-1 CLA, as assemblies cannot act both as reference and query) and 50 unidirectional comparisons for CLA-SLA (10 CLA x 5 SLA). From the data relative to chromosomes, the average alignment coverage between the 21 chromosomes was calculated and compared between the CLA-CLA and CLA-SLA groups. Regarding the statistical analysis of the second hypothesis, only data from comparisons CLA-CLA was used and chromosome alignment coverage was compared between the 21 chromosomes.

Boxplots were generated plotting either the comparison type (CLA-CLA or CLA-SLA) or the analyzed chromosome against the percentage of coverage (Fig. 4). Analysis of residuals was conducted among the data to check prerequisites of the analysis of variance (ANOVA). Given the heteroscedasticity and non-normal distribution of the data, the Wilcoxon rank sum exact test was applied to the first hypothesis and Kruskal-Wallis to the second. To check differences in alignment coverage between chromosome 5B and the rest of the chromosomes, the Wilcoxon test was used.

2.2.5. Search for KASP marker targeting the haplotype RDMa-h2

The platform ‘Search sequence capture probes’ (https://www.cerealsdb.uk.net/cerealgenomics/CerealsDB/sequence_capture.php) was

consulted to search for all codominant SNP marker probes in chromosome 5B on the database. After selecting options and clicking on submit form, the website returned the names, sequences and the SNP positions of 497 probes. After that, a fasta was displayed and used as query for BLAST search with IWGSC Chinese Spring reference genome as database. Probes seemed to be clearly specific to chromosome 5B and the alignment bitscores in this chromosome tended to be higher than in its homoeologues. The BLAST output was generated in tabular format 6 and the table was read in RStudio, where only probes with chr5B alignments within the synonymous region in Chinese Spring chromosome 5B between 663,000,000 and 671,000,000 (equivalent to Lancer 655,000,000-663,000,000) were kept. These remaining probes were afterwards BLASTed against the pseudomolecules to filter out non-diagnostic markers (markers with alleles that differ between LongReach Lancer and Paragon). Locus-specificity was also pre-analyzed as recommended in Makhoul et al. (2020) by selecting the markers whose alignments in chromosome 5B showed a higher bitscore value than its counterparts in the homoeologous chromosomes in IWGSC Chinese Spring.

2.2.6. Identification of SNPs through high-stringency variant calling from Illumina paired-end short reads

The forward and reverse raw fastq files containing Illumina paired-end short reads generated by whole-genome shotgun (WGS) from three Chinese accessions reported to contain the RDM-positive allele T of the SNP marker Tdurum_con48959_1172 (Table 1) were downloaded from the website of the Genome Sequence Archive in BIG Data Center, accession number CRA001870 (Hao et al., 2020) The accession numbers of the three cultivars are: Guangtoushanmai, CRR062160; Yanmai 5, CRR062073 and Jimai 20, CRR062080.

Subsequently, a pipeline was developed on Ubuntu bash for variant calling from Illumina paired-end short reads and applied to each pair of forward and reverse fastq files. Pre-processing was performed with the tool cutadapt (Kechin et al., 2017). Reads under 30 bp, quality threshold

30 and more than 10% N were removed. The adaptor sequences of the kit TruSeq kit used for library preparation were found on the official Illumina webpage, 'Illumina Adapter Sequences' (<https://support-docs.illumina.com/SHARE/AdapterSeq/Content/SHARE/AdapterSeq/TruSeq/CDIndexes.htm>). Minimap2 (Li, 2018) was selected to map illumina paired-end reads against a fasta file containing the region between 652634893 and 653655509 bp on chromosome 5B of Chinese Spring RefSeq v1.0, region located between the SNP markers BS00029852_51 and BS00110293_51 (Table 1). The output SAM file was converted to BAM and quality-filtered with threshold 50 using samtools view (Li et al., 2019). To get rid of duplicated reads, these were first collated, e.g. grouped together by names, using samtools collate; followed by samtools fixmate to fill in mate coordinates and size fields; samtools sort to order reads and finally samtools markdup to mark duplicate alignments. After duplicate marking, the BAM file was filtered positively for paired reads (-f 1) and negatively for unmapped reads and mates (-F 12) using samtools view and refiltered for quality using the threshold 50. Alignment coverage was checked with mosdepth (Pedersen and Quinlan, 2018). Before variant calling, RG tags were added to the BAM file with picard AddOrReplaceReadGroups (<http://broadinstitute.github.io/picard/>) and a sequence dictionary was created with gatk-launch CreateSequenceDictionary for the fasta reference file. These steps were necessary for the next step involving gatk HaplotypeCaller (Ren et al., 2019) to call variants between the bam file generated from the short reads and the fasta file of the chromosome 5B fragment in Chinese Spring RefSeqv 1.0. Output was generated in GVCF format and the SNP marker Tdurum_con48959_1172 was found successfully when searching in each GVCF file by position (Table 1).

Each of the three individual GVCF files contained around 250k recorded variants. After merging the files, the number decreased to 413302, from which 139919 were SNPs. Multiallelic

SNPs were filtered out and only SNPs were retained in the merged file with bcftools view (Li, 2018), option --min-af 1:alt1 (minimum frequency of the alternative allele 1) --types=snps (only keep SNPs). This step reduced the number of SNPs to 7311. Minimum quality and read depth thresholds were set to 10000 and 100 respectively after thorough analysis applying intervals of values and checking the outcome number of SNPs. This filtering step produced 6105 SNPs. Output was cleared out of sites that had not been called in each of the Chinese varieties with grep, -v "./.", producing 568 SNPs. Finally, the same was conducted to eliminate SNPs with indels as alleles, leaving 542 SNPs.

The 101 bp flanking sequences of each SNP were extracted from the reference genome with samtools faidx. To generate the command lines for this operation, an R script (R Core Team, 2021) was generated that took in the text file and calculated the intervals between the sites located 50 bp downstream and upstream each SNP, printing these intervals inside the command lines that were later executed in Ubuntu. The 542 flanking sequences were merged and used as query for BLAST (Zhang et al., 2000) with each of the RDMb-h2-carrying assemblies as database (Julius, Norin, Stanley). The maximum number of hits per alignment was set to one so that only one hit with the maximum bitscore was produced on chromosome 5B for each SNP. Output tables were edited with grep to get rid of hits off chromosome 5B. Most of the aligned sequences between the SNP 101 bp flanking sequences from Chinese Spring genome and targets in chromosome 5B from these RDMb-h2-carrying assemblies were 100% identical.

The text files files were imported to R and merged into a single table. SNP flanking sequences that had produced <100% identical alignments were re-BLASTed to check if the missaligned nucleotide was located in the SNP allele. 4 SNPs were different in the allele position containing the allele for RDMb-h3 so finally 538 SNPs were exported to a table (Supplementary File 2).

2.3. RESULTS AND DISCUSSION

2.3.1. Two wheat assemblies contain a SNP-haplotype that increases root dry mass

In order to retrieve genomic information from genotypes carrying favorable RDM haplotypes, the flanking sequences of the 9 and 6 SNP chip markers in Table 1 that define each of two RDM haploblocks, RDMA and RDMB respectively, were searched in wheat genome databases. BLAST revealed that the RDM-associated haplotype h2 in RDMA was present in the assemblies LongReach Lancer and Paragon, meanwhile Hap-5B-RDMA-h1 was carried by the remaining

| Hap-5B-RDMA haplotypes | | ALLELES | | | | | | | | ASSEMBLIES | |
|------------------------|--|------------------|--------------------|----------------------|-------------------|-----------------|----------------------|-------------------|-------------|----------------------|---|
| | | Kukri_c46570_214 | RAC875_c24226_1356 | RAC875_c18088_2222 | RAC875_c18088_950 | BobWhite_c43_86 | Excalibur_c25522_755 | RAC875_c12293_588 | GENE-2890_- | Excalibur_c60554_394 | |
| RDMA-h1 | | A | A | G | C | G | C | G | G | G | ARI, JAG, JUL, LDM, MAC, NOR, STA, MAT, CSP, CAD, CLA, ROB, WEE |
| RDMA-h2 | | G | C | A | T | A | T | A | A | A | LAC, PAR |
| Hap-5B-RDMB haplotypes | | ALLELES | | | | | | | | ASSEMBLIES | |
| RDMB-h1 | | BS00022477_51 | BS00029852_51 | Tdurum_con48959_1172 | BS00110293_51 | BS00022231_51 | IACX6288 | | | | ARI, JUL, LAC, LDM, MAC, MAT, CAD, CLA, PAR, ROB, WEE |
| RDMB-h2 | | G | T | G | G | A | C | | | | CSP, JAG, NOR, STA |
| RDMB-h3 | | G | T | T | G | G | C | | | | |
| RDMB-h8 | | A | G | T | G | G | C | | | | |

Table 2: SNP-haplotype pattern for the RDM haploblocks among the 15 wheat genome assemblies: ARI, ArinaLrFor; JAG, Jagger; JUL, Julius; LAC, Lancer; LDM, Landmark; MAC, Mace; NOR, Norin61; STA, Stanley; MAT, SY-Mattis; CSP, Chinese Spring; CAD, Cadenza; CLA, Claire; PAR, Paragon; ROB, Robigus; WEE, Weebill. Haplotypes with positive effects on RDM are filled in green. Haplotype RDMB-h3 was not found to be present in any of the assemblies.

13 lines. Regarding the adjacent haplotype block RDMb, the RDM-associated haplotype h3 was not found in any of the assemblies, which contained RDMb-h1 or h2 instead (Table 2).

2.3.2. A sequence-based haplotype matches with the SNP-based RDMA-h2 according to crop-haplotypes.com

To confirm if the SNP-haplotypes identified by BLAST match with the sequence-based haplotypes produced by Brinton et al. among the 15 wheat assemblies, crop-haplotypes.com was used. Shared haplotypes were predicted by the website between the RDMA-h2-carrying assemblies, Lancer and Paragon in the proximity of the SNP-haplloblock RDMA at the three different bin sizes (5-, 2.5- and 1-Mbp), as shown in Figure 3a and 3b. However, the first two predictions did not overlap with the SNP-haplloblock and only the prediction at 1-Mbp did, as it started from 656 Mbp in Lancer's chromosome 5B. Haplloblock length increased at 1-Mbp compared with 5- and 2.5-Mbp as an effect of the increase in resolution. Between 5- and 2.5-Mbp, the lengths remains identical although the haplloblock position changes, which can be seen in Figure 1b with the green bar moving to the right from 5- to 2.5-Mbp. Therefore, the reported sequence-based haplotype was hypothesized to be the same as the SNP-haplotype RDMA-h2. It was also considered that the haplloblock could be larger than the original SNP-haplloblock. Nevertheless, since only one of the three predictions overlapped the SNP-haplloblock, it was considered convenient to validate this hypothesis following the original Brinton's method.

Regarding the 13 assemblies carrying RDMA-h1, the website revealed a complex haplotype pattern in this block. This means that many different shared and unique haplotypes, e.g. haplotypes that are not shared between at least two lines, were observed around the RDMb block position (Chinese Spring, 652,600,000-653,700,000 bp), as can be appreciated in Figure 3c. If the RDMA-h1 is actually conserved between the assemblies, it could have been undetected because of (1) its length under the minimum bin size, 1-Mbp, or (2) that the small region is not

identical-by-state ($\geq 99.99\%$ median sequence identity) between these assemblies but nearly-identical ($<99.99\%$). The latter hypothesis would correspond with the higher sensibility of Brinton's method to detect small differences between alleles, although in this case, these differences would not be relevant for breeding, since this haplotype is not associated with RDM.

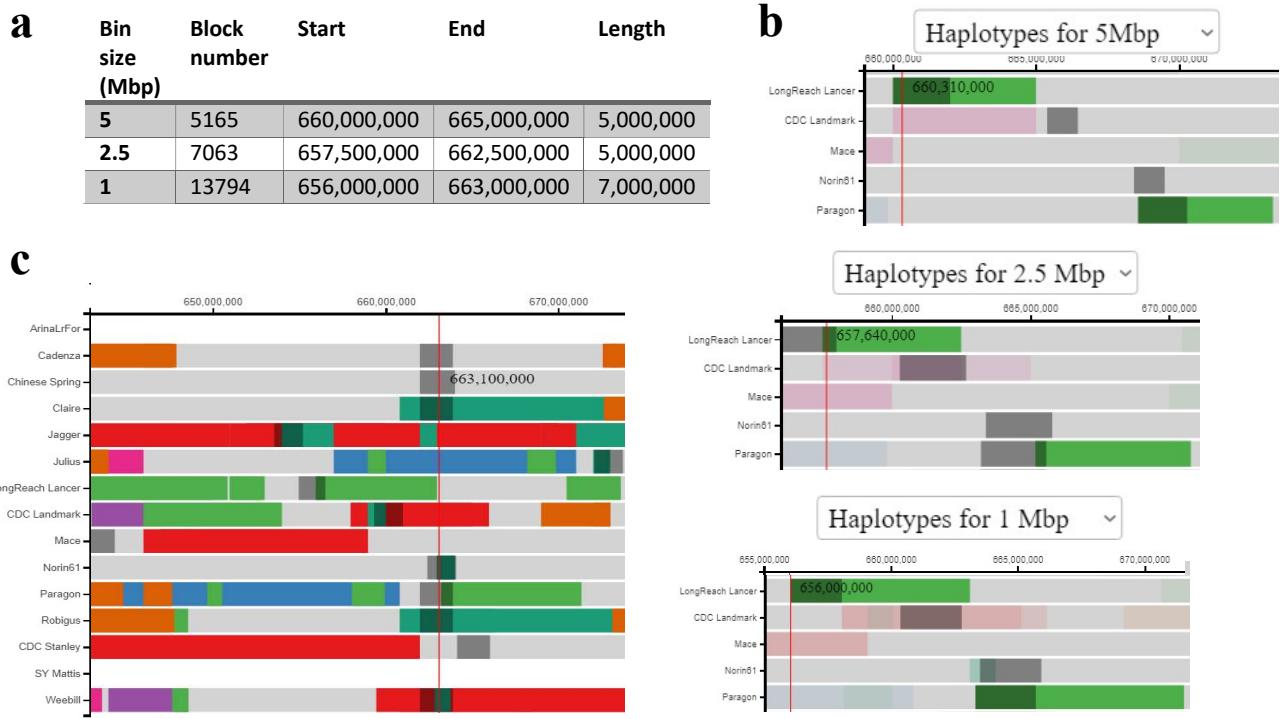


Figure 3: Haplotype predictions on crop-haplotypes.com. (a) Haplotype shared between LongReach Lancer and Paragon in chromosome 5B near the SNP-haplotype RDMA. Coordinates relative to Lancer. (b) Graphical visualization from the data in a. Only the haplotype prediction at bin size 1-Mbp includes part of the SNP-haplotype RDMA (655,700,000-656,600,000). Coordinates relative to Lancer. (c) Haplotypes among the 15 wheat genome assemblies under 1-Mbp bin size at the start position of RDMA in IWGSC Chinese Spring (663,100,000). The highlighted dark grey boxes are the homologous of the region around 663 Mbp in Chinese Spring in every assembly.

2.3.3. Using scaffolds as query in NUCmer pairwise alignments reduces alignment coverage

The observations during the analysis of the RDMA haplotype raised two questions: (1) whether using scaffold-level assemblies as query had a negative effect on alignment coverage by NUCmer pairwise comparisons, which could affect the reliability of haplotype predictions and (2) if the differences in alignment coverage depend on the analyzed chromosome.

Regarding the first hypothesis, the boxplot showed an apparently much lower alignment coverage in comparisons CLA-SLA than in comparisons CLA-CLA (Fig. 4a). This observation

was proved significant using the Wilcoxon rank sum exact test (p-value 7.431e-12). Two reasons explaining why using scaffolds in NUCmer pairwise comparisons produces low alignment coverage and thus inaccurate predictions are: (1) lower sequencing quality than their chromosome-level counterparts and (2) the lack of chromosome-level organization, since this would facilitate the aligner to find the right target, which would otherwise be challenged to find the this sequence among homologous sequences in the whole scaffold set without any indication. This problem is not indicated in Brinton et al. and they only rely on the gene-based BLAST pairwise comparisons to alleviate it despite their reported insufficiency to fully differentiate between haplotypes. Some solutions to improve the alignment quality of these comparisons would be: (1) to further assemble the scaffold-level genomes to the chromosome-level or (2) to sort out the scaffolds into chromosomes by for example BLASTing them against the reference genome.

Regarding the second hypothesis, the boxplot was sort by increasing median value and chromosome 5B was observed to have the lowest median, so the hypothesis was reshaped to ask if alignment coverage in chromosome 5B differed to those from other chromosomes (Fig. 4b). The analyzed chromosome seemed to have a signficant effect on alignment coverage (p-value 2.2e-16), according to Kruskal-Wallis test. Additionally, the alignment coverage of chromosome 5B in NUCmer comparisons between chromosomes was significantly lower than in every other chromosome (p-value 0.0325), as Wilcoxon test revealed. These differences in alignment coverage between chromosomes could be due to differences in sequencing efforts

that could translate in some chromosomes having presenting more sequencing errors or Ns that difficult the aligning process.

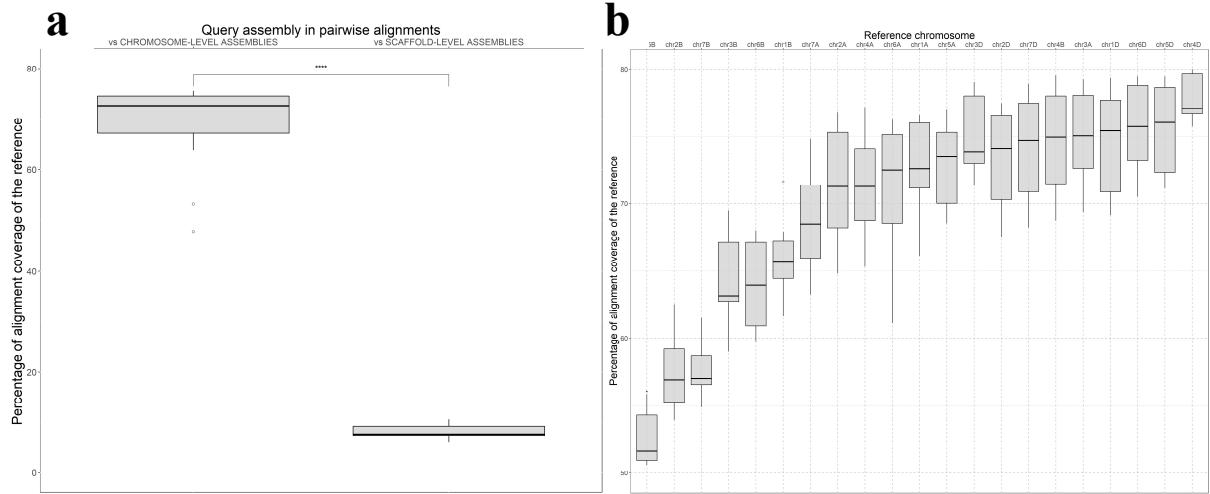


Figure 4: Boxplots showing differences in the percentage of alignment coverage of the reference chromosome, between (a) comparisons in which the query is a chromosome-level assembly (left) or a scaffold-level assembly (right), which have very significantly lower alignment coverage (p -value X, **), and between (b) comparisons chromosome-chromosome in different analyzed chromosomes, among which the chromosome 5B has a significantly lower alignment coverage than any of them (p value X, *). Significance thresholds: * ($P = 0.05$), ** ($P = 0.01$), *** ($P = 0.001$), **** ($P = 0.0001$).**

2.3.4. Checking alignment properties is essential for accurate haploblock prediction

When analyzing the dotplot generated with the R function ‘plot_aln_and_bins’, a large discontinuity at 663.0-664.5 Mbp containing only two <90.0% identical alignments was noticeable (Fig. 5a). Despite of this, the region between 660 and 665 Mbp had been assigned to a haploblock at 5-Mbp bin size given that the median percentage of identity of the 5-Mbp bin where it is located laid over the cut-off 99.99%. Discontinuities contain missing information from alignments and can increase false positive and negative haplotype predictions affecting the power to detect haploblocks and the precision to make right predictions, since the algorithm from Brinton’s method does not take low alignment coverage into account.

This negative situation could be common in NUCmer pairwise alignments with low alignment coverage, as reported previously about the comparisons between chromosome- and scaffold-level assemblies. In this case, discontinuities could be due to sequencing gaps or alignment errors due to the difficulty of aligning assemblies with different coordinate systems

(chromosome vs scaffolds). The latter would explain why these discontinuities are often associated with alignments of relatively low percentage of identity. However, given the differences in alignment coverage between analyzed chromosomes, these discontinuities could be also expected in comparisons between reference-quality assemblies, especially under heterogeneous distribution.

Despite this wrong assignment at 5-Mbp bin size, reducing bin size to 2.5-Mbp can overturn this error by excluding the segment between 662.5 and 665.0 Mbp containing this discontinuity (Fig. 5b). However, this prediction extends over another discontinuity between 658,0 and 659,0 Mbp containing only one alignment of 99.75% identity that remains when the bin size is further decreased to 1-Mbp (Fig. 5a), due to the ‘two consecutive errors’ rule by Brinton et al. (2020) discussed in the introduction. It seems difficult to assess if this error is actually a false positive, e.g. a wrong haploblock assignation in a region that differs between the compared assemblies, or if the low coverage and identity are simply due to technical factors such as sequencing errors.

A different scenario in which considering alignment properties is relevant would be when haplotypes are called based on only a few alignments in comparisons in which haploblocks are not expected, illustrated by one case in Norin61’s chromosome 5B (Fig. 6). Here, this assembly mainly shares haplotypes only with Mace at bin sizes 5- and 2.5-Mbp, but at 1-Mbp, many new haplotypes appear mainly in comparisons with scaffold-level assemblies. One of these blocks, at 537-538 Mbp (Norin61) was called due to the presence of only one alignment with $\geq 99.99\%$ identity in a 1-Mbp bin meanwhile the region around this bin contains many empty bins or alignments with low percentage of identity. Furthermore, Norin61 is with Chinese Spring one of the pre-1950 line and thus clearly different to the rest of the assemblies (Brinton et al., 2020). These factors lead to think that this is a false positive prediction, which was probably caused by misalignment with a very similar off-target region

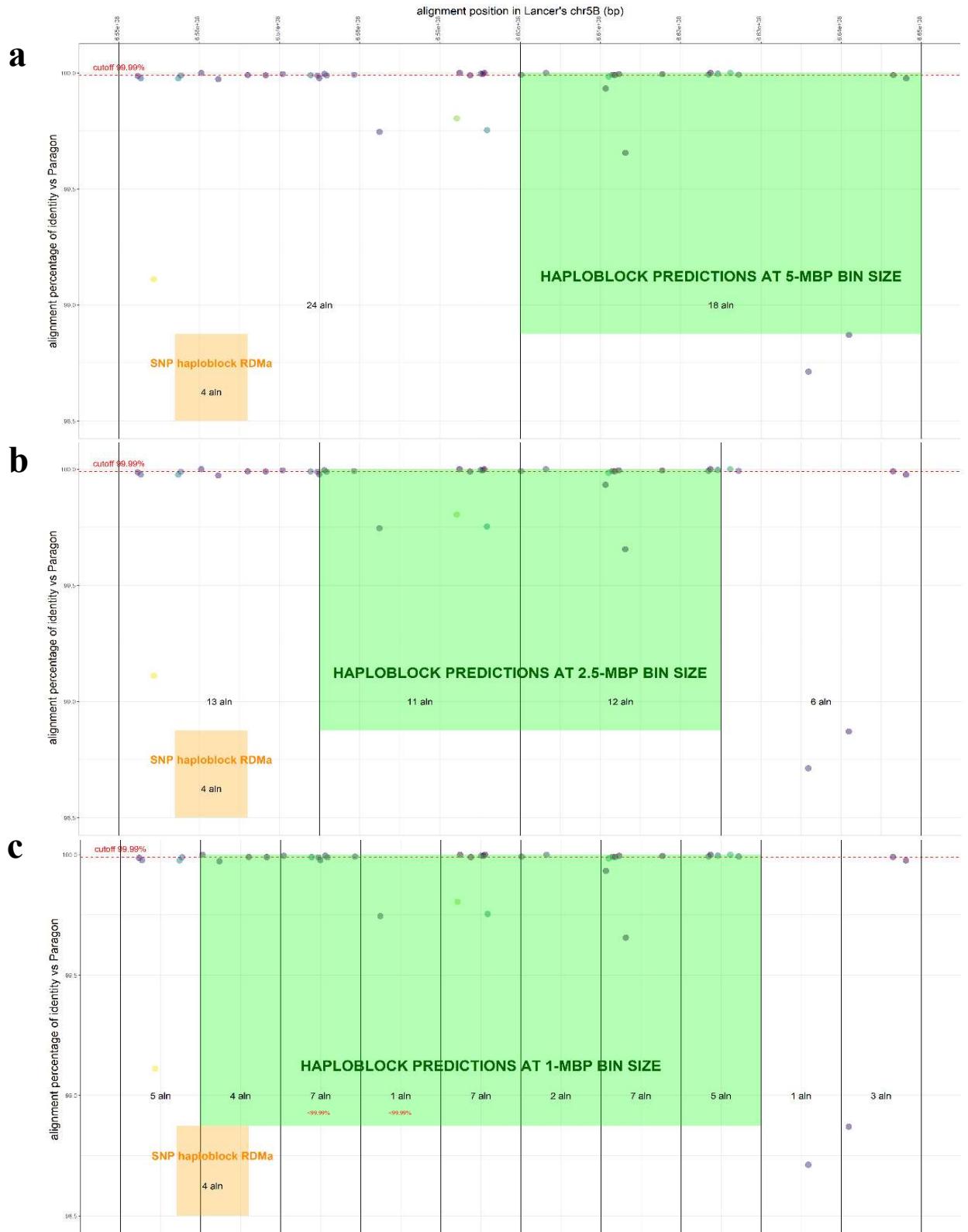


Figure 5: Effects of changing bin size on haploblock predictions on crop-haplotypes.com in the region between 655 and 665 Mbp of Lancer's chromosome 5B, at (a) 5-Mbp, (b) 2.5-Mbp and (c) 1-Mbp bin size. Dots indicate Lancer's alignment position in the x axis and percentage of identity against Paragon's scaffolds in the y axis. Simultaneously, the bins of fixed size used in Brinton's method for haplotype detection are marked with vertical lines and the number of alignments within each bin is shown between consecutive lines. Haplloblock predictions generated by Brinton's method and that can be observed on crop-haplotypes.com are highlighted in green and titled and the SNP-haplloblock where RDMa-h2, shared between Lancer and Paragon, is located is also highlighted in orange. This graph was generated by the function 'plot_aln_and_bins' in the R script 'functions_hbpa.R'.

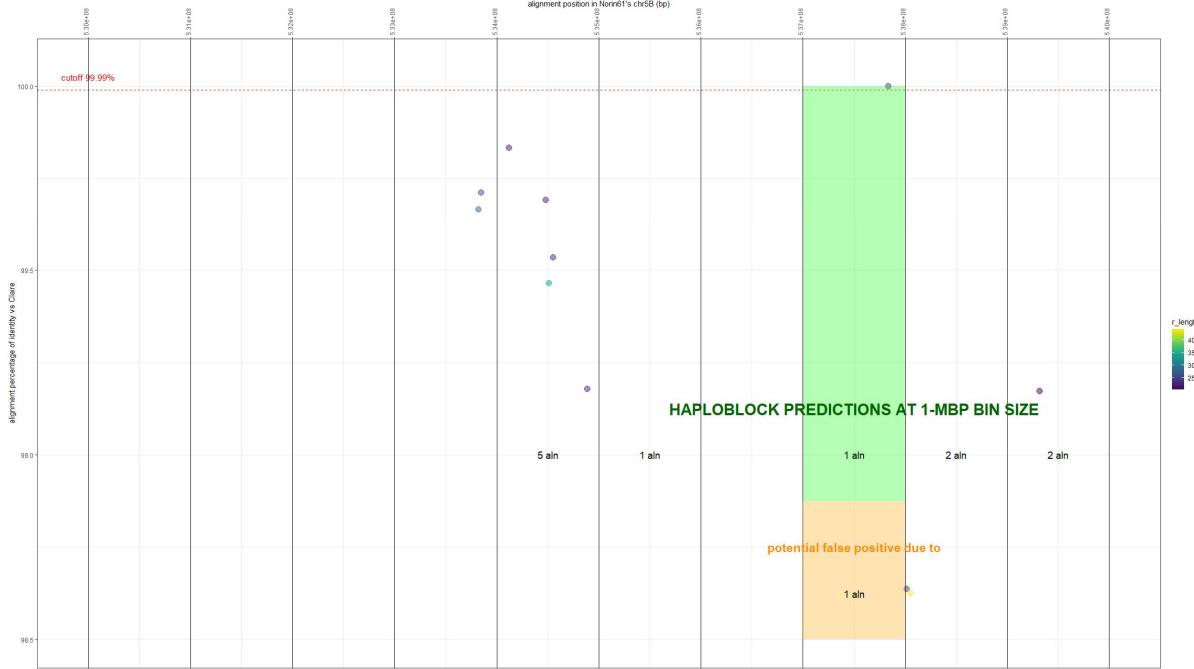


Figure 6: Potential wrong assignment of the Norin’s chromosome 5B region between 537 and 538 Mbp to a 1-Mbp-long haploblock at 1-Mbp bin size on crop-haplotypes.com. Dots indicate Norin’s alignment position in the x axis and percentage of identity against Claire’s scaffolds in the y axis. Simultaneously, the bins of fixed size used in Brinton’s method for haplotype detection are marked with vertical lines and the number of alignments within each bin is shown between consecutive lines (no message is printed if bins are empty). Haplloblock predictions generated by Brinton’s method and that can be observed on crop-haplotypes.com is highlighted in green and titled and the same haploblock is also highlighted in orange. This graph was generated by the function ‘plot_aln_and_bins’ in the R script ‘functions_hbp.R’.

2.3.5. Low alignment coverage and discontinuities hamper validation and redefinition of the haploblock RDMA

The chromosome-scale analysis pointed out that the NUCmer pairwise comparison between Lancer and Paragon was defined by low-quality alignment properties. Only 11.64% of the 702438406 bp of Lancer’s chromosome 5B are covered by alignments against the scaffolds from Paragon, which is very low in comparison with its average coverage 55.81% in alignments between Lancer and any chromosome-level assembly. This lower alignment coverage was due to shorter alignment average lengths (~31 Kbp versus ~47 Kbp) and lower number of alignments (~2.7 K vs ~8.3 K) when Lancer is aligned versus Paragon than when it is aligned versus any CLA. From these values, the average number of alignments that we can expect in 1-Mbp bins was 3.8 in Lancer-Paragon, compared with 11.9 in Lancer against any chromosome.

As reported previously, these features are characteristic not only of this comparison but generally of any comparison between chromosomes and scaffolds.

When the region between 655,000,000 and 665,000,000 bp in Lancer's chromosome 5B was analyzed with the function 'plot_aln_and_bins', the existence of two discontinuities at 658,0-659,0 Mbp and 663,0-664,5 Mbp was evident. At 1-Mbp bin size, the haploblock prediction avoids the second discontinuity and overlaps with the SNP-haploblock RDMA (655,700,000-656,600,000). However, the haploblock prediction still spans across the first discontinuity, which cuts it into the two regions at 656-658 and 659-663 Mbp, which seem to have a good balance between alignment density and high percentage of identity.

Therefore, in order to redefine haploblock RDMA based on the information obtained from sequence comparison, we could adopt one of the following four approaches (Fig. 7):

1. To completely invalidate the haploblock predictions due to their low-quality alignment properties and therefore not attempt to redefine it. This conservative approach would suppose to keep the SNP-haploblock RDMA as the target region for RDM breeding in wheat: 655,700,000-656,600,000.
2. To consider discontinuities within haploblocks as wrong assignations explained by differences between the lines. This would mean to expand the haploblock in upstream direction but stopping before the first discontinuity: 656,000,000-658,000,000.
3. To redefine the haploblock as the region after the discontinuity, considering that this might be where the causal gene actually lies and that the SNP-haploblock RDMA is only linked by proximity with this region: 659,000,000-663,000,000.
4. To understand discontinuities within haploblocks as identical regions that are incorrectly reported as non-identical due to sequencing or alignment errors. This less conservative

approach would mean to completely validate the prediction at 1-Mbp bin size: 656,000,000-663,000,000.

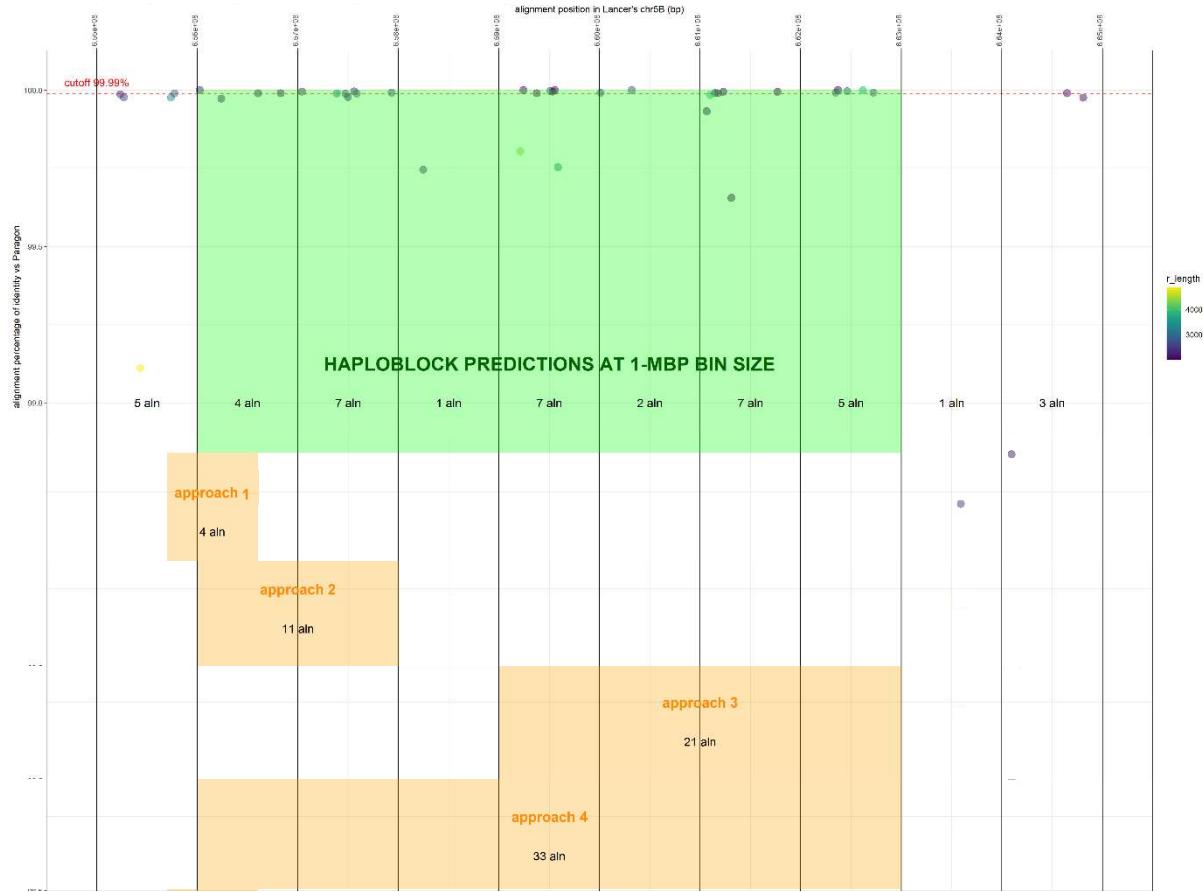


Figure 7: Comparison between the four approaches for haploblock redefinition. This graph was edited from four dotplots generated by the function ‘plot_aln_and_bins’, in the script ‘functions_hbpa.R’.

2.3.6. New KASP markers from an array could target the potentially larger haploblock RDMA-h2

Since RDMA could be larger than the original SNP region (0.9-Mbp-long), as found out in previous steps, molecular markers were searched to target a new region that could be associated with the haploblock. 5 codominant SNP probes were discovered among a raw list of 39 Axiom SNP probes flanking a region between 658,499,200 and 662,451,059 in LongReach Lancer’s chromosome 5B following these restrictive search criteria (Table 3). 34 more probes were added to a raw list of markers located in the candidate region. Ideally, these probes could be

used to develop KASP markers to genotype wheat plants for association studies, which could help by increasing trait-marker correlation for the haplotype Hap-5B-RDMA-h2.

| Axiom SNP marker name | AX-94878590_116_Co-dominant_pos-61_v1 | AX-95255236_263_Co-dominant_pos-61_v1 | AX-95173034_464_Co-dominant_pos-61_v1 | AX-94636419_373_Co-dominant_pos-61_v1 | AX-94994996_195_Co-dominant_pos-61_v1 |
|--------------------------------|---|--|---|---|--|
| SNP probe | ATGTCCTTACATGATGGCTCG AAGGTTGCTTCTTGTGGAACTTG TGCCATTGATGAAAGAAGC GGTTTGATGGCTTCTTCATGGCTT GGTTTGATGGCTTCTTCATGGCTT GGCTTA | TGTTGGGCAAACTCTTCACT TGCGGATCTTCTTCATGGCTT TGCGGATCTTCTTCATGGCTT TGCGGATCTTCTTCATGGCTT TGCGGATCTTCTTCATGGCTT GTGTTG | CTGCACTACTTAAAGCTTGCTG GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT ATATG | GTGCACTACTTAAAGCTTGCTG GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT AACAC | GTGCACTACTTAAAGCTTGCTG GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT GTTTGAATCTTCTTCATGGCTT GCCT |
| BLAST position (IWGSC) | 666.209.099 | 666.209.246 | 667.839.964 | 668.473.987 | 670.831.214 |
| BLAST position (Lancer) | 658.499.200 | 658.499.347 | 659.290.270 | 659.893.416 | 662.451.059 |
| LONGREACH LANCER | C | G | T | G | G |
| PARAGON | C | G | T | G | G |
| ArinaLrFor | C | G | T | C | T |
| Jagger | C | G | C | G | T |
| Mace | C | G | C | G | T |
| Norin61 | C | G | C | G | G |
| Julius | C | G | T | C | T |
| Stanley | T | A | - | C | G |
| CDC Landmark | C | G | C | G | T |
| SY Mattis | C | G | C | G | T |
| Claire | C | G | C | G | G |
| Weebill | C | G | C | G | T |
| Cadenza | C | G | C | G | T |
| Robigus | C | G | C | G | G |
| IWGSC - chr5B | T | A | T | C | G |
| IWGSC - Chr5A | C | g | - | g | g |
| IWGSC - Chr5D | C | g | - | g | g |

Table 3: Candidate codominant SNP markers to target the redefined potential RDMA haploblock. Identifier, SNP probe, position in Lancer and Chinese Spring and allele sequences in 15 wheat genome assemblies are provided in the table. SNP probes were obtained from the Axiom Pannel at cerealsdb.uk.net. Notice that the RDMA-h2-carrying lines Lancer and Paragon share the same allele sequence and among the other lines, some of them share the same allele (red) and some of them possess the other allele (green) of the same markers. In the last three rows, the allele sequence is displayed for Chinese Spring in its three homoeologues chromosomes 5 (5A, 5B, 5D), showing how alleles differ between chromosome 5B and the other two off-target chromosomes. Lower case indicates that this hit does not have the highest bitscore among the three hits in the homoeologous chromosomes.

2.3.7. Genome sequences of RDMb-h3-carrying varieties were identified and provided new SNP markers

As the RDM haplotype RDMb-h3 was not detected in any of the available wheat genome assemblies, it needed to be searched in public databases. In 2020, Hao et al. conducted the resequencing of 145 Chinese cultivars and captured more than 60.92 million SNPs, based on Chinese Spring RefSeq v1.0. Upon request, the SNP marker Tdurum_con48959_1172 was search and found among the raw SNP set generated in the experiment in three Chinese varieties: one landrace, Guangtoushanmai and two modern cultivars, Yangmai 5, Jimai 20. However, this SNP was not found among the highly reliable SNP set, unlike the other 5 markers that define RDMb. This finding provided for the first time genome sequences of the very infrequent haplotype RMDb-h3 (Voss Fels et al., 2017). Furthermore, high-stringency mapping of the Illumina paired-end short reads and variant calling between Chinese Spring (RDMb-h2) and the three Chinese varieties (RDMb-h3) provided 538 SNPs in the highly-reliable SNP set (Supplementary File 2) obtained with very stringent filtering. Tdurum_con48959_1172 was found in all raw GVCF files but removed after filtering due to lower quality and read depth. These SNP sequences could be confirmed by Sanger sequencing to develop diagnostic KASP markers to select RDMb-h3. Furthermore, these genome sequences could be further analyzed to identify the underlying genes causing higher RDM or to find out if any of these accessions also carries the haplotype RDMA-h2, since the combination RDMA-h2/RDMb-h3 is associated with the highest phenotypic effects (Voss-Fels et al., 2017). Also, it could be studied if the sequence similarity between RDMb-h2 and RDMb-h3 is due to a recombination event between the markers Tdurum_con48959_1172 and BS00029852_51. This could have been possible in this study if long reads had been available for these three varieties since assembling these sequences enables the analysis by comparison of potential recombination breakpoints between RDMb-h2 and h3.

3. GENERAL DISCUSSION

Pangenomics provide a groundbreaking opportunity to understand and harness genetic diversity in wheat. According to Rasheed et al. (2017), genetic resources should be prioritized in the post-genomics era for the improvement of wheat production in order to double yield by 2050.

The on-going international efforts on wheat pan-genome sequencing succeeded with the publication of reference-quality wheat assemblies and scaffold-level assemblies (Walkowiak et al., 2020). The sequence-based haploblock prediction method developed by Brinton et al. (2020) will contribute to increasing precision of wheat breeding by redefining traditional SNP-based haploblocks. This breakthrough could not have been possible without the recent abovementioned wheat pangenomic resources.

In this master thesis, new contributions were made to Brinton's method after the realization that alignment properties must be considered for calling haploblocks from sequence comparison. The newly-developed R script allows target-scale analysis of shorter chromosome fragments of interest and provides graphic and numeric information about the alignment properties in this region that can be used for the validation of the haploblock predictions observed in crop-haplotypes.com. Furthermore, it enables the simple manipulation of parameters, such as bin size or cutoff, to personalize the analysis given the clear usage instructions provided along the script and the high degree of detail provided by the graphs. The script is not only constricted to the original raw data but can instead be used to analyze new data sets. Since wheat pangenomics does not have a long history and there are still not enough available resources, these tools could be very important for precision breeding.

Furthermore, this work achieved to find genome sequences for the haploblock RDMb-h3 for the first time. The genomic resources were used to identify new diagnostic SNP markers that

could potentially be used for breeding wheat with higher root dry mass and are still subject to further interesting analyses.

4. SUMMARY

Haploblocks can be defined by comparing DNA sequences in wheat instead of by following the traditional method based on SNP markers. Sequence-based methods have the advantage of higher precision in finding regions that are causal for traits and higher accuracy in detecting small differences between highly identical alleles. SNP-based haplotypes can be redefined to sequence-based haplotypes if sequence information from the haplotype-carrying varieties is available. In this article, a SNP-based haploblock associated with higher root dry mass was analyzed among 15 wheat genome assemblies. To achieve this, a new method based on pairwise alignments between assemblies was used to confirm that two wheat lines carrying one of the favorable SNP-based haplotype also share a sequence-based haplotype at this region. Subsequently, the fundaments of this method were explored in an attempt to redefine the SNP-based haploblock to a more precise sequence-based region. The relevance of alignment properties in haploblock predictions was highlighted and bioinformatic resources were developed to include the analysis of these properties into the method to detect haplotype from pairwise alignments with increased precision. Additionally, genome sequences were found in public databases for another favorable haplotype and used to identify new SNPs to target this haplotype.

5. REFERENCES

Barrett, J. C., Fry, B., Maller, J., & Daly, M. J. (2005). Haplovew: analysis and visualization of LD and haplotype maps. *Bioinformatics* (Oxford, England), 21(2), 263–265.
<https://doi.org/10.1093/bioinformatics/bth457>

Bernardo, R. N. (2010). Breeding for Quantitative Traits in Plants. Woodbury, MN: Stemma Press.

Boose, D., Harrison, S., Clement, S., & Meyer, S. (2011). Population genetic structure of the seed pathogen Pyrenophora semeniperda on Bromus tectorum in western North America. *Mycologia*, 103(1), 85–93. <https://doi.org/10.3852/09-310>

Brinton, J., Ramirez-Gonzalez, R. H., Simmonds, J., Wingen, L., Orford, S., Griffiths, S., 10 Wheat Genome Project, Haberer, G., Spannagl, M., Walkowiak, S., Pozniak, C., & Uauy, C. (2020). A haplotype-led approach to increase the precision of wheat breeding. *Communications biology*, 3(1), 712. <https://doi.org/10.1038/s42003-020-01413-2>

Casebow, R., Hadley, C., Uppal, R., Addisu, M., Loddo, S., Kowalski, A., Griffiths, S., & Gooding, M. (2016). Reduced Height (Rht) Alleles Affect Wheat Grain Quality. *PloS one*, 11(5), e0156056. <https://doi.org/10.1371/journal.pone.0156056>

Collard, B. C., & Mackill, D. J. (2008). Marker-assisted selection: an approach for precision plant breeding in the twenty-first century. Philosophical transactions of the Royal Society of London. Series B, *Biological sciences*, 363(1491), 557–572. <https://doi.org/10.1098/rstb.2007.2170>

El Baidouri, M., Murat, F., Veyssiere, M., Molinier, M., Flores, R., Burlot, L., Alaux, M., Quesneville, H., Pont, C., & Salse, J. (2017). Reconciling the evolutionary origin of bread wheat (*Triticum aestivum*). *The New phytologist*, 213(3), 1477–1486. <https://doi.org/10.1111/nph.14113>

Ewert, F., Rounsevell, I., Reginster, I. , Metzger M. & Leemans, R. (2005). Future scenarios of European agricultural land use. I. Estimating changes in crop productivity. *Agriculture, Ecosystems and Environment* 107: 107-116.

Gabriel, S. B., Schaffner, S. F., Nguyen, H., Moore, J. M., Roy, J., Blumenstiel, B., Higgins, J., DeFelice, M., Lochner, A., Faggart, M., Liu-Cordero, S. N., Rotimi, C., Adeyemo, A., Cooper, R., Ward, R., Lander, E. S., Daly, M. J., & Altshuler, D. (2002). The structure of haplotype

blocks in the human genome. *Science* (New York, N.Y.), 296(5576), 2225–2229.

<https://doi.org/10.1126/science.1069424>

Giraldo, P., Benavente, E., Manzano-Agugliaro, F., Gimenez, E. (2019) Worldwide Research

Trends on Wheat and Barley: *A Bibliometric Comparative Analysis*. *Agronomy*: 9(7):352.

<https://doi.org/10.3390/agronomy9070352>

Hao, C., Jiao, C., Hou, J., Li, T., Liu, H., Wang, Y., Zheng, J., Liu, H., Bi, Z., Xu, F., Zhao, J., Ma,

L., Wang, Y., Majeed, U., Liu, X., Appels, R., Maccaferri, M., Tuberosa, R., Lu, H., & Zhang,

X. (2020). Resequencing of 145 Landmark Cultivars Reveals Asymmetric Sub-genome

Selection and Strong Founder Genotype Effects on Wheat Breeding in China. *Molecular plant*,

13(12), 1733–1751. <https://doi.org/10.1016/j.molp.2020.09.001>

Hartung, F., & Schiemann, J. (2014). Precise plant breeding using new genome editing techniques:

opportunities, safety and regulation in the EU. *The Plant journal : for cell and molecular*

biology, 78(5), 742–752. <https://doi.org/10.1111/tpj.12413>

Hong, F., Gao, L., Han, H. L., Wang, P., Wang, J., Wei, D., & Liu, Y. (2019). Population Genetics

of Bactrocera minax (Diptera: Tephritidae) in China Based on nad4 Gene Sequence. *Insects*,

10(8), 236. <https://doi.org/10.3390/insects10080236>

International Wheat Genome Sequencing Consortium (IWGSC) (2014). A chromosome-based draft

sequence of the hexaploid bread wheat (*Triticum aestivum*) genome. *Science* (New York,

N.Y.), 345(6194), 1251788. <https://doi.org/10.1126/science.1251788>

Kechin, A., Boyarskikh, U., Kel, A., & Filipenko, M. (2017). cutPrimers: A New Tool for Accurate

Cutting of Primers from Reads of Targeted Next Generation Sequencing. *Journal of*

computational biology : a journal of computational molecular cell biology, 24(11), 1138–

1143. <https://doi.org/10.1089/cmb.2017.0096>

Li H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* (Oxford, England), 34(18), 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., & 1000 Genome Project Data Processing Subgroup (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics* (Oxford, England), 25(16), 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>

Lin, M., Cai, S., Wang, S., Liu, S., Zhang, G., & Bai, G. (2015). Genotyping-by-sequencing (GBS) identified SNP tightly linked to QTL for pre-harvest sprouting resistance. TAG. *Theoretical and applied genetics. Theoretische und angewandte Genetik*, 128(7), 1385–1395. <https://doi.org/10.1007/s00122-015-2513-1>

Makhoul, M., Rambla, C., Voss-Fels, K. P., Hickey, L. T., Snowdon, R. J., & Obermeier, C. (2020). Overcoming polyploidy pitfalls: a user guide for effective SNP conversion into KASP markers in wheat. TAG. *Theoretical and applied genetics. Theoretische und angewandte Genetik*, 133(8), 2413–2430. <https://doi.org/10.1007/s00122-020-03608-x>

Matsumoto, H., & Kiryu, H. (2013). MixSIH: a mixture model for single individual haplotyping. *BMC genomics*, 14 Suppl 2(Suppl 2), S5. <https://doi.org/10.1186/1471-2164-14-S2-S5>

Miura, H., Parker, B. B., & Snape, J. W. (1992). The location of major genes and associated quantitative trait loci on chromosome arm 5BL of wheat. TAG. *Theoretical and applied genetics. Theoretische und angewandte Genetik*, 85(2-3), 197–204. <https://doi.org/10.1007/BF00222860>

Pedersen, B. S., & Quinlan, A. R. (2018). Mosdepth: quick coverage calculation for genomes and exomes. *Bioinformatics* (Oxford, England), 34(5), 867–868. <https://doi.org/10.1093/bioinformatics/btx699>

Platten, J. D., Cobb, J. N., & Zantua, R. E. (2019). Criteria for evaluating molecular markers: Comprehensive quality metrics to improve marker-assisted selection. *PLoS one*, 14(1), e0210529. <https://doi.org/10.1371/journal.pone.0210529>

Poland, J. A., Brown, P. J., Sorrells, M. E., & Jannink, J. L. (2012). Development of high-density genetic maps for barley and wheat using a novel two-enzyme genotyping-by-sequencing approach. *PLoS one*, 7(2), e32253. <https://doi.org/10.1371/journal.pone.0032253>

Powder K. E. (2020). Quantitative Trait Loci (QTL) Mapping. Methods in molecular biology (Clifton, N.J.), 2082, 211–229. https://doi.org/10.1007/978-1-0716-0026-9_15

Qian, L., Hickey, L. T., Stahl, A., Werner, C. R., Hayes, B., Snowdon, R. J., & Voss-Fels, K. P. (2017). Exploring and Harnessing Haplotype Diversity to Improve Yield Stability in Crops. *Frontiers in plant science*, 8, 1534. <https://doi.org/10.3389/fpls.2017.01534>

R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

Rasheed, A., Mujeeb-Kazi, A., Ogbonnaya, F. C., He, Z., & Rajaram, S. (2018). Wheat genetic resources in the post-genomics era: promise and challenges. *Annals of botany*, 121(4), 603–616. <https://doi.org/10.1093/aob/mcx148>

Ren, S., Ahmed, N., Bertels, K., & Al-Ars, Z. (2019). GPU accelerated sequence alignment with traceback for GATK HaplotypeCaller. *BMC genomics*, 20(Suppl 2), 184. <https://doi.org/10.1186/s12864-019-5468-9>

Reynolds, M., Foulkes, J., Furbank, R., Griffiths, S., King, J., Murchie, E., Parry, M., & Slafer, G. (2012). Achieving yield gains in wheat. *Plant, cell & environment*, 35(10), 1799–1823. <https://doi.org/10.1111/j.1365-3040.2012.02588.x>

Reynolds, M., Foulkes, J., Furbank, R., Griffiths, S., King, J., Murchie, E., Parry, M., & Slafer, G. (2012). Achieving yield gains in wheat. *Plant, cell & environment*, 35(10), 1799–1823. <https://doi.org/10.1111/j.1365-3040.2012.02588.x>

Rizzi, R., Cairo, M., Makinen, V., Tomescu, A. I., & Valenzuela, D. (2019). Hardness of Covering Alignment: Phase Transition in Post-Sequence Genomics. *IEEE/ACM transactions on computational biology and bioinformatics*, 16(1), 23–30. <https://doi.org/10.1109/TCBB.2018.2831691>

Voss-Fels, K. P., Qian, L., Parra-Londono, S., Uptmoor, R., Frisch, M., Keeble-Gagnère, G., Appels, R., & Snowdon, R. J. (2017). Linkage drag constrains the roots of modern wheat. *Plant, cell & environment*, 40(5), 717–725. <https://doi.org/10.1111/pce.12888>

Voss-Fels, K. P., Snowdon, R. J., & Hickey, L. T. (2018). Designer Roots for Future Crops. *Trends in plant science*, 23(11), 957–960. <https://doi.org/10.1016/j.tplants.2018.08.004>

Voss-Fels, K., & Snowdon, R. J. (2016). Understanding and utilizing crop genome diversity via high-resolution genotyping. *Plant biotechnology journal*, 14(4), 1086–1094. <https://doi.org/10.1111/pbi.12456>

Walkowiak, S., Gao, L., Monat, C., Haberer, G., Kassa, M. T., Brinton, J., Ramirez-Gonzalez, R. H., Kolodziej, M. C., Delorean, E., Thambugala, D., Klymiuk, V., Byrns, B., Gundlach, H., Bandi, V., Siri, J. N., Nilsen, K., Aquino, C., Himmelbach, A., Copetti, D., Ban, T., ... Pozniak, C. J. (2020). Multiple wheat genomes reveal global variation in modern breeding. *Nature*, 588(7837), 277–283. <https://doi.org/10.1038/s41586-020-2961-x>

Wall, J. D., & Pritchard, J. K. (2003). Haplotype blocks and linkage disequilibrium in the human genome. *Nature reviews. Genetics*, 4(8), 587–597. <https://doi.org/10.1038/nrg1123>

Warschefsky, E., Penmetsa, R. V., Cook, D. R., & von Wettberg, E. J. (2014). Back to the wilds: tapping evolutionary adaptations for resilient crops through systematic hybridization with crop

wild relatives. *American journal of botany*, 101(10), 1791–1800.

<https://doi.org/10.3732/ajb.1400116>

Zhang, Z., Schwartz, S., Wagner, L., & Miller, W. (2000). A greedy algorithm for aligning DNA sequences. *Journal of computational biology : a journal of computational molecular cell biology*, 7(1-2), 203–214. <https://doi.org/10.1089/10665270050081478>

Zheng, Y., Roberts, R. J., & Kasif, S. (2004). Identification of genes with fast-evolving regions in microbial genomes. *Nucleic acids research*, 32(21), 6347–6357.

<https://doi.org/10.1093/nar/gkh935>

6. APPENDIX

The following pages contain the supplementary files 1 and 2

Haplotype-based Pangenome Analysis in Wheat

As the writer of this script I would like to express appreciation and give all credit to the authors of the following article for creating original raw data and scripts required for this R script:

Brinton, J., Ramirez-Gonzalez, R. H., Simmonds, J., Wingen, L., Orford, S., Griffiths, S., 10 Wheat Genome Project, Haberer, G., Spannagl, M., Walkowiak, S., Pozniak, C., & Uauy, C. (2020). A haplotype-led approach to increase the precision of wheat breeding. *Communications biology*, 3(1), 712. <https://doi.org/10.1038/s42003-020-01413-2> (<https://doi.org/10.1038/s42003-020-01413-2>)

This script was written in R 4.1.0 by Jose Antonio Montero Tena as part of his master thesis in 2021. Free use and modification of this script for personal interest is encouraged.

Introduction

SNP-haplotypes discovered from marker arrays can be redefined as sequence-based haplotypes, which are identical-by-state sequences sharing 100 percentage of identity between the genomes and discovered by sequence comparison. Brinton et al developed a method in 2020 that mainly defined haploblocks as physical regions with ≥ 99.99 pident across fixed-size bins of 5-, 2.5- or 1-Mbp between NUCmer pairwise comparisons of the assemblies in the 10+ Wheat Genome Project. This method will subsequently be referred to as the Brinton's method. The resulting sequence-based haplotypes can be either shared haplotypes or unique haplotypes, e.g. alleles that are carried by one single assembly. Additionally, Brinton and colleagues published [crop-haplotypes.com](#), a website that provides an interactive graphic visualization of the shared haplotypes between the wheat genome assemblies. These assemblies were chromosome-level genome assemblies of nine wheat assemblies (ArinaLrFor, Jagger, Julius, Lancer, Landmark, Mace, Norin61, Stanley, SY-Mattis) and the Chinese Spring RefSev1.0 assembly alongside scaffold-level assemblies of five additional assemblies (Cadenza, Claire, Paragon, Robigus and Weebil).

Aim

The aim of this script is to use Brinton's method and resources to conduct both a chromosome-scale and small-scale analysis of the pairwise comparisons between assemblies. The small-scale analysis option is original to this script and allows the analysis of target regions within chromosomes. The aim of the small-scale analysis is to assist with mapping physical start and end coordinates of the predicted haploblock both in the reference and query genomes after considering alignment coverage, in other words to avoid false positive haplotype predictions. Also, this script aims to identify the genes contained in mapped haploblock, both in the haplotype-carrying genomes and in the IWGSC Chinese Spring Ref v1.1. Brinton et al called haplotypes both from mummer and gene-based BLAST pairwise alignments and selected for matching positions the longer predictions. Both types of analysis are provided in this script, although haplotype predictions obtained these two methods, NUCmer and BLAST, do not always match.

Background

The SNP-haplotype Hap-5B-RDMA-h2, associated with increased root dry mass (RDM), is taken as an example. This haplotype was first found by GWAS upon 9 SNP marker alleles located in the wheat chromosome 5B and matched with the assemblies of the assemblies LongReach Lancer and Paragon by BLAST. SNP-RDMA-h2 BLAST position in Lancer's chromosome 5B (655760000-656600000 bp) was observed in [crop-haplotypes.com](#) in order to determine if the SNP-haploblock region matched with the region of any sequence-based haploblock, discovered with Brinton's method, in this area. The following blocks were observed in the pairwise comparison Lancer-Paragon under each bin size (bin size - start - end - length): 5-Mbp - 660.000.000 - 665.000.000 - 5.000.000 // 2.5-Mbp - 657.500.000 - 662.500.000 - 5.000.000 // 1-Mbp - 656.000.000 - 663.000.000 - 7.000.000 (start and end according to Lancer's coordinate system as Lancer acts as the reference in this comparison). Data from the website suggested that there might be a match between the SNP-based haplotype and a sequence-based haplotype in Lancer and Paragon. Despite of this, Brinton's method must be applied on this redefined region to understand if the prediction made on [crop-haplotypes.com](#) is reliable. Additionally, other research questions as the exact start and end position of the sequence-based haploblock or the underlying genes can be answered in this script.

Running the script

The script is designed so that you only have to edit only the upcoming parameters. You can run the script by lines (recommended for first time) or simply pressing 'Source' in the upper right corner of this window. The running time is expected to be between 5 and 7 minutes. Please, be patient and only look for solutions if error messages interrupt the pipeline. If errors do not stop coming up, you can contact the script editor at the e-mail address jmonterotena@gmail.com (<mailto:jmonterotena@gmail.com>).

Defining parameters

On which chromosome are your haploblocks located? Copy the same format

```
chromosome <- "5B"
```

Which is the reference assembly in the pairwise comparison that you are using? Alignment coordinates will apply for this assembly's chromosome. Can NEVER be a scaffold-level assembly (Cadenza, Claire, Paragon, Robigus, Weebil). Follow the next naming guidelines: "arinalfor", "jagger", "julius", "lancer", "landmark", "mace", "norin61", "stanley", "sy_mattis", "chinese" (works with capital letters)

```
reference_assembly <- "Lancer"
```

Which is the query assembly? This genome is aligned against the reference chromosome. Can be a scaffold-level assembly. Follow the next naming guidelines: "cadenza", "claire", "paragon", "robigus", "weebil" (works with capital letters)

```
query_assembly <- "Paragon"
```

Specify a start of the region highlighted in the small-scale analysis to look for haploblocks after this position. ADVICE: use multiples of 5 so that the bins match with those on [crop-haplotypes](#)

```
zoom_start <- 655000000
```

Specify an end of the region highlighted in the small-scale analysis to look for haploblocks before this position. ADVICE: use multiples of 5 so that the bins match with those on crop-haplotypes

```
zoom_end <- 665000000
```

Do you have previous information on haploblock locations in this region (e.g. SNP-based haploblocks from SNP chips)? If you have evidence of a haplotype within your zoom region that you want to compare with the sequence-based haplotypes obtained in this script, write its start coordinate in the reference chromosome. If not, simply write NA. The target region must be within the zoom region

```
target_start <- 655700000
```

Do you have previous information on haploblock locations in this region (e.g. SNP-based haploblocks from SNP chips)? If you have evidence of a haplotype within your zoom region that you want to compare with the sequence-based haplotypes obtained in this script, write its end coordinate in the reference chromosome. If not, simply write NA. The target region must be within the zoom region

```
target_end <- 656600000
```

You do not need to worry about the next step yet. Complete section one and then come back to define the **new coordinates of your haploblocks**. Genes will be extracted from this position. If you have no interest in redefining your region, simply write 'target_start' or 'zoom_start', to keep with the previous coordinates

```
selected_start <- 655760000
```

You do not need to worry about the next step yet. Complete section one and then come back to define the **new coordinates of your haploblocks**. Genes will be extracted until this position. If you have no interest in redefining your region, simply write 'target_end' or 'zoom_end', to keep with the previous coordinates

```
selected_end <- 662740000
```

Running functions

Source the script 'functions.hbpa.r' to read the packages and functions required for this script.

1. Calling haplotypes from raw data containing >= 20-Kbp-long NUCmer pairwise alignments between Lancer and Paragon

Requirements:

- Raw delta files in 'all_20_kb_filtered_delta_CHROMOSOME_tables.rds'; in this case the chromosome is 5B (downloadable from https://opendata.earlham.ac.uk/wheat/under_license/toronto/Brinton_et.al_2020-05-20-Haplotypes-for-wheat-breeding/nucmer/ (https://opendata.earlham.ac.uk/wheat/under_license/toronto/Brinton_et.al_2020-05-20-Haplotypes-for-wheat-breeding/nucmer/)). Brinton et al filtered out alignments under 20000Kbp of length in order to get rid of non-syntenic retrotransposons
- Script 'functions.hbpa.r' (downloadable from <https://github.com/MonteroJLU/Haplotype-based-Pangenome-Wheat-Analysis.git> (<https://github.com/MonteroJLU/Haplotype-based-Pangenome-Wheat-Analysis.git>))
- Text file 'table_assembly_chr_length.txt' (downloadable from <https://github.com/MonteroJLU/Haplotype-based-Pangenome-Wheat-Analysis.git> (<https://github.com/MonteroJLU/Haplotype-based-Pangenome-Wheat-Analysis.git>))

1.1. Chromosome-scale analysis

1.1.1. Download and save the required documents in the working directory

1.1.2. Run the script 'functions_hbpa.r'

1.1.3. Read the rds file into a data frame:

```
aln_library <- readRDS(file = paste0("all_20_kb_filtered_delta_", chromosome, "_tables.rds"))
```

1.1.4. Create a subset of the data frame containing only the alignments from the pairwise comparison Lancer-Paragon, where Lancer acts as the the reference genome:

```
aln_subset <- aln_library[grepl(paste0("^", tolower(reference_assembly), sep = ""), aln_library$comparison) & grepl(tolower(query_assembly), aln_library$comparison),]  
head(aln_subset, 3)
```

```

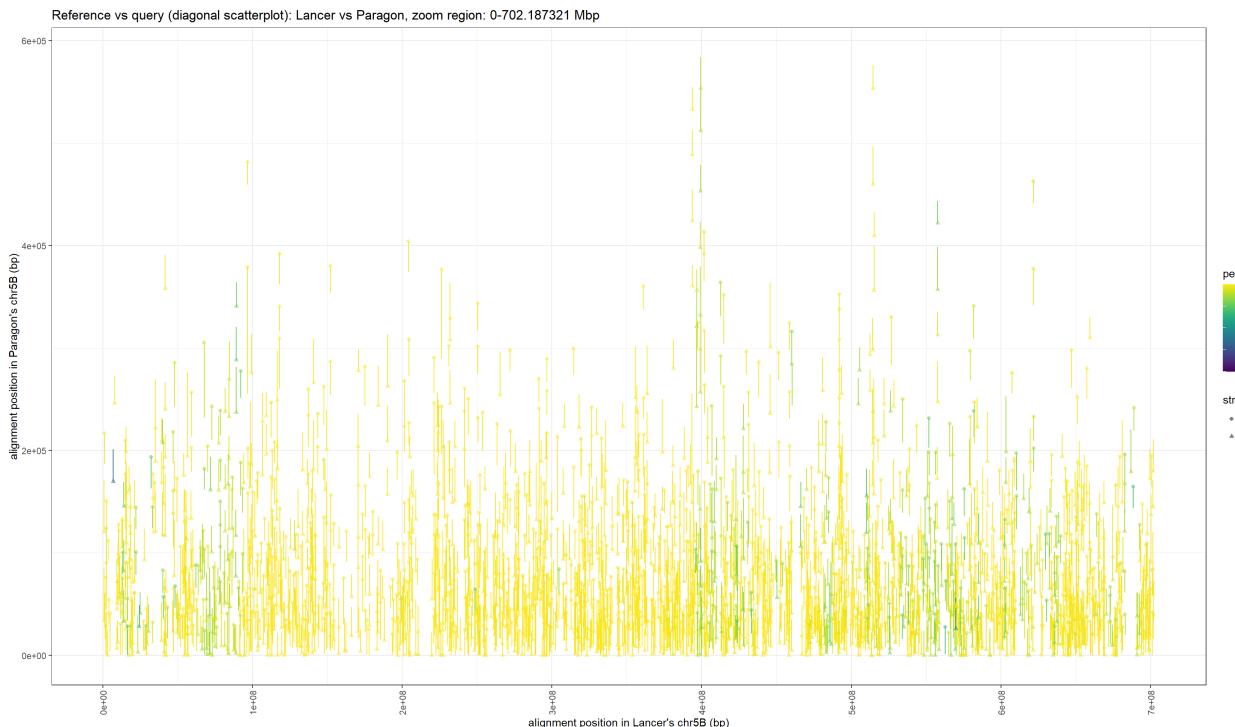
##          rs      re      qs      qe error          qid      rid
## 110088 398797814 398831949    527  34667    454 par_scaffold_000027 chr5B_lac
## 210074 398846085 398881400   48451  83752    257 par_scaffold_000027 chr5B_lac
## 310044 398885774 398923047 104871 142145   268 par_scaffold_000027 chr5B_lac
##          strand r_length perc_id perc_id_factor   r_mid   q_mid
## 110088     +    34135 98.66999      <100 398814882 17597.0
## 210074     +    35315 99.27226      <100 398863743 66101.5
## 310044     +    37273 99.28098      <100 398904411 123508.0
##          comparison chrom
## 110088 lancer_v_paragon.5B_filtered_L20Kb_rq.delta    5B
## 210074 lancer_v_paragon.5B_filtered_L20Kb_rq.delta    5B
## 310044 lancer_v_paragon.5B_filtered_L20Kb_rq.delta    5B

```

Description of the column headers: rs: reference start, re: reference end, qs: query start, qe: query end, error: number of unmatches, qid: query identification, rid: reference identification, strand: forward or reverse strand, r_length: length of the alignment in the reference, perc_id: percentage of identity, perc_id_factor: factor that summarises perc_id, r_mid: midpoint the alignment in reference ((r_end-r_start)/2), q_mid: query midpoint, comparison: assemblies compared, chrom: chromosome.

1.1.5. Scatter-plot the alignment midpoints (X: r_mid, Y: q_mid)

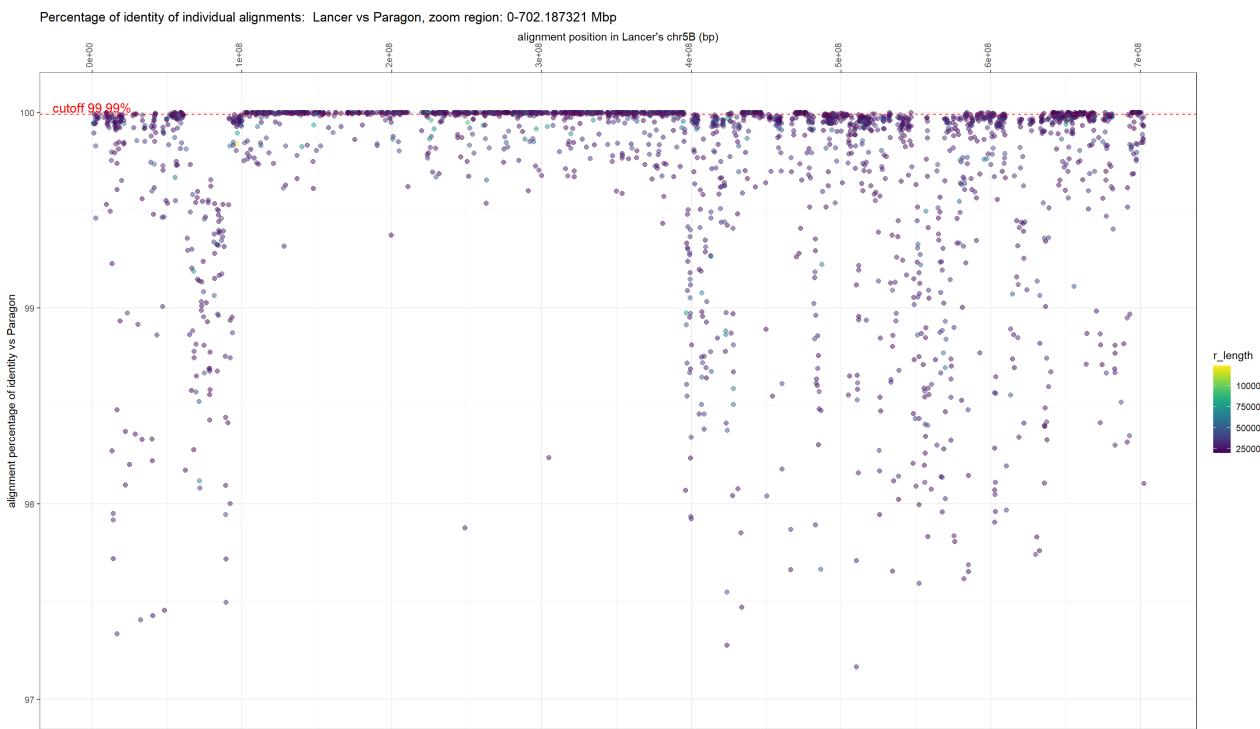
```
plot_diagonal_scatterplot(aln_subset, cap_lower = 90.00, cap_upper = 100, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 100000000)
```



Notice unexpected vertical lines in the scatter-plot due to the use of a scaffold-level assembly as query (CLA-SLA comparison). However, this function works well for CLA-CLA comparisons.

1.1.6. Dot-plot the alignments to show percentage of identity and alignment length (X: r_mid, Y: perc_id)

```
plot_aln_pid_and_length(aln_subset, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 1000000
00, dot_size = 2)
```



1.1.7. Check for alignment properties in the data frame (percentage of alignment coverage of reference chromosome, average aln length and av aln number)

Chromosome length

```
table_chrlength_v_assembly <- read.table(file = "table_assembly_chr_length.txt", sep = "\t", header = TRUE)
chr_length <- table_chrlength_v_assembly$sequence_length[grep(tolower(reference_assembly), table_chrlength_v_assembly$assembly_name) & grep(paste0("chr", unique(aln_subset$chrom)), table_chrlength_v_assembly$sequence_name)]
print(paste0("chr", unique(aln_subset$chrom), " is ", chr_length, " bp-long in ", reference_assembly))
```

```
## [1] "chr5B is 702438406 bp-long in Lancer"
```

Average alignment length in this pairwise comparison

```
print(paste0(round(mean(aln_subset$r_length), 0), " is the average alignment length in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_subset$chrom), " comparison")))
```

```
## [1] "30858 is the average alignment length in Lancer-Paragon chr5B comparison"
```

Number of alignments in this pairwise comparison

```
print(paste0(nrow(aln_subset), " is the number of alignments in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_subset$chrom), " comparison")))
```

```
## [1] "2650 is the number of alignments in Lancer-Paragon chr5B comparison"
```

Alignment coverage in this pairwise comparison

```
print(paste0(round((sum(aln_subset$r_length)/chr_length*100), 0), "% is the alignment coverage in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_subset$chrom), " comparison")))
```

```
## [1] "12% is the alignment coverage in Lancer-Paragon chr5B comparison"
```

Alignment coverage overview in different comparisons by reference assembly and comparison type (with/without scaffold-level assemblies)

```
coverage <- COV(aln_library, chr_length = table_chrlength_v_assembly)
print(coverage)
```

```

## $`Average % coverage with CHROMOSOME-LEVEL ASSEMBLIES as query`
## [1] 47.74529
##
## $`By assembly`'
## arinalrfor    jagger    julius    lancer    landmark      mace    norin61
## 28.44218  50.59484  51.33529  55.81078  50.99109  55.28362  53.99118
## stanley sy_mattis   chinese
## 51.89199  28.44447  50.66749
##
## $`Average % coverage with SCAFFOLD-LEVEL ASSEMBLIES as query`
## [1] 9.390724
##
## $`By assembly`'
## arinalrfor    jagger    julius    lancer    landmark      mace    norin61
## 4.216677 10.113839 10.478023 12.165450  9.566494 11.940812 11.252819
## stanley sy_mattis   chinese
## 10.103706 4.302350  9.767073

```

Average alignment length overview in different comparisons by reference assembly and comparison type (with/without scaffold-level assemblies)

```

length <- LENGTH(aln_library)
print(length)

```

```

## $`Average aln length with CHROMOSOME-LEVEL ASSEMBLIES as query`
## [1] 45635.03
##
## $`By assembly`'
## arinalrfor    jagger    julius    lancer    landmark      mace    norin61
## 50573.98  42245.55  43894.80  46985.47  42621.59  46726.19  45787.59
## stanley sy_mattis   chinese
## 44081.69  50841.05  42592.41
##
## $`Average aln length with SCAFFOLD-LEVEL ASSEMBLIES as query`
## [1] 30425.71
##
## $`By assembly`'
## arinalrfor    jagger    julius    lancer    landmark      mace    norin61
## 31092.67  29791.85  30193.23  31045.12  29593.04  30915.88  30557.44
## stanley sy_mattis   chinese
## 29983.10  31462.44  29622.29

```

Average number of alignments overview in different comparisons by reference assembly and comparison type (with/without scaffold-level assemblies)

```

number <- NUMBER(aln_library)
print(number)

```

```

## $`Average aln number with CHROMOSOME-LEVEL ASSEMBLIES as query`
## [1] 7253.244
##
## $`By assembly`'
## arinalrfor    jagger    julius    lancer    landmark      mace    norin61
## 2703.778  8426.556  8480.333  8343.778  8500.111  8172.778  8399.667
## stanley sy_mattis   chinese
## 8404.222  2617.667  8483.556
##
## $`Average aln number with SCAFFOLD-LEVEL ASSEMBLIES as query`
## [1] 2129.46
##
## $`By assembly`'
## arinalrfor    jagger    julius    lancer    landmark      mace    norin61
## 652.0     2388.6    2516.4    2752.6    2296.8    2668.0    2623.2
## stanley sy_mattis   chinese
## 2405.8    639.8    2351.4

```

Expected number of alignments per bin in this pairwise comparison at different bin sizes

```

bin_size <- c(5000000, 2500000, 1000000)
names(bin_size) <- c("bin size: 5-Mbp", "bin size: 2.5-Mbp", "bin size: 1-Mbp")
for (i in 1:3){
  print("average expected number of alignments per bin across the chromosome")
  print(round(nrow(aln_subset)/(chr_length/bin_size[i]), 1))
}

```

```

## [1] "average expected number of alignments per bin across the chromosome"
## bin size: 5-Mbp
##          18.9
## [1] "average expected number of alignments per bin across the chromosome"
## bin size: 2.5-Mbp
##          9.4
## [1] "average expected number of alignments per bin across the chromosome"
## bin size: 1-Mbp
##          3.8

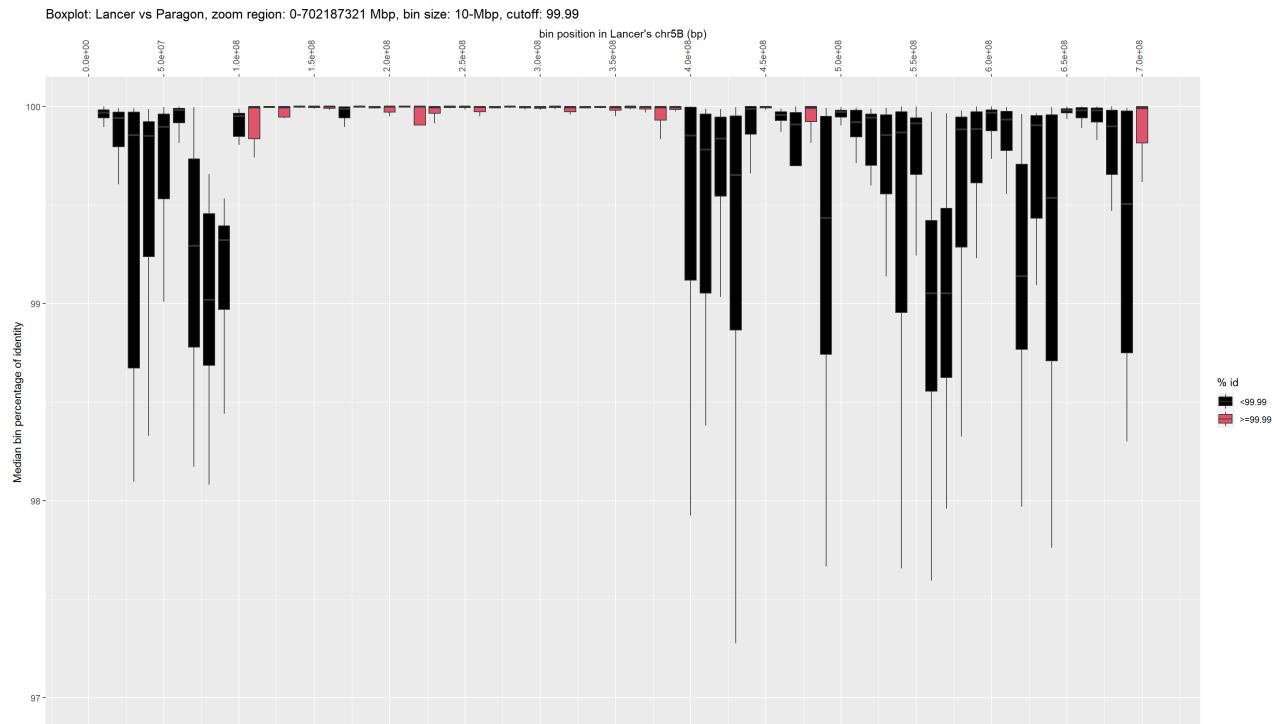
```

1.1.8. Boxplot the alignments (X: bin, Y: perc_id_median)

```

plot_boxplots_bin_median(aln_subset, bin_size = 10000000, bin_start = 0, bin_end = max(aln_subset$re), cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 50000000, show_outliers = FALSE)

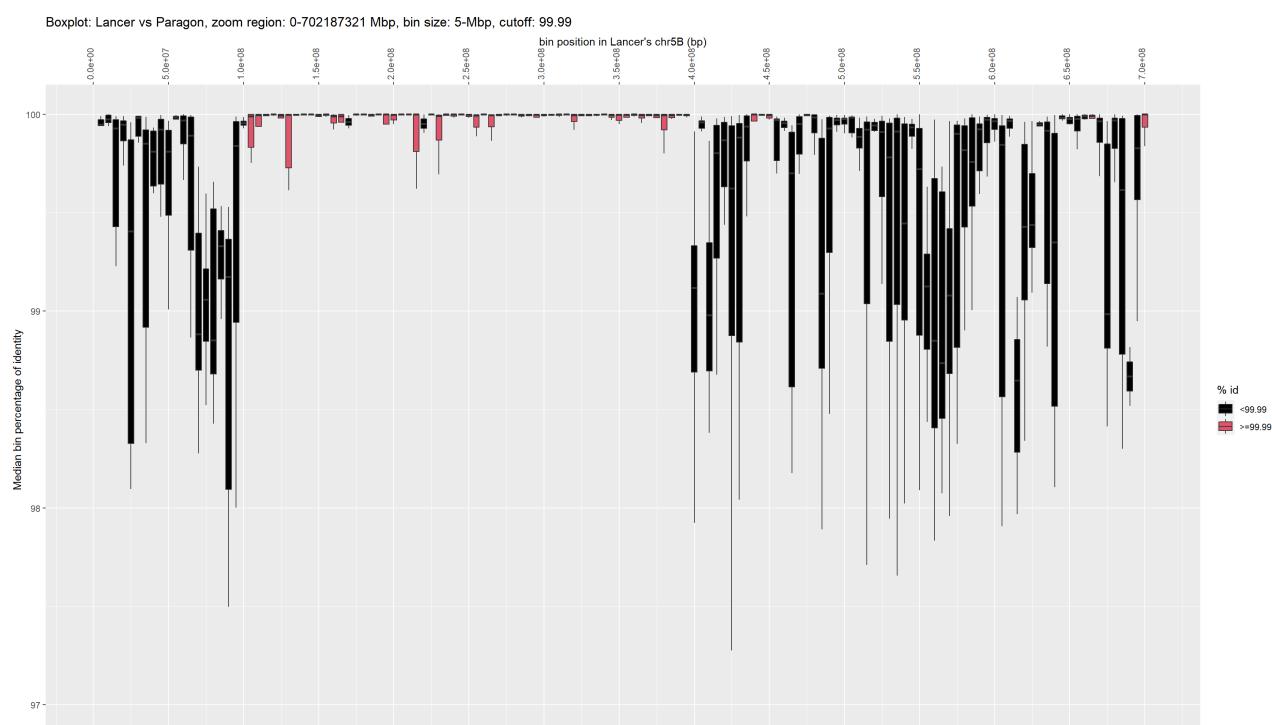
```



```

plot_boxplots_bin_median(aln_subset, bin_size = 5000000, bin_start = 0, bin_end = max(aln_subset$re), cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 50000000, show_outliers = FALSE)

```

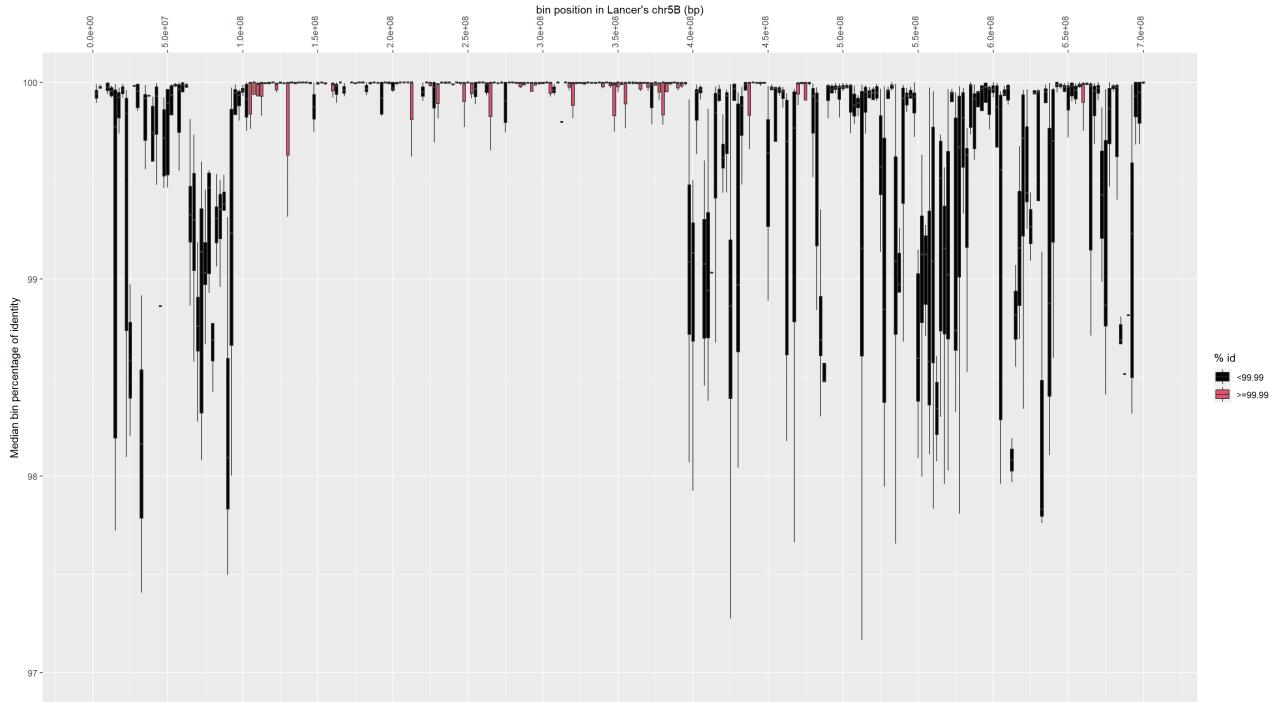


```

plot_boxplots_bin_median(aln_subset, bin_size = 2500000, bin_start = 0, bin_end = max(aln_subset$re), cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 50000000, show_outliers = FALSE)

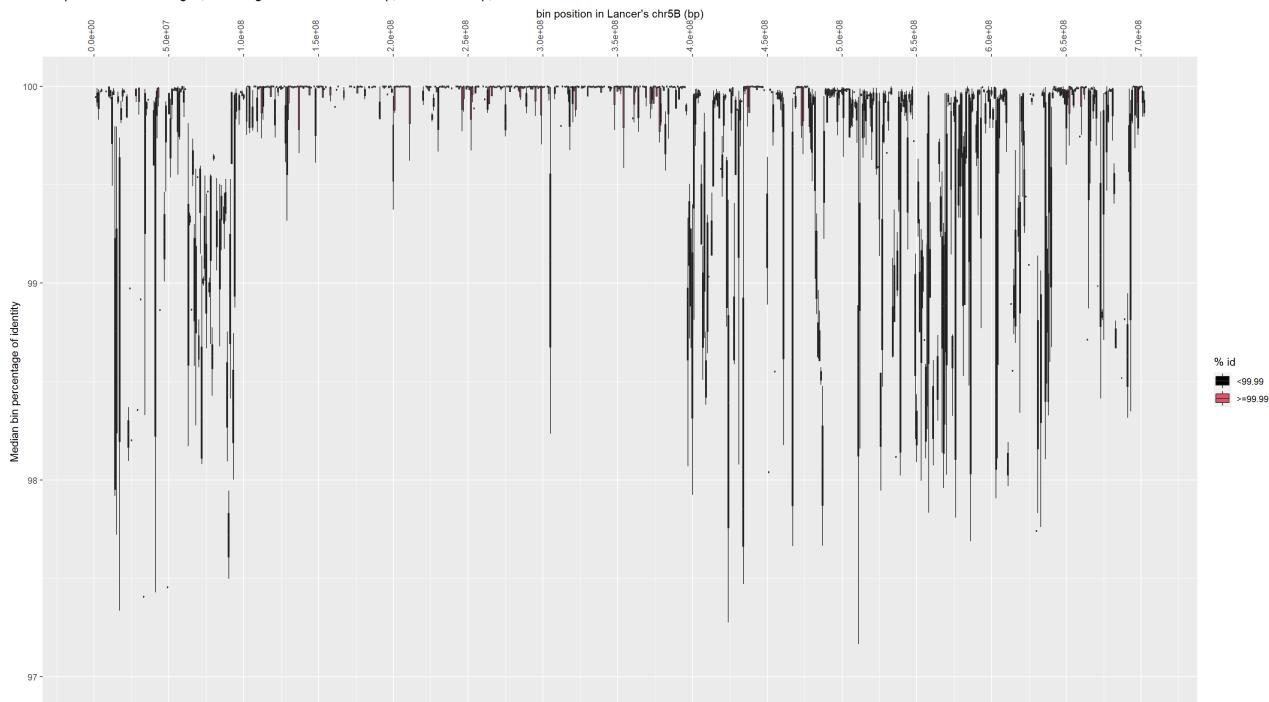
```

Boxplot: Lancer vs Paragon, zoom region: 0-702187321 Mbp, bin size: 2.5-Mbp, cutoff: 99.99



```
plot_boxplots_bin_median(aln_subset, bin_size = 1000000, bin_start = 0, bin_end = max(aln_subset$re), cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 5000000, show_outliers = FALSE)
```

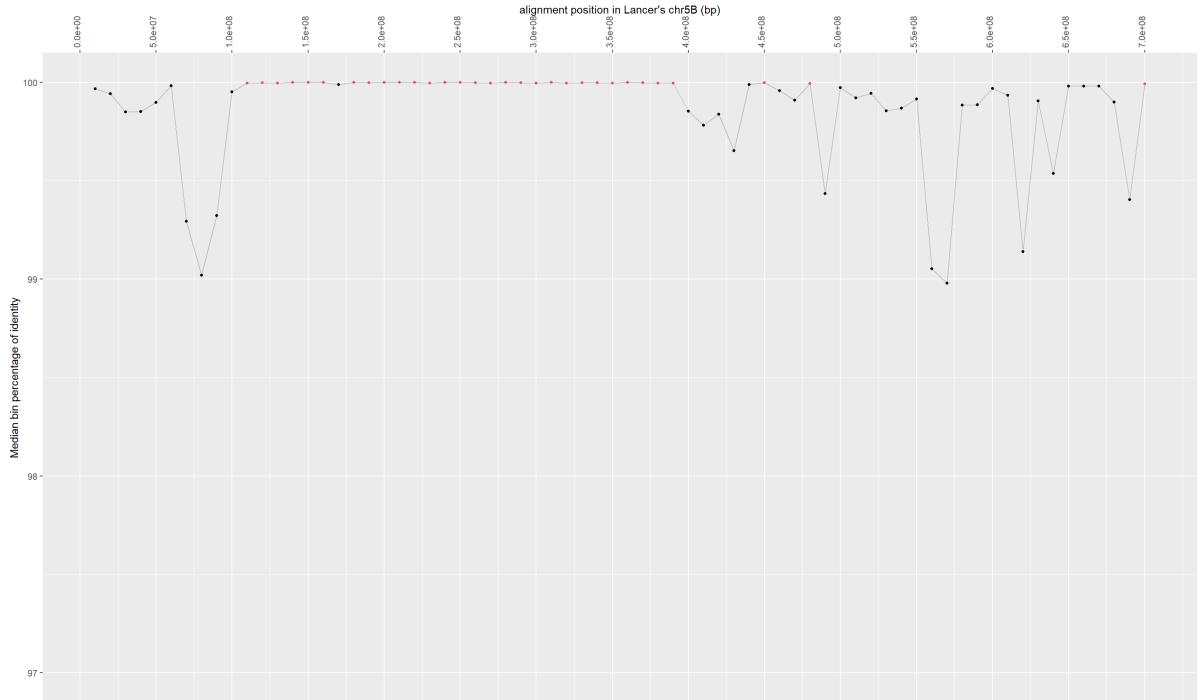
Boxplot: Lancer vs Paragon, zoom region: 0-702187321 Mbp, bin size: 1-Mbp, cutoff: 99.99



1.1.9. Plot the median percentage of identity across the chromosome or median line (X: r_mid, Y: perc_id_median)

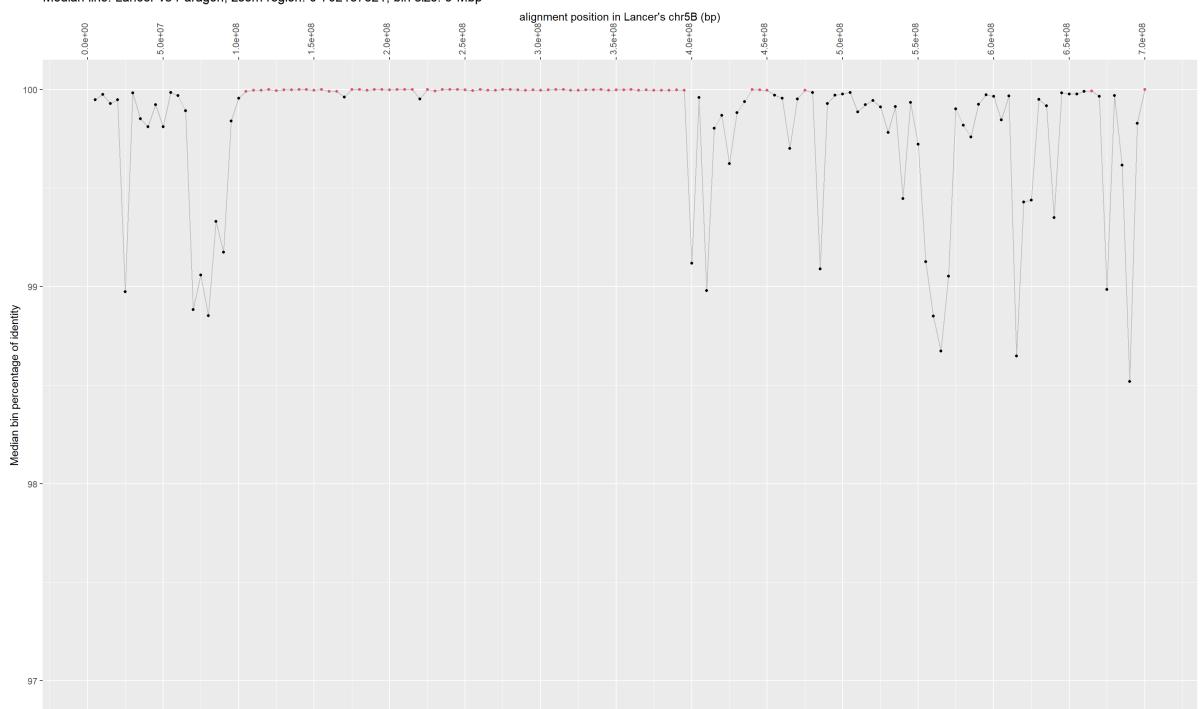
```
medians_aln_subset_10Mbp <- plot_line_bin_median(aln_subset, bin_size = 1000000, bin_start = 0, bin_end = max(aln_subset$re), cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 5000000)
```

Median line: Lancer vs Paragon, zoom region: 0-702187321, bin size: 10-Mbp



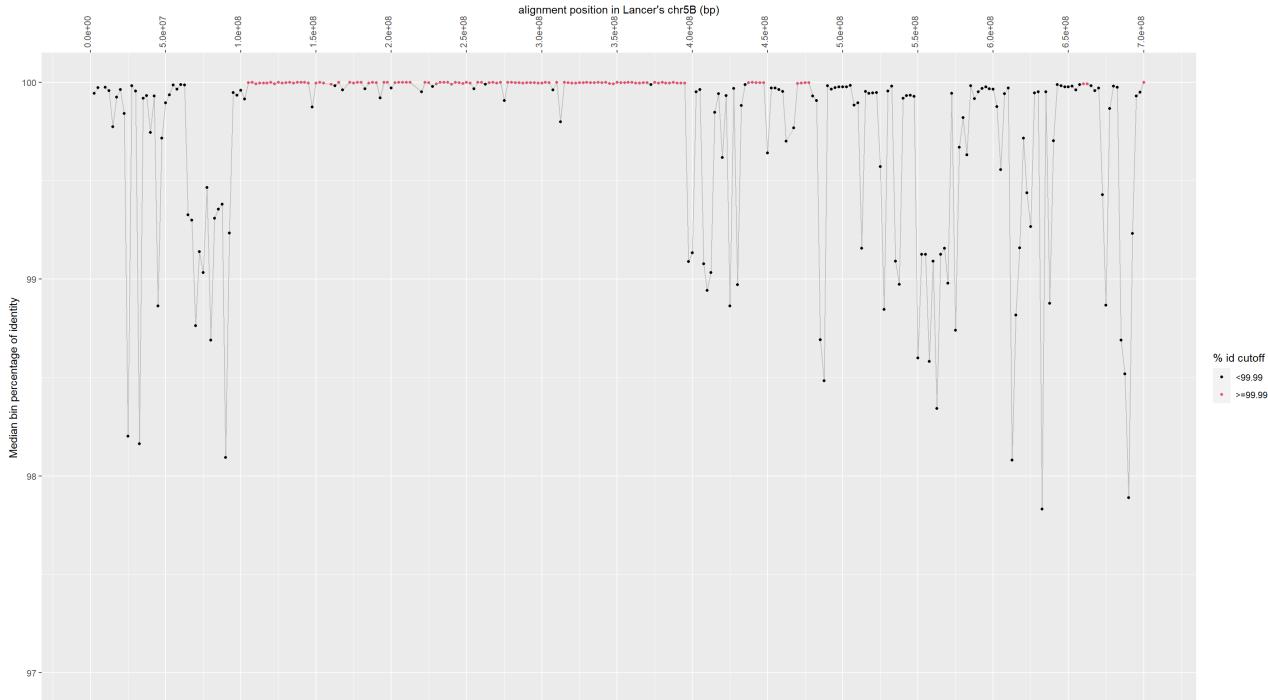
```
medians_aln_subset_5Mbp <- plot_line_bin_median(aln_subset, bin_size = 5000000, bin_start = 0, bin_end = max(aln_subset$re),  
cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 50000000)  
medians_aln_subset_5Mbp
```

Median line: Lancer vs Paragon, zoom region: 0-702187321, bin size: 5-Mbp



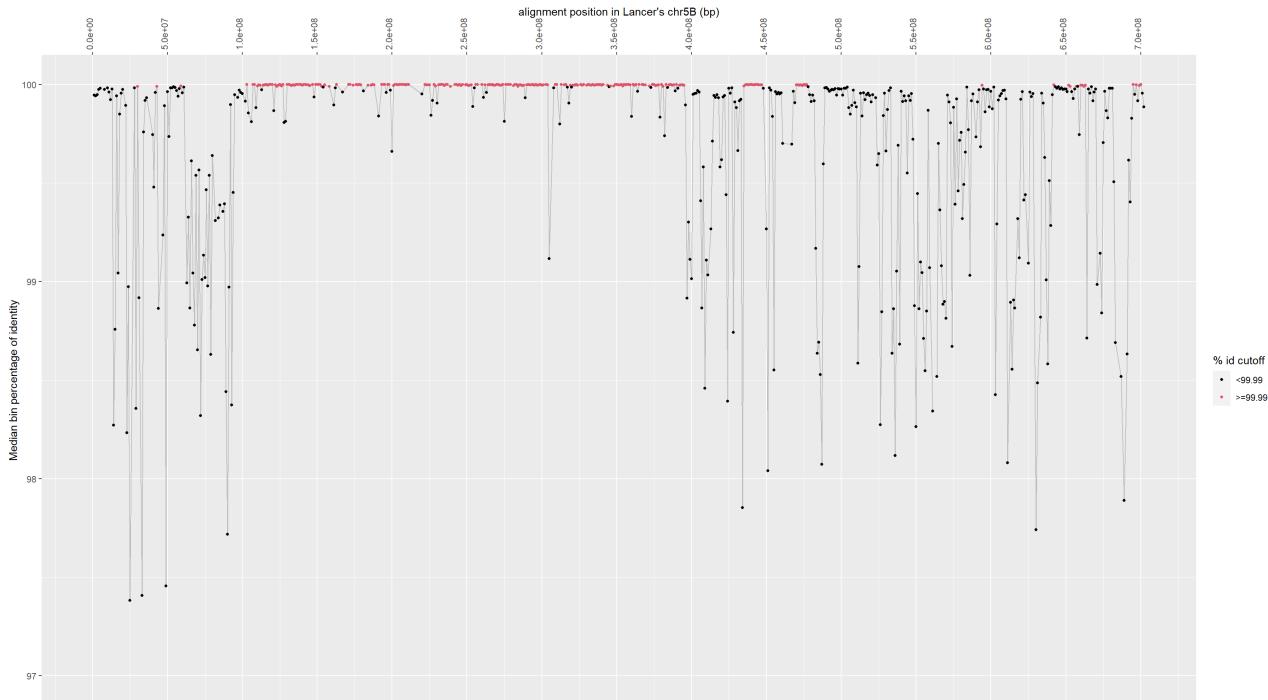
```
medians_aln_subset_2.5Mbp <- plot_line_bin_median(aln_subset, bin_size = 2500000, bin_start = 0, bin_end = max(aln_subset$re),  
cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 50000000)  
medians_aln_subset_2.5Mbp
```

Median line: Lancer vs Paragon, zoom region: 0-702187321, bin size: 2.5-Mbp



```
medians_aln_subset_1Mbp <- plot_line_bin_median(aln_subset, bin_size = 1000000, bin_start = 0, bin_end = max(aln_subset$re),
cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 50000000)
medians_aln_subset_1Mbp
```

Median line: Lancer vs Paragon, zoom region: 0-702187321, bin size: 1-Mbp



1.1.10. Print information about the haploblock predictions in the pairwise comparison across the reference chromosome

```
medians_aln_subset_10Mbp_bin_info <- assign_blocks_mummer(median_cutoffs = medians_aln_subset_10Mbp[["data"]], original_file =
= aln_subset)
medians_aln_subset_5Mbp_bin_info <- assign_blocks_mummer(median_cutoffs = medians_aln_subset_5Mbp[["data"]], original_file =
= aln_subset)
medians_aln_subset_2.5Mbp_bin_info <- assign_blocks_mummer(median_cutoffs = medians_aln_subset_2.5Mbp[["data"]], original_file =
= aln_subset)
medians_aln_subset_1Mbp_bin_info <- assign_blocks_mummer(median_cutoffs = medians_aln_subset_1Mbp[["data"]], original_file =
= aln_subset)
medians_aln_subset_10Mbp_block_summary <- block_summary(medians_aln_subset_10Mbp_bin_info, bin_size = 10000000, reference_na
me = reference_assembly, query_name = query_assembly )
medians_aln_subset_5Mbp_block_summary <- block_summary(medians_aln_subset_5Mbp_bin_info, bin_size = 5000000, reference_name =
= reference_assembly, query_name = query_assembly )
medians_aln_subset_2.5Mbp_block_summary <- block_summary(medians_aln_subset_2.5Mbp_bin_info, bin_size = 2500000, reference_n
ame = reference_assembly, query_name = query_assembly )
medians_aln_subset_1Mbp_block_summary <- block_summary(medians_aln_subset_1Mbp_bin_info, bin_size = 1000000, reference_name =
= reference_assembly, query_name = query_assembly )
print(medians_aln_subset_10Mbp_block_summary)
```

```
## bin_size      comparison block_no block_start block_end
## 1 10-Mbp Lancer->Paragon      1   1.0e+08  3.9e+08
## 2 10-Mbp Lancer->Paragon      2   4.4e+08  4.8e+08
## 3 10-Mbp Lancer->Paragon      3   6.9e+08  7.0e+08
```

```
print(medians_aln_subset_5Mbp_block_summary)
```

```
## bin_size      comparison block_no block_start block_end
## 1 5-Mbp Lancer->Paragon      1   1.00e+08  3.95e+08
## 2 5-Mbp Lancer->Paragon      2   4.35e+08  4.50e+08
## 3 5-Mbp Lancer->Paragon      3   4.70e+08  4.75e+08
## 4 5-Mbp Lancer->Paragon      4   6.60e+08  6.65e+08
## 5 5-Mbp Lancer->Paragon      5   6.95e+08  7.00e+08
```

```
print(medians_aln_subset_2.5Mbp_block_summary)
```

```
## bin_size      comparison block_no block_start block_end
## 1 2.5-Mbp Lancer->Paragon      1  102500000 395000000
## 2 2.5-Mbp Lancer->Paragon      2  435000000 447500000
## 3 2.5-Mbp Lancer->Paragon      3  467500000 477500000
## 4 2.5-Mbp Lancer->Paragon      4  657500000 662500000
## 5 2.5-Mbp Lancer->Paragon      5  697500000 700000000
```

```
print(medians_aln_subset_1Mbp_block_summary)
```

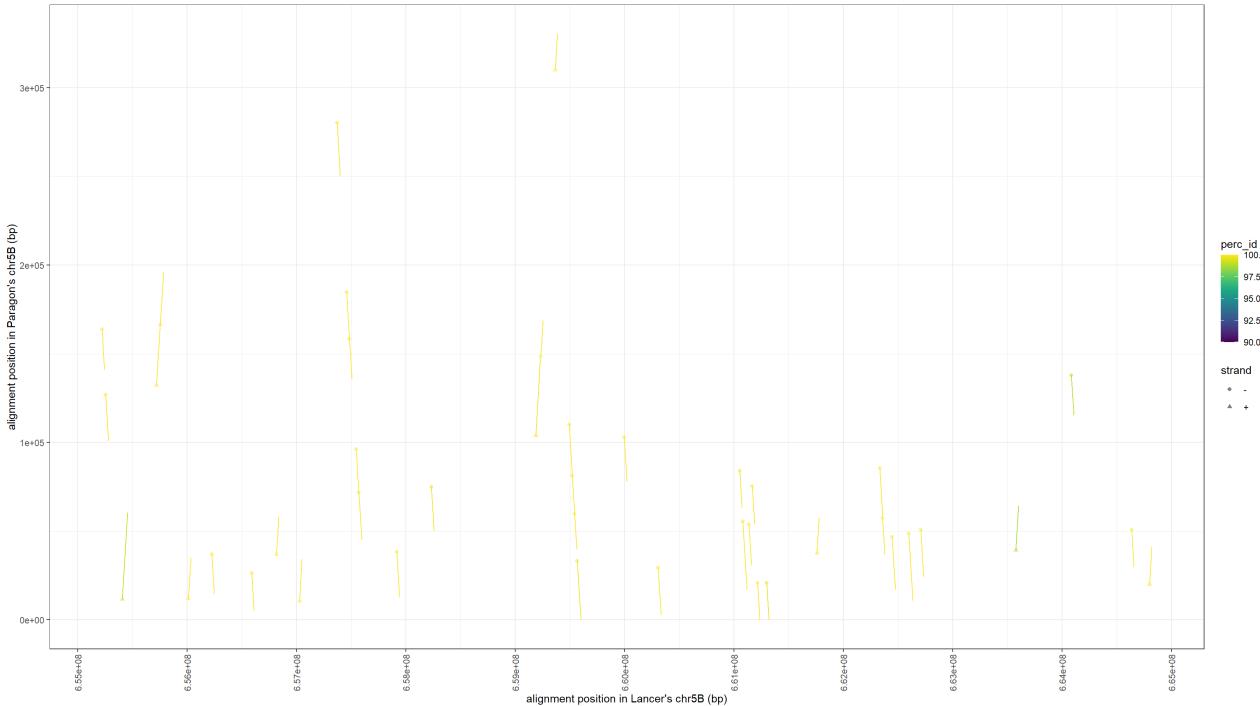
```
## bin_size      comparison block_no block_start block_end
## 1 1-Mbp Lancer->Paragon      1   2.90e+07  3.00e+07
## 2 1-Mbp Lancer->Paragon      2   4.20e+07  4.30e+07
## 3 1-Mbp Lancer->Paragon      3   5.80e+07  5.90e+07
## 4 1-Mbp Lancer->Paragon      4   1.02e+08  3.95e+08
## 5 1-Mbp Lancer->Paragon      5   4.34e+08  4.47e+08
## 6 1-Mbp Lancer->Paragon      6   4.69e+08  4.77e+08
## 7 1-Mbp Lancer->Paragon      7   5.93e+08  5.94e+08
## 8 1-Mbp Lancer->Paragon      8   6.41e+08  6.42e+08
## 9 1-Mbp Lancer->Paragon      9   6.51e+08  6.53e+08
## 10 1-Mbp Lancer->Paragon     10  6.56e+08  6.63e+08
## 11 1-Mbp Lancer->Paragon     11  6.94e+08  7.00e+08
```

Block number 10 reveals the haplotype prediction obtained at crop-haplotypes.com at 1-Mbp bin size between Lancer and Paragon that could match with the SNP-based haplotype discovered previously. This information can be contrasted with other haplotypes shared between Lancer and Paragon graphically with different graphs (boxplot, median line...) back in this script.

1.2. Small-scale analysis

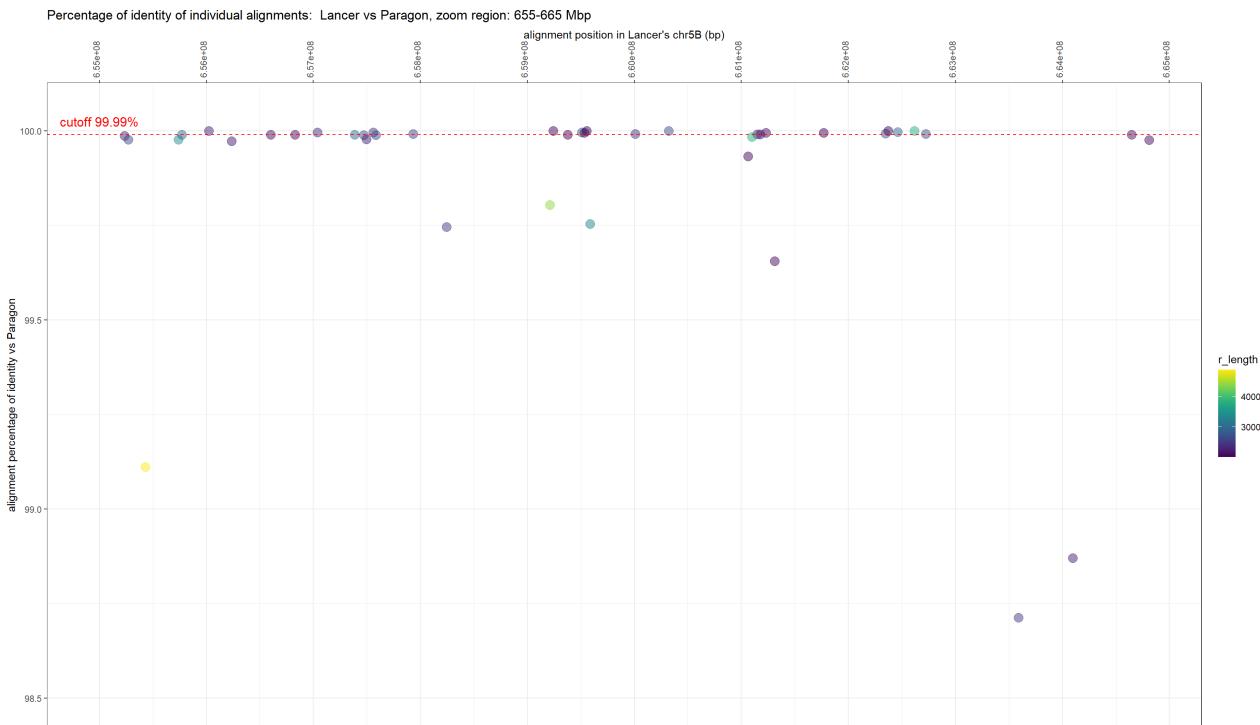
1.2.1. Scatter-plot the alignment midpoints across the zoom region (X: r_mid, Y: q_mid)

```
plot_diagonal_scatterplot(aln_subset, xmin = zoom_start, xmax = zoom_end, cap_lower = 90.00, cap_upper = 100, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 1000000)
```



1.2.2. Dot-plot the alignments to show percentage of identity and alignment length in the zoom region (X: r_mid, Y: perc_id)

```
graph <- plot_aln_pid_and_length(data = aln_subset, xmin = zoom_start, xmax = zoom_end, ymin = 98.5, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 1000000, dot_size = 4)
graph
```



1.2.3. Check for alignment properties in the zoom region

Average alignment length in the target region

```
aln_target <- aln_subset[(aln_subset$r_mid >= target_start) & (aln_subset$r_mid <= target_end),]
print(paste0(round(mean(aln_target$r_length), 0), " is the average alignment length for the target region between ", target_start/1e06, " and ", target_end/1e06, " Mbp in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_target$chrom)), " comparison"))
```

```
## [1] "27206 is the average alignment length for the target region between 655.7 and 656.6 Mbp in Lancer-Paragon chr5B comparison"
```

Number of alignments in the target region

```
print(paste0(nrow(aln_target), " is the number of alignments for the target region between ", target_start/1e06, " and ", target_end/1e06, " Mbp in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_target$chrom)), " comparison"))
})
```

```
## [1] "4 is the number of alignments for the target region between 655.7 and 656.6 Mbp in Lancer-Paragon chr5B comparison"
```

Alignment coverage in the target region

```
print(paste0(round((sum(aln_target$r_length)/(target_end-target_start)*100), 0), "% is the alignment coverage for the target region between ", target_start/1e06, " and ", target_end/1e06, " Mbp in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_target$chrom)), " comparison"))
```

```
## [1] "12% is the alignment coverage for the target region between 655.7 and 656.6 Mbp in Lancer-Paragon chr5B comparison"
```

Average alignment length in the zoom region

```
aln_zoom <- aln_subset[(aln_subset$r_mid >= zoom_start) & (aln_subset$r_mid <= zoom_end),]
print(paste0(round(mean(aln_zoom$r_length), 0), " is the average alignment length for the zoom region between ", zoom_start/1e06, " and ", zoom_end/1e06, " Mbp in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_subset$chrom)), " comparison"))
```

```
## [1] "25963 is the average alignment length for the zoom region between 655 and 665 Mbp in Lancer-Paragon chr5B comparison"
```

Number of alignments in the zoom region

```
print(paste0(nrow(aln_zoom), " is the number of alignments for the zoom region between ", zoom_start/1e06, " and ", zoom_end/1e06, " Mbp in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_subset$chrom)), " comparison"))
```

```
## [1] "42 is the number of alignments for the zoom region between 655 and 665 Mbp in Lancer-Paragon chr5B comparison"
```

Alignment coverage in the zoom region

```
print(paste0(round((sum(aln_zoom$r_length)/(zoom_end-zoom_start)*100), 0), "% is the alignment coverage for the zoom region between ", zoom_start/1e06, " and ", zoom_end/1e06, " Mbp in ", reference_assembly, "-", query_assembly, paste0(" chr", unique(aln_subset$chrom)), " comparison"))
```

```
## [1] "11% is the alignment coverage for the zoom region between 655 and 665 Mbp in Lancer-Paragon chr5B comparison"
```

Expected number of alignments per bin across the zoom region

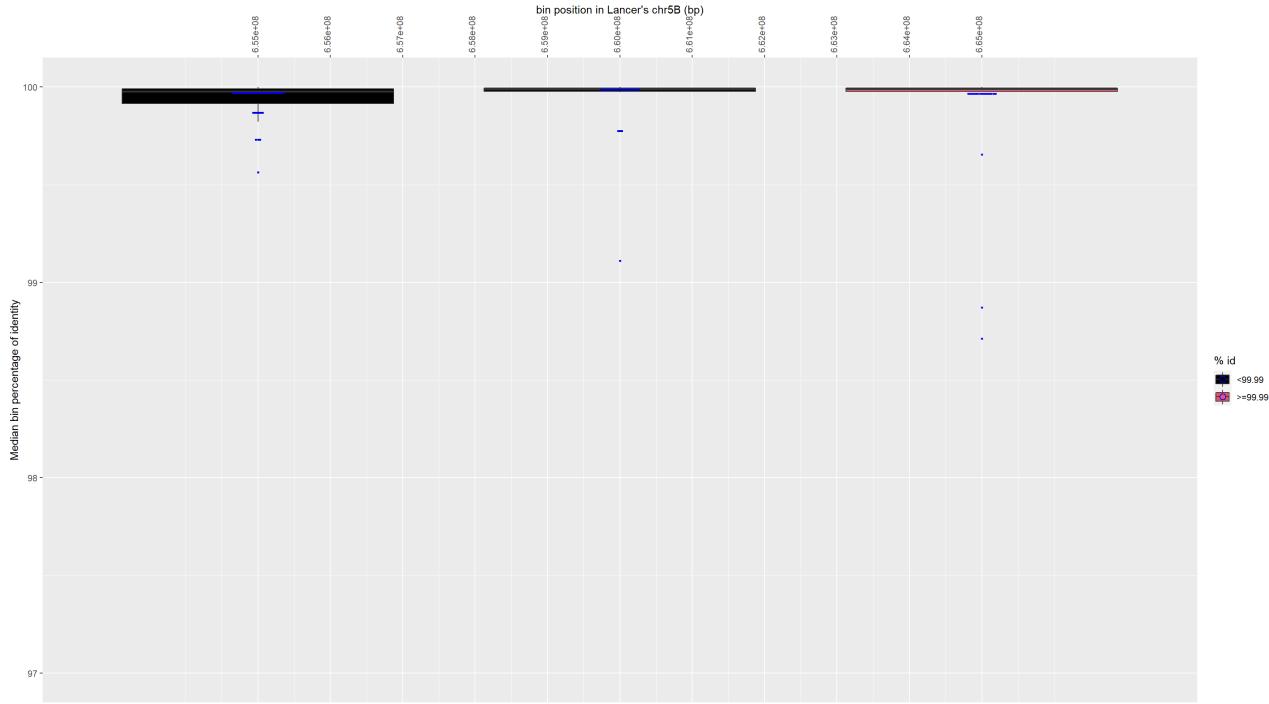
```
bin_size <- c(5000000, 2500000, 1000000)
names(bin_size) <- c("bin size: 5-Mbp", "bin size: 2.5-Mbp", "bin size: 1-Mbp")
for (i in 1:3){
  print("average expected number of alignments per bin across zoom region")
  print(nrow(aln_zoom)/((zoom_end-zoom_start)/bin_size[i]))
}
```

```
## [1] "average expected number of alignments per bin across zoom region"
## bin size: 5-Mbp
##          21
## [1] "average expected number of alignments per bin across zoom region"
## bin size: 2.5-Mbp
##          10.5
## [1] "average expected number of alignments per bin across zoom region"
## bin size: 1-Mbp
##          4.2
```

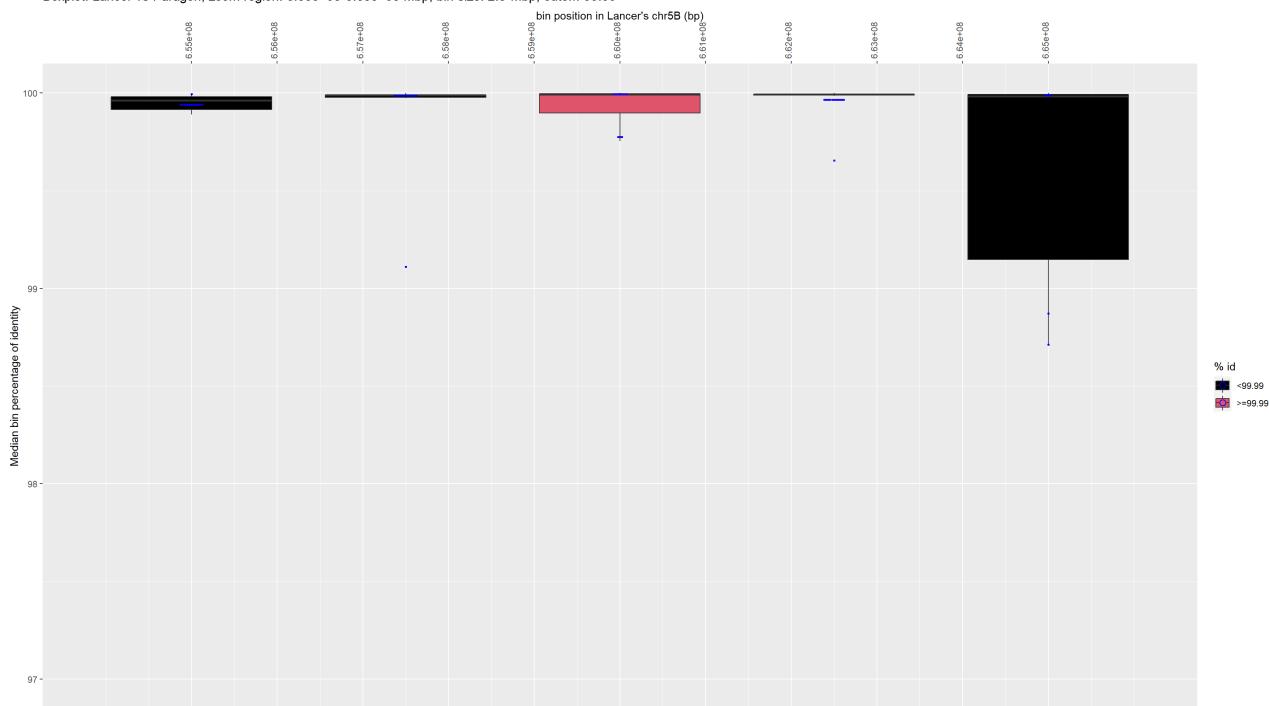
1.2.4. Boxplot the bin median percentage of identity in the zoom region to check for outliers

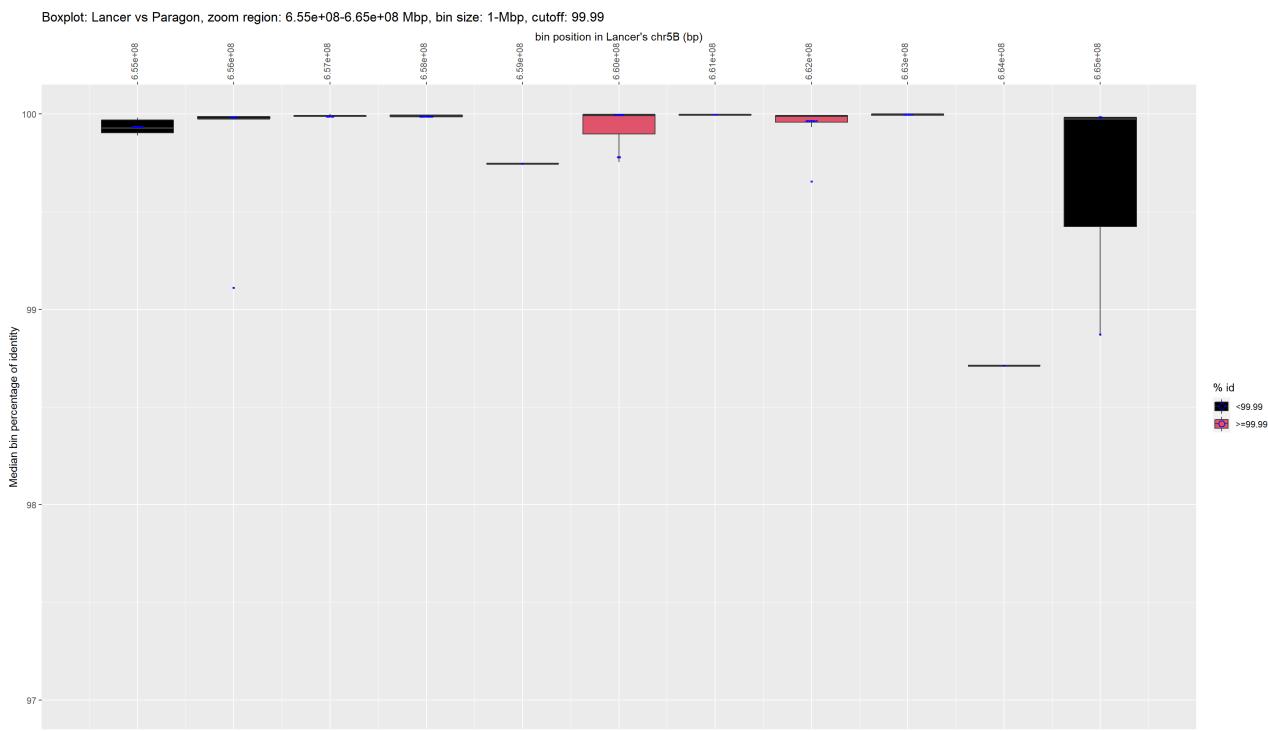
```
plot_boxplots_bin_median(aln_subset, bin_size = 5000000, bin_start = zoom_start, bin_end = zoom_end, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 1000000, show_outliers = TRUE)
```

Boxplot: Lancer vs Paragon, zoom region: 6.55e+08-6.65e+08 Mbp, bin size: 5-Mbp, cutoff: 99.99



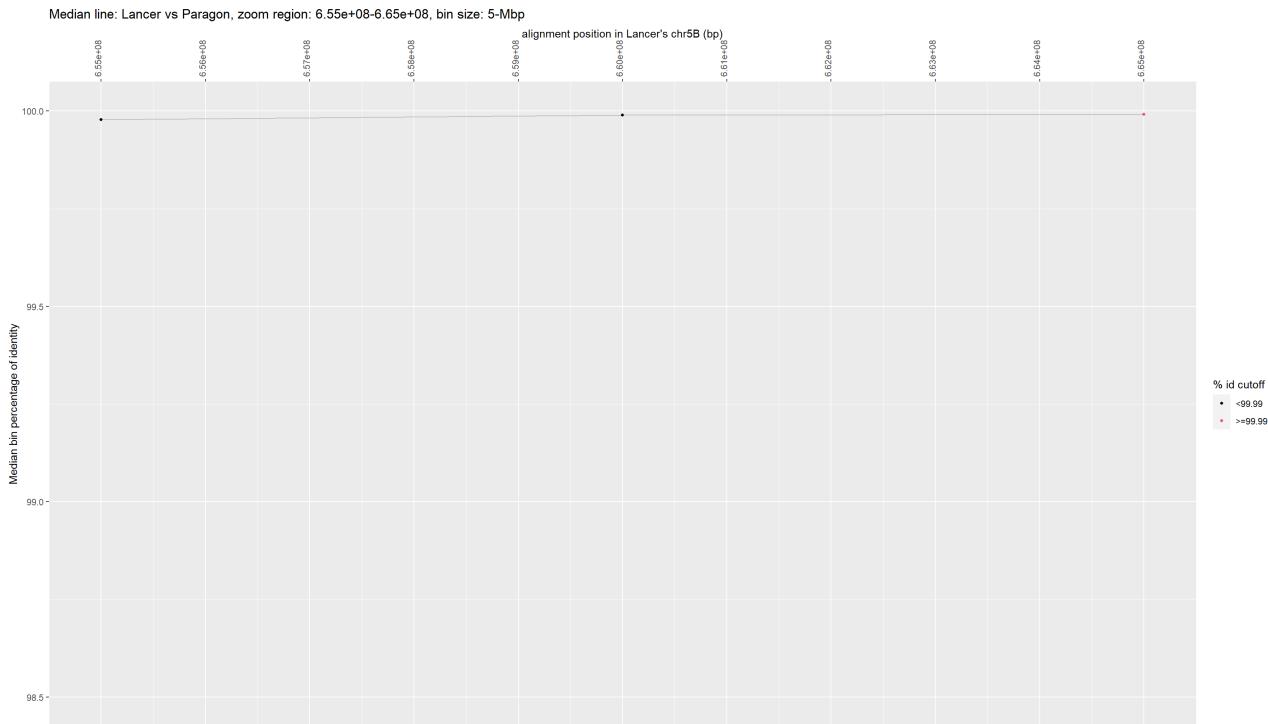
Boxplot: Lancer vs Paragon, zoom region: 6.55e+08-6.65e+08 Mbp, bin size: 2.5-Mbp, cutoff: 99.99



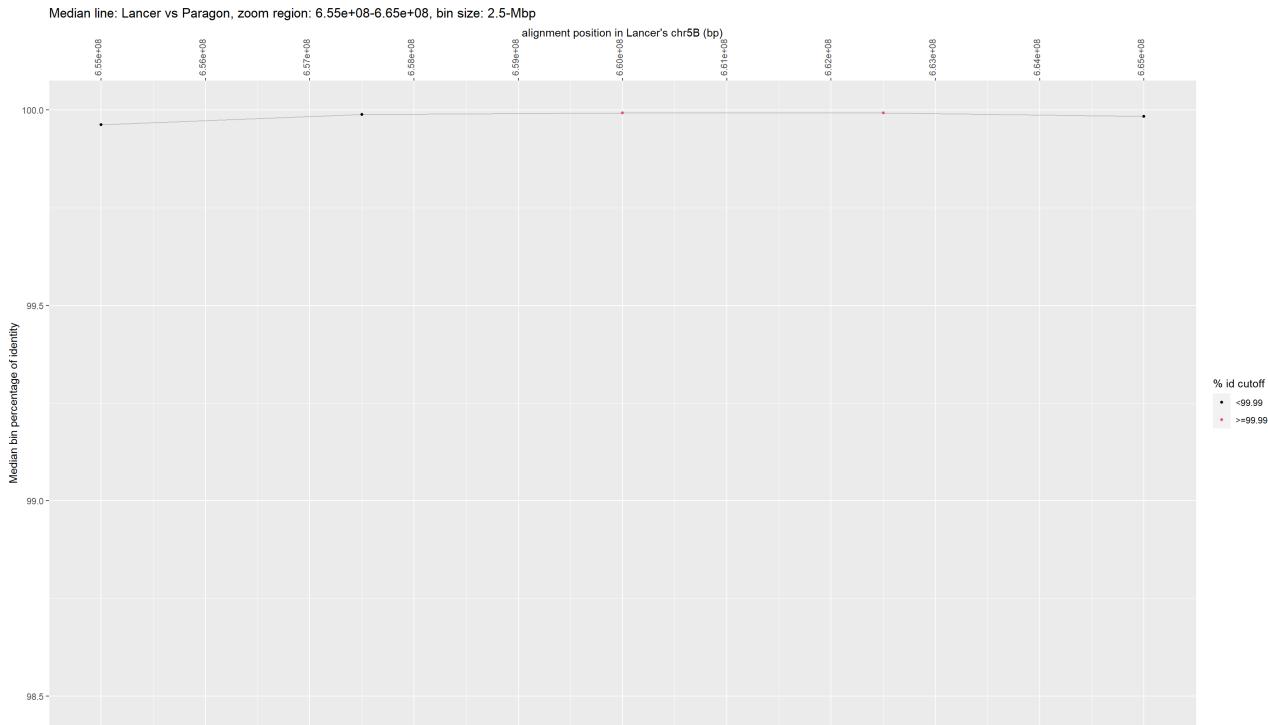


1.2.5. Plot the median line in the zoom region to see the haploblock predictions at different bin sizes

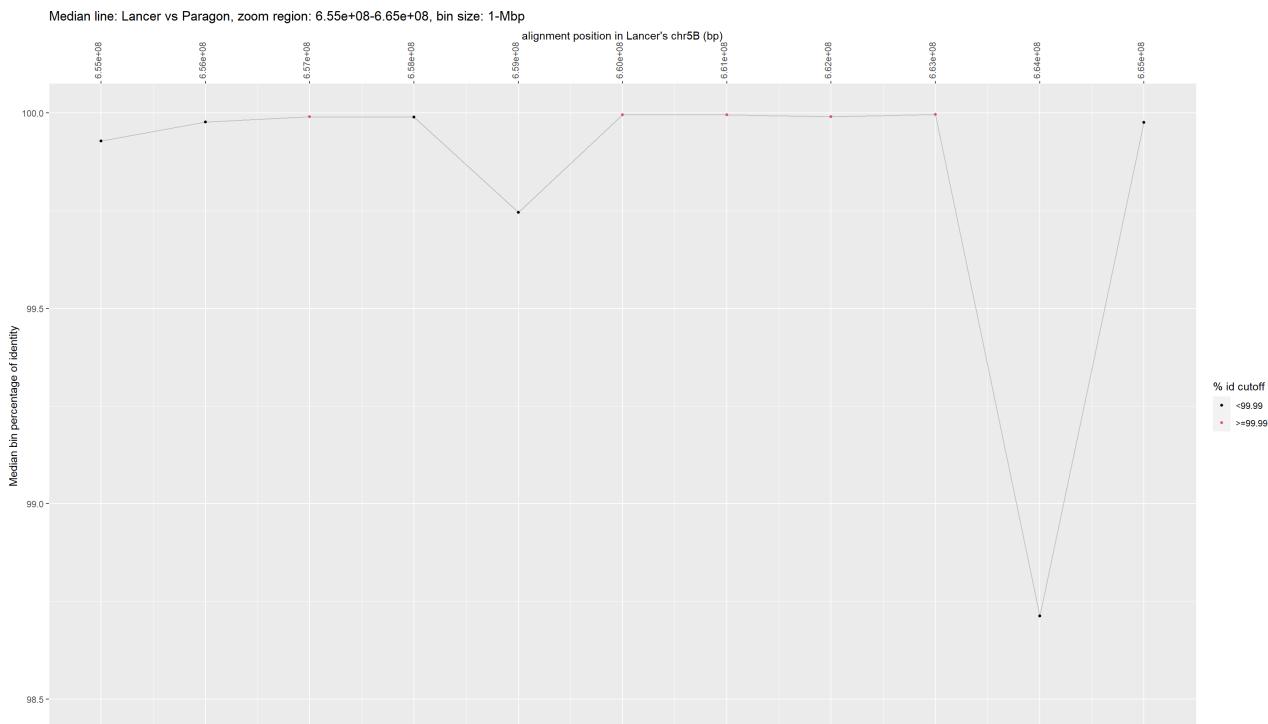
```
plot_line_bin_median(aln_subset, bin_size = 5000000, bin_start = zoom_start, bin_end = zoom_end, ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 1000000)
```



```
plot_line_bin_median(aln_subset, bin_size = 2500000, bin_start = zoom_start, bin_end = zoom_end, ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly , x_label_gap = 1000000)
```



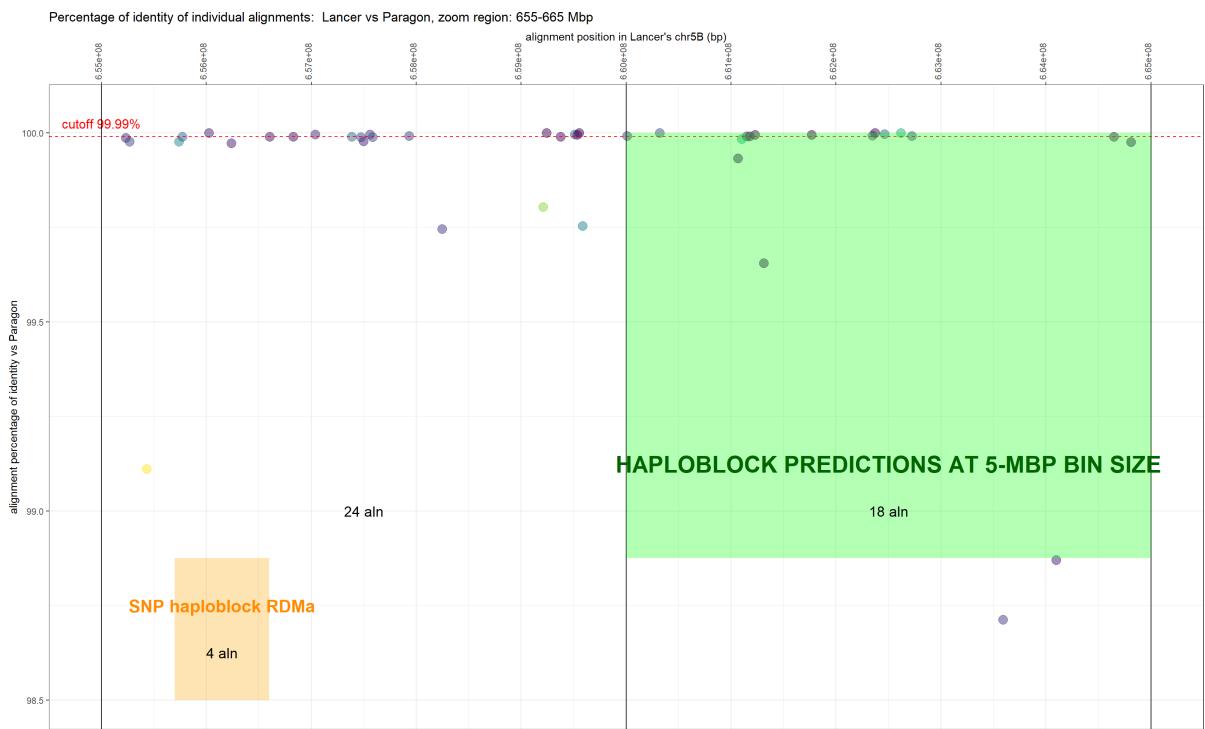
```
plot_line_bin_median(aln_subset, bin_size = 1000000, bin_start = zoom_start, bin_end = zoom_end, ymin = 98.5, cut_off = 99.9, reference_name = reference_assembly, query_name = query_assembly, x_label_gap = 1000000)
```



1.2.6. Dot-plot the zoom region with the haploblock predictions and the target region and make decisions regarding the start and end limits of the sequence-based haploblock

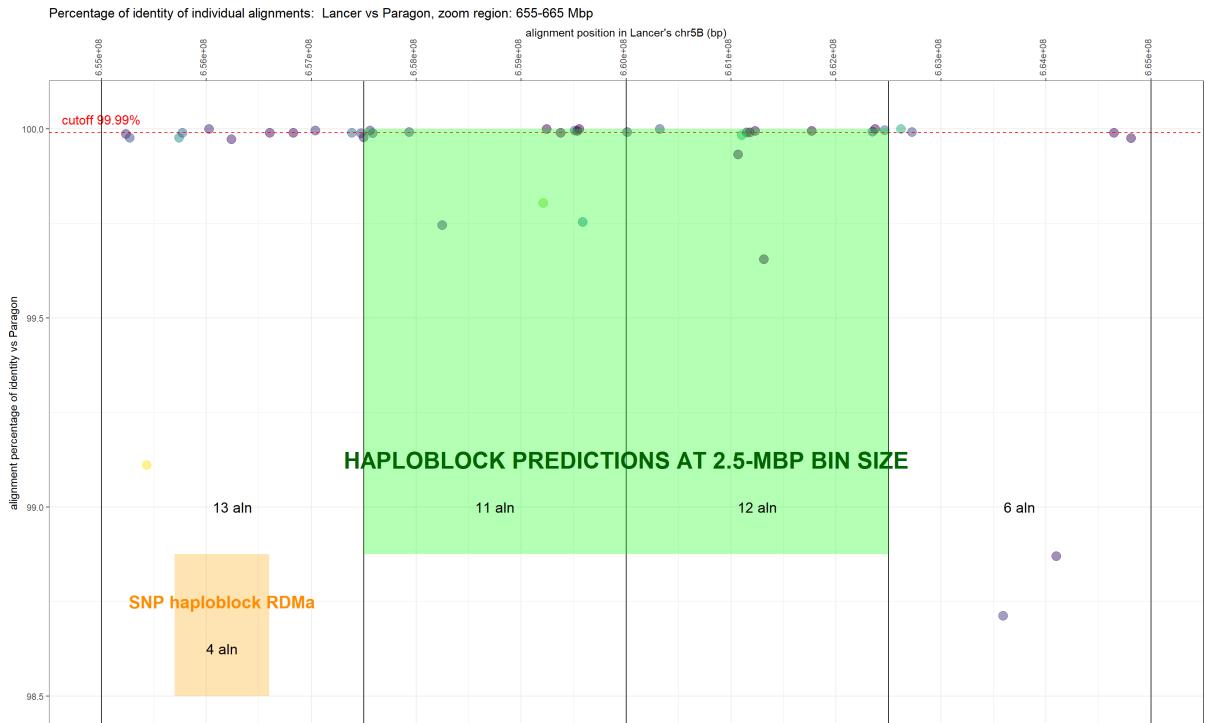
```
target <- data.frame(target_start, target_end)
plot_aln_and_bins(aln_subset = aln_subset, bin_size = 5000000, zoom_start = zoom_start, zoom_end = zoom_end, highlighted_target = target, target_text = "SNP haploblock RDMA", fill_target = "orange", color_target_text = "darkorange", fill_prediction = "green", color_prediction_text = "darkgreen", ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly, dot_size = 4, x_label_gap = 1000000)
```

```
## [1] "BINS AT 5-MBP BIN SIZE"
##      bin perc_id median cut_off block_no bin_start bin_end aln_number
## 1 6.55e+08    99.97763 <99.99 NO_BLOCK  6.50e+08 6.55e+08      37
## 2 6.60e+08    99.98986 <99.99 NO_BLOCK  6.55e+08 6.60e+08      24
## 3 6.65e+08    99.99159 >=99.99        1  6.60e+08 6.65e+08      18
## [1] "BLOCK SUMMARY AT 5-MBP BIN SIZE"
##      bin_size comparison block_no block_start block_end
## 1      5-Mbp Lancer->Paragon       1     6.6e+08 6.65e+08
```



```
plot_aln_and_bins(aln_subset = aln_subset, bin_size = 2500000, zoom_start = zoom_start, zoom_end = zoom_end, highlighted_target = target, target_text = "SNP haploblock RDMA", fill_target = "orange", color_target_text = "darkorange", fill_predictions = "green", color_prediction_text = "darkgreen", ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly, dot_size = 4, x_label_gap = 1000000)
```

```
## [1] "BINS AT 2.5-MBP BIN SIZE"
##      bin perc_id median cut_off block_no bin_start   bin_end aln_number
## 1 6550000000  99.96227 <99.99 NO_BLOCK 6525000000 6550000000     13
## 2 6575000000  99.98867 <99.99 NO_BLOCK 6550000000 6575000000     13
## 3 6600000000  99.99211 >=99.99          1 6575000000 6600000000     11
## 4 6625000000  99.99228 >=99.99          1 6600000000 6625000000     12
## 5 6650000000  99.98336 <99.99 NO_BLOCK 6625000000 6650000000      6
## [1] "BLOCK SUMMARY AT 2.5-MBP BIN SIZE"
##      bin_size      comparison block_no block_start block_end
## 1 2.5-Mbp Lancer->Paragon          1 6575000000 6625000000
```

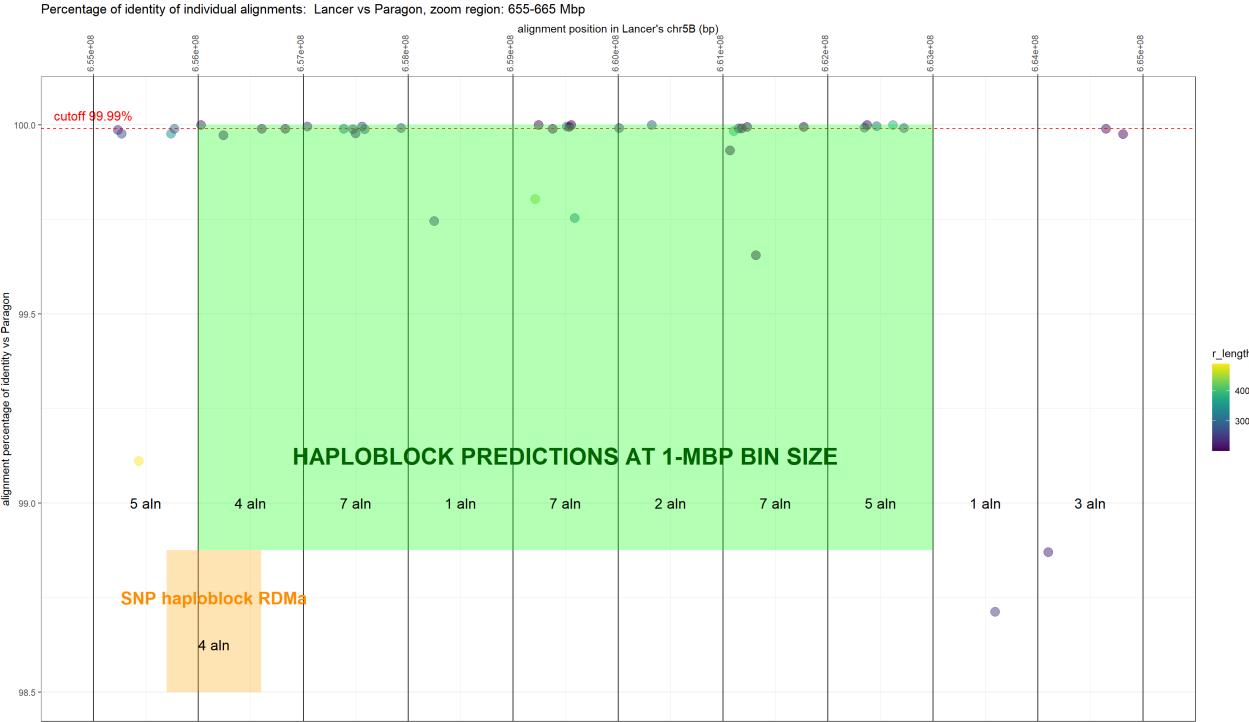


```
plot_aln_and_bins(aln_subset = aln_subset, bin_size = 1000000, zoom_start = zoom_start, zoom_end = zoom_end, highlighted_target = target, target_text = "SNP haploblock RDMA", fill_target = "orange", color_target_text = "darkorange", fill_predictions = "green", color_prediction_text = "darkgreen", ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly, dot_size = 4, x_label_gap = 1000000)
```

```

## [1] "BINS AT 1-MBP BIN SIZE"
##      bin_perc_id median cut_off block_no bin_start bin_end aln_number
## 1 6.55e+08      99.92852 <99.99 NO_BLOCK 6.54e+08 6.55e+08      6
## 2 6.56e+08      99.97730 <99.99 NO_BLOCK 6.55e+08 6.56e+08      5
## 3 6.57e+08      99.99048 >=99.99      1 6.56e+08 6.57e+08      4
## 4 6.58e+08      99.98996 <99.99 NO_BLOCK 6.57e+08 6.58e+08      7
## 5 6.59e+08      99.74574 <99.99 NO_BLOCK 6.58e+08 6.59e+08      1
## 6 6.60e+08      99.99528 >=99.99      1 6.59e+08 6.60e+08      7
## 7 6.61e+08      99.99595 >=99.99      1 6.60e+08 6.61e+08      2
## 8 6.62e+08      99.99072 >=99.99      1 6.61e+08 6.62e+08      7
## 9 6.63e+08      99.99664 >=99.99      1 6.62e+08 6.63e+08      5
## 10 6.64e+08     98.71248 <99.99 NO_BLOCK 6.63e+08 6.64e+08      1
## 11 6.65e+08     99.97621 <99.99 NO_BLOCK 6.64e+08 6.65e+08      3
## [1] "BLOCK SUMMARY AT 1-MBP BIN SIZE"
##      bin_size comparison block_no block_start block_end
## 1    1-Mbp Lancer->Paragon      1 6.56e+08 6.63e+08

```



Define new parameters:

```

selected_start <- 655760000 # Genes will be extracted from this position. If you have no interest in redefining your region,
# simply write 'target_start' or 'zoom_start', to keep with the previous coordinates
selected_end <- 662740000 # Genes will be extracted until this position. If you have no interest in redefining your region,
# simply write 'target_end' or 'zoom_end', to keep with the previous coordinates
target_text <- "Potential new haploblock" # Text to print on the selected region

selected_haploblock <- data.frame(selected_start, selected_end)
print(selected_haploblock)

```

```

## selected_start selected_end
## 1 655760000 662740000

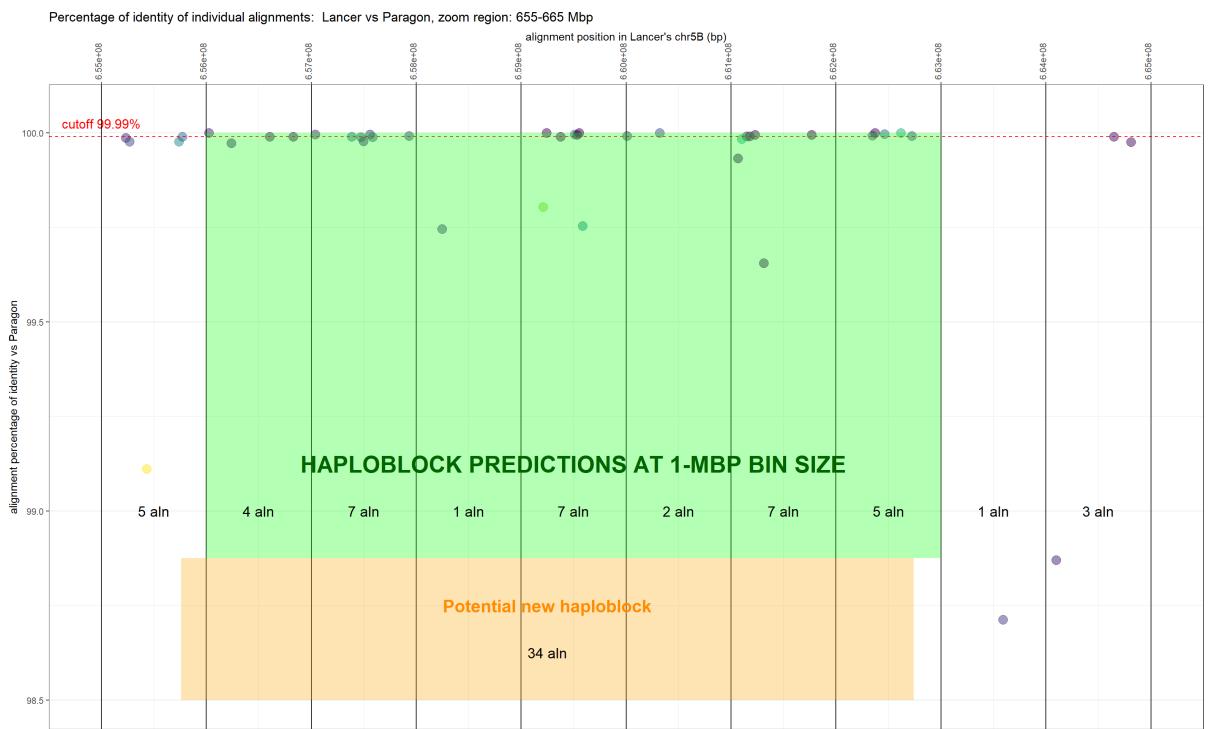
```

1.2.7. Plot summary graphs

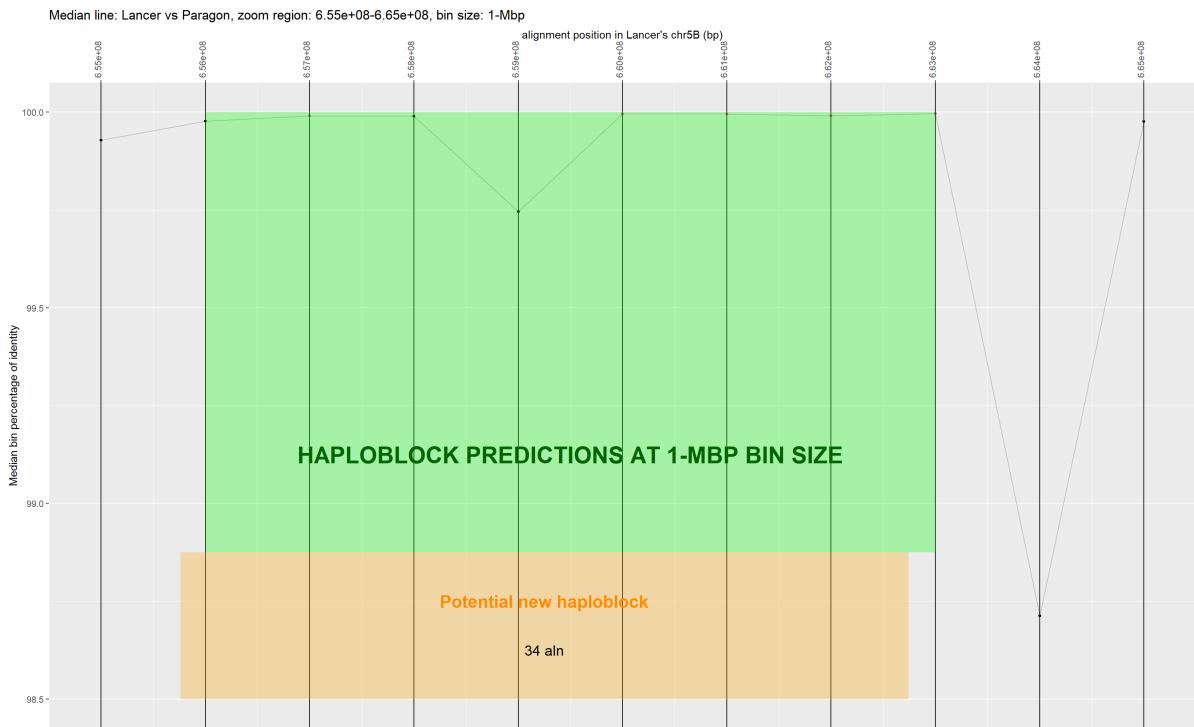
```

plot_aln_and_bins(print_tables = FALSE, aln_subset = aln_subset, bin_size = 1000000, zoom_start = zoom_start, zoom_end = zoom_end,
highlighted_target = selected_haploblock, target_text = target_text, fill_target = "orange", color_target_text = "darkorange",
fill_predictions = "green", color_prediction_text = "darkgreen", ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly,
query_name = query_assembly, dot_size = 4, x_label_gap = 1000000)

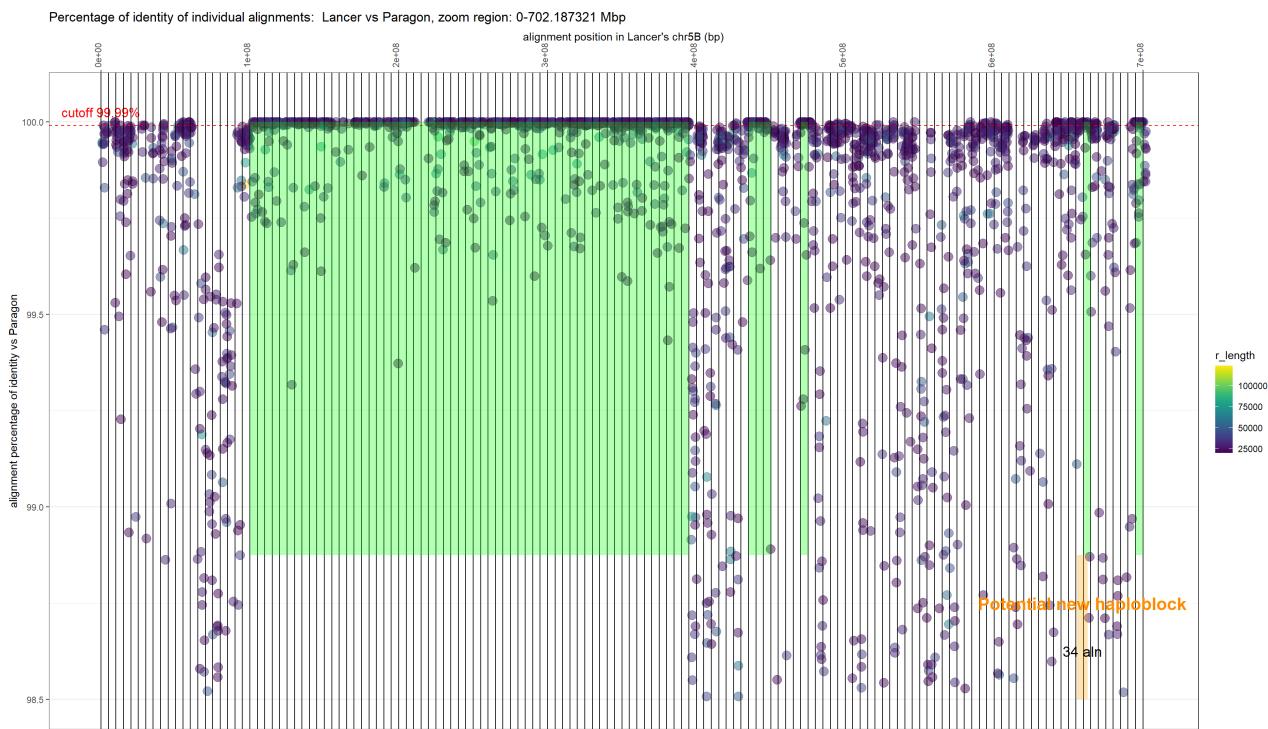
```



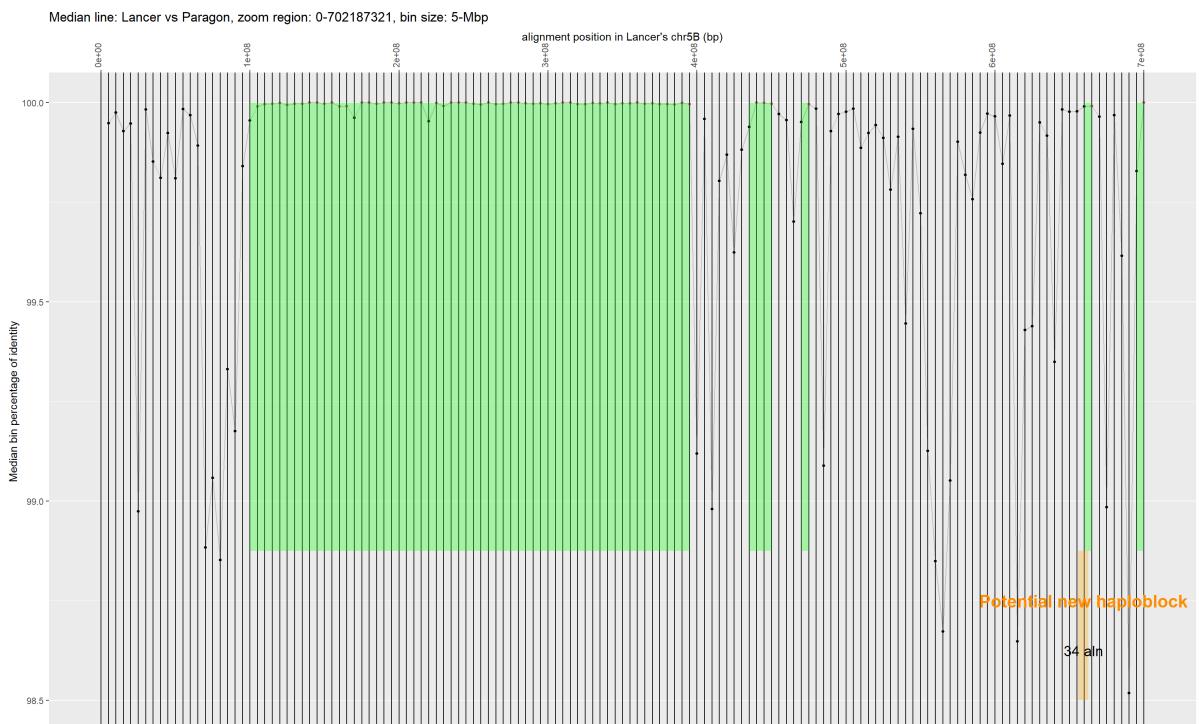
```
plot_bins_and_selected_region(print_tables = FALSE, aln_subset = aln_subset, bin_size = 1000000, zoom_start = zoom_start, zoom_end = zoom_end, highlighted_target = selected_haploblock, target_text = target_text, fill_target = "orange", color_target_text = "darkorange", fill_predictions = "green", color_prediction_text = "darkgreen", ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly, dot_size = 4, x_label_gap = 1000000)
```



```
plot_aln_and_bins(print_tables = FALSE, aln_subset = aln_subset, bin_size = 5000000, zoom_start = 0, zoom_end = max(aln_subset$re), highlighted_target = selected_haploblock, target_text = target_text, fill_target = "orange", color_target_text = "darkorange", fill_predictions = "green", ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly, dot_size = 4, x_label_gap = 10000000, prediction_text = FALSE, aln_text = F)
```



```
plot_bins_and_selected_region(print_tables = FALSE, aln_subset = aln_subset, bin_size = 5000000, zoom_start = 0, zoom_end = max(aln_subset$re), highlighted_target = selected_haploblock, target_text = target_text, fill_target = "orange", color_target_text = "darkorange", fill_predictions = "green", color_prediction_text = "darkgreen", ymin = 98.5, cut_off = 99.99, reference_name = reference_assembly, query_name = query_assembly, dot_size = 4, x_label_gap = 10000000, prediction_text = FALSE)
```



2. Identify genes in the haploblock

Requirements:

- File 'projectedGenes__Triticum_aestivum_REFERENCEassembly_v1.0.gff', in this case the reference is LongReach Lancer (downloadable from https://webblast.ipk-gatersleben.de/downloads/wheat/gene_projection/ (https://webblast.ipk-gatersleben.de/downloads/wheat/gene_projection/))
- File 'geneid_2_chinese.sourceid.txt' (downloadable from https://webblast.ipk-gatersleben.de/downloads/wheat/gene_projection/ (https://webblast.ipk-gatersleben.de/downloads/wheat/gene_projection/))

2.1. Download the files and save them in the working directory

2.2. Read the gff file containing the gene model projection for the reference assembly

```
ref_gff <- read.table ( file = "projectedGenes__Triticum_aestivum_LongReach_Lancer_v1.0.gff", sep = "\t" , header = F, stringsAsFactors = F)
```

2.3. Create subsets for gene projections only and edit the table

```
ref_gff_gene_only <- ref_gff[ref_gff$V3 == "gene",]  
colnames(ref_gff_gene_only) <- c("chr", "annotation", "biotype", "start", "end", "score", "strand", "info", "ref_id")  
ref_gff_gene_only$var_id <- gsub("ID=", "", ref_gff_gene_only$ref_id)
```

2.4. Create a subset for the haploblock

```
my_genes <- ref_gff_gene_only[grep(chromosome, ref_gff_gene_only$chr),]  
my_genes <- my_genes[(my_genes$start >= selected_start) & (my_genes$end <= selected_end),]
```

2.5. Download 'geneid_2_chinese.sourceid.txt' from the last link and read the table

```
cs_id <- read.table(file = "geneid_2_chinese.sourceid.txt", sep="\t", header=T, stringsAsFactors = F)
```

2.6. Add a column to the subset with the IDs of the gene sources in Chinese Spring for each of Lancer's genes

```
my_genes <- cs_id.filler(data = my_genes, library = cs_id, chr = chromosome, ref.var = tolower(reference_assembly))
```

2.7. Extract the names of the genes in chromosome 5B separated by commas

```
my_gene_sources <- as.character(my_genes$cs_id[!grepl("^source", my_genes$cs_id)])  
my_gene_sources_text <- paste(my_gene_sources, collapse = ", ")  
write.table(x = my_gene_sources_text, file = "my_gene_sources.txt", sep = "", row.names = F, col.names = F, quote = F)  
  
my_variety_genes <- as.character(my_genes$var_id[!grepl("^source", my_genes$cs_id)])  
my_variety_genes_text <- paste(my_variety_genes, collapse = ", ")  
write.table(x = my_variety_genes_text, file = "my_variety_genes.txt", sep = "", row.names = F, col.names = F, quote = F)  
  
print_my_genes <- data.frame( "var_id" = my_genes$var_id, "chinese_id" = my_genes$cs_id)  
write.csv2(x = print_my_genes, file = "my_genes_and_sources.csv", row.names = F, quote = F)  
  
head(print_my_genes, 10)
```

```
##           var_id      chinese_id  
## 1 TraesLAC5B01G535100 TraesCS5B02G495900  
## 2 TraesLAC5B01G535200 TraesCS5B02G496000  
## 3 TraesLAC5B01G535300 TraesCS5B02G496100  
## 4 TraesLAC5B01G535400 TraesCS5B02G496200  
## 5 TraesLAC5B01G535500 TraesCS5B02G496300  
## 6 TraesLAC5B01G535600 TraesCS5B02G496400  
## 7 TraesLAC5B01G535700 TraesCS5B02G496600  
## 8 TraesLAC5B01G535800 TraesCS5B02G496700  
## 9 TraesLAC5B01G535900 TraesCS5B02G496900  
## 10 TraesLAC5B01G536000 TraesCS5B02G497100
```

3. Check if the haplotype prediction was also made by gene-based BLAST pairwise alignments

Requirements:

- File 'varieties_all_identities_2000bp.tar.gz' (downloadable from https://opendata.earlham.ac.uk/wheat/under_license/toronto/Brinton_et.al_2020-05-20-Haplotypes-for-wheat-breeding/pairwise_blast/ (https://opendata.earlham.ac.uk/wheat/under_license/toronto/Brinton_et.al_2020-05-20-Haplotypes-for-wheat-breeding/pairwise_blast/))
- File 'iwgsc_refseq_v1.2_gene_annotation.zip' (downloadable from https://urgi.versailles.inrae.fr/download/iwgsc/IWGSC_RefSeq_Annotations/v1.2/ (https://urgi.versailles.inrae.fr/download/iwgsc/IWGSC_RefSeq_Annotations/v1.2/))

3.1. Download the zip file, decompress it in the working directory and read the tables

```
HC_gtf <- read.table("IWGSC_v1.2_HC_20200615.gff3", sep = "\t", header = FALSE, stringsAsFactors = FALSE)  
LC_gtf <- read.table("IWGSC_v1.2_LC_20200615.gff3", sep = "\t", header = FALSE, stringsAsFactors = FALSE)  
ALL_gtf <- rbind(HC_gtf, LC_gtf)
```

3.2. Extract the BLAST alignments and put them in table with Chinese Spring genes and their location in the IWGSC genome (long process)

```
BLAST_library <- read_pairwise_position(blast_path_gz = "varieties_all_identities_2000bp.tab.gz", gtf = ALL_gtf, write_table = "BLAST_library.tab")
```

3.3. Extract a subset with Lancer-Paragon comparison in chr5B

```
BLAST_subset <- BLAST_library[grep1(tolower(reference_assembly), "->", tolower(query_assembly), sep = ""), BLAST_library$aln_type) & grep1(chromosome, BLAST_library$chr),]
```

3.4. Use the vector with the genes identified in the previous step to extract another subset with only the genes in our haploblock

```
BLAST_subset <- BLAST_subset[grep1(paste(my_gene_sources, collapse = "|"), BLAST_subset$transcript),]  
write.csv2(x = BLAST_subset, file = "BLAST_subset.csv", row.names = F, quote = F)
```

Brinton et al (2020) only retained only gene projections consistent with the expected chromosome. This explains why this subset contains less genes than the total amount of annotated genes in our region.

```
colnames(BLAST_subset)[1] <- "cs_transcript"  
colnames(BLAST_subset)[14] <- "cs_start"  
colnames(BLAST_subset)[15] <- "cs_end"  
  
BLAST_subset$var_transcript <- my_genes$var_id[grep1(paste(BLAST_subset$cs_transcript, collapse = "|"), my_genes$cs_id)]  
BLAST_subset$var_start <- my_genes$start[grep1(paste(BLAST_subset$var_transcript, collapse = "|"), my_genes$var_id)]  
BLAST_subset$var_end <- my_genes$end[grep1(paste(BLAST_subset$var_transcript, collapse = "|"), my_genes$var_id)]
```

Notice that some genes were filtered out if they contained more than one projection in the expected chromosome, so the amount of genes shown in this table is expected to be smaller than those found in lancer's projected genes in the region. The genes that were excluded can be seen here:

```
my_gene_sources[my_gene_sources %in% BLAST_subset$cs_transcript == FALSE]  
  
## [1] "TraesCSS5B02G496400" "TraesCSS5B02G496700" "TraesCSS5B02G497300"  
## [4] "TraesCSS5B02G497500" "TraesCSS5B02G500500" "TraesCSS5B02G500600"  
## [7] "TraesCSS5B02G500700" "TraesCSS5B02G502700" "TraesCSS5B02G503400"  
## [10] "TraesCSS5B02G503500" "TraesCSS5B02G504400" "TraesCSS5B02G552700"  
  
my_genes$var_id[grep1(paste(my_gene_sources[my_gene_sources %in% BLAST_subset$cs_transcript == FALSE], collapse = "|"), my_genes$cs_id)]  
  
## [1] "TraesLAC5B01G535600" "TraesLAC5B01G535800" "TraesLAC5B01G536100"  
## [4] "TraesLAC5B01G536300" "TraesLAC5B01G539700" "TraesLAC5B01G539800"  
## [7] "TraesLAC5B01G539900" "TraesLAC5B01G542300" "TraesLAC5B01G542700"  
## [10] "TraesLAC5B01G542900" "TraesLAC5B01G543700" "TraesLAC5B01G544200"
```

We can extract another list containing only the genes that were present among the BLAST alignments

```
my_gene_sources_filtered_by_Brinton <- my_gene_sources[my_gene_sources %in% BLAST_subset$cs_transcript == TRUE]  
my_variety_genes_filtered_by_Brinton <- my_genes$var_id[grep1(paste(my_gene_sources %in% BLAST_subset$cs_transcript == TRUE), collapse = "|"), my_genes$cs_id]  
my_gene_sources_filtered_by_Brinton_text <- paste(my_gene_sources_filtered_by_Brinton, collapse = ", ")  
write.table(x = my_gene_sources_filtered_by_Brinton_text, file = "my_gene_sources_filtered_by_Brinton_text.txt", sep = "", row.names = F, col.names = F, quote = F)  
  
my_variety_genes_filtered_by_Brinton_text <- paste(my_variety_genes_filtered_by_Brinton, collapse = ", ")  
write.table(x = my_variety_genes_filtered_by_Brinton_text, file = "my_variety_genes_filtered_by_Brinton_text.txt", sep = "", row.names = F, col.names = F, quote = F)  
  
print_blast <- data.frame("var_id" = my_variety_genes_filtered_by_Brinton_text, "chinese_id" = my_gene_sources_filtered_by_Brinton_text)  
write.csv2(x = print_blast, file = "my_genes_and_sources_filtered_by_Brinton.csv", row.names = F, quote = F)
```

3.5. Calculate the percentage of identity in windows of 20 genes where genes containing Ns are filtered out

```
BLAST_subset <- BLAST_subset[ BLAST_subset$Ns_total == 0, ]  
  
window_BLAST_subset <- edited_calculate_pid_windows(aln_data = BLAST_subset)  
blocks_BLAST_subset <- assign_blocks(window_BLAST_subset)
```

blocks_BLAST_subset

```

## start end cs_start cs_end var_start var_end pident_mean
## 1 20 663134978 666196835 655761539 658486203 99.99828
## 2 21 663234734 666207172 656020702 658496928 99.99761
## 3 22 663796017 666209291 656553758 658499392 99.99761
## 4 23 664022132 666372136 656600250 658525735 99.99742
## 5 24 664672987 668463343 657035254 659882405 99.99742
## 6 25 664682920 668474547 657044824 659893484 99.99609
## 7 26 664726808 668519308 657104030 659938479 99.99609
## 8 27 664949927 668583715 657329620 660019299 99.99498
## 9 28 664979682 669217640 657359517 660818003 99.99587
## 10 29 664998106 669232964 657377616 660827247 99.99460
## 11 30 665016776 669292146 657396273 660884339 99.99254
## 12 31 665024295 669306368 657403038 660900818 99.99387
## 13 32 665170918 669452267 657534486 661060959 99.99280
## 14 33 665539582 669483645 657903967 661091774 99.99280
## 15 34 665560612 669898631 657925163 661273027 99.99280
## 16 35 665719477 669909637 658082849 661430757 99.99382
## 17 36 665767971 670099306 658131416 661722411 99.99146
## 18 37 665938590 670116476 658242152 661739162 99.99146
## 19 38 665939237 670118309 658242780 661741296 99.99146
## 20 39 666195640 670130079 658485922 661752676 99.99146
## 21 40 666202129 670246635 658492724 661761415 99.99146
## 22 41 666207065 670305619 658497166 661773746 99.99213
## 23 42 666371429 670694936 658525028 662378160 99.99213
## 24 43 668461124 670716758 659880577 662400081 99.99449
## aln_type start_cs_transcript end_cs_transcript start_var_transcript
## 1 lancer->paragon TraesCS5B02G495900 TraesCS5B02G500000 TraesLAC5B01G535100
## 2 lancer->paragon TraesCS5B02G496300 TraesCS5B02G500100 TraesLAC5B01G535500
## 3 lancer->paragon TraesCS5B02G496600 TraesCS5B02G500200 TraesLAC5B01G535700
## 4 lancer->paragon TraesCS5B02G497100 TraesCS5B02G500300 TraesLAC5B01G536000
## 5 lancer->paragon TraesCS5B02G497700 TraesCS5B02G501100 TraesLAC5B01G536500
## 6 lancer->paragon TraesCS5B02G497800 TraesCS5B02G501200 TraesLAC5B01G536600
## 7 lancer->paragon TraesCS5B02G497900 TraesCS5B02G501300 TraesLAC5B01G536700
## 8 lancer->paragon TraesCS5B02G498000 TraesCS5B02G501400 TraesLAC5B01G537100
## 9 lancer->paragon TraesCS5B02G498100 TraesCS5B02G501600 TraesLAC5B01G537200
## 10 lancer->paragon TraesCS5B02G498200 TraesCS5B02G501800 TraesLAC5B01G537300
## 11 lancer->paragon TraesCS5B02G498300 TraesCS5B02G501900 TraesLAC5B01G537400
## 12 lancer->paragon TraesCS5B02G498400 TraesCS5B02G502100 TraesLAC5B01G537500
## 13 lancer->paragon TraesCS5B02G498500 TraesCS5B02G502200 TraesLAC5B01G537600
## 14 lancer->paragon TraesCS5B02G498700 TraesCS5B02G502300 TraesLAC5B01G537800
## 15 lancer->paragon TraesCS5B02G498800 TraesCS5B02G503200 TraesLAC5B01G537900
## 16 lancer->paragon TraesCS5B02G498900 TraesCS5B02G503300 TraesLAC5B01G538000
## 17 lancer->paragon TraesCS5B02G499100 TraesCS5B02G503700 TraesLAC5B01G538200
## 18 lancer->paragon TraesCS5B02G499300 TraesCS5B02G503800 TraesLAC5B01G538400
## 19 lancer->paragon TraesCS5B02G499400 TraesCS5B02G503900 TraesLAC5B01G538500
## 20 lancer->paragon TraesCS5B02G500000 TraesCS5B02G504000 TraesLAC5B01G539000
## 21 lancer->paragon TraesCS5B02G500100 TraesCS5B02G504100 TraesLAC5B01G539100
## 22 lancer->paragon TraesCS5B02G500200 TraesCS5B02G504200 TraesLAC5B01G539200
## 23 lancer->paragon TraesCS5B02G500300 TraesCS5B02G504600 TraesLAC5B01G539300
## 24 lancer->paragon TraesCS5B02G501100 TraesCS5B02G504700 TraesLAC5B01G540400
## end_var_transcript block_no
## 1 TraesLAC5B01G539000 NA
## 2 TraesLAC5B01G539100 NA
## 3 TraesLAC5B01G539200 NA
## 4 TraesLAC5B01G539300 NA
## 5 TraesLAC5B01G540400 NA
## 6 TraesLAC5B01G540500 NA
## 7 TraesLAC5B01G540600 NA
## 8 TraesLAC5B01G540700 NA
## 9 TraesLAC5B01G541300 NA
## 10 TraesLAC5B01G541500 NA
## 11 TraesLAC5B01G541600 NA
## 12 TraesLAC5B01G541800 NA
## 13 TraesLAC5B01G541900 NA
## 14 TraesLAC5B01G542000 NA
## 15 TraesLAC5B01G542400 NA
## 16 TraesLAC5B01G542600 NA
## 17 TraesLAC5B01G543100 NA
## 18 TraesLAC5B01G543200 NA
## 19 TraesLAC5B01G543300 NA
## 20 TraesLAC5B01G543400 NA
## 21 TraesLAC5B01G543500 NA
## 22 TraesLAC5B01G543600 NA
## 23 TraesLAC5B01G543800 NA
## 24 TraesLAC5B01G543900 NA

```

According to the previous data, no haplotype is predicted by BLASTing the gene sequences from this set between Lancer and Paragon (block_no = NA). In this case, the NUCmer-based prediction should be taken as the reference, as Brinton et al. (2020) explained that genic sequence alone was insufficient to fully differentiate between haplotypes.

Supplementary file 2: List of potential SNP markers to target the haplotype RDMb-h3, associated with higher root dry mass. SNPs were identified by high-stringency mapping of Illumina paired-end short reads and variant calling against Chinese Spring from three resequenced Chinese varieties by Hao et al. (2020). The SNP flanking sequence corresponds to those of the h2-carrying assemblies. The alleles expected for the RDMb-h2 and RDMb-h3 blocks are indicated alongside with the start and end coordinates in each of the h2-carrying assemblies.

| | | | | | | | | | | |
|-----------------|--|---|---|-----------|-----------|------------|------------|-----------|------------|------------|
| RDmB-h3:SNP_311 | ACAGGCAAGAACATTCTCGGCCAACTGAAGGCCGGAGAGGGTGTCTGAAAGAGTCAGAAAGAACAGTCAGATGCCAGATGCCAGCTGGAA | A | T | 653313146 | 653313246 | 644596395 | 653610092 | 653610192 | 655189300 | 65518940 |
| RDmB-h3:SNP_312 | AGTGGCCACAGCTTGTAAGTCCTGGCCACTTCTCTGTTGATCTGCAGGCTCTGGCAGCTGAAAGTCAGGCTCTGGCAGCTGGCTTC | G | T | 653314807 | 653314907 | 644598056 | 653611753 | 653611835 | 655190961 | 65519106 |
| RDmB-h3:SNP_313 | TCTCATGGGCCGATTACTGCCTTATGTCGAGGTTGTTCTGCATCGGCCGGTGTATCGAGATTCGCTCGGCCGCTTCGGTTG | A | G | 653315862 | 653315962 | 644599111 | 653612082 | 653612088 | 655192016 | 65519211 |
| RDmB-h3:SNP_314 | GAGGCAACCGCGGAGGTGACCCGTCGACGACATGAGGAGTACCGCTCTGAGCCCCTGACTCACAGCAGAGGAAGAGCTCGCCGCGGAA | T | C | 653321007 | 653321107 | 644604256 | 653617953 | 653618053 | 655197161 | 65519726 |
| RDmB-h3:SNP_315 | CGCTGAGAACCTGGACATCTCGCATGAAACCTCTGAGATCCGGGTTTACCCAGGGAGTCTGGCATGCTGAGCTGAGCAAGGAAAC | T | G | 653323546 | 653323646 | 644606795 | 653620492 | 653620592 | 655199700 | 65519980 |
| RDmB-h3:SNP_316 | TTTACCACTCGGAGTCTGCCAATGTCGATGCTGCCCCAAAGAGATAATTCTTCGATGTCAGCAGGATTACATGGGGCTTCGGCC | T | C | 653323743 | 653323843 | 644606992 | 653620689 | 653620709 | 655199897 | 65519999 |
| RDmB-h3:SNP_317 | GGCAATCTGCTCATGGTCCCCAAAGAGATAATTCTTCGATGTCAGCAGGATTACATGGGGCTTCGGCCGAAAGATAATTCTTC | C | T | 653323763 | 653323863 | 644607012 | 653620709 | 653620809 | 655199917 | 65520009 |
| RDmB-h3:SNP_318 | TCGATGTCACCGATTCAACGATATCAATGGGGCTGCCGAAAGATAATTCTTCCTCCCGATCAGGAAATAGTCGACTCACTGGGGATGT | T | C | 653323801 | 653323901 | 644607050 | 653620747 | 653620847 | 655199955 | 65520009 |
| RDmB-h3:SNP_319 | TTTACATCACCCTTGGCTGCTCTGTTGACCTGATTCGGGCAAAAGTCAGGCGAACCTGGCTGATGAGGATGATTGAAAGTGTCTCAC | A | T | 653323797 | 653324079 | 644607228 | 653620925 | 653621025 | 655200133 | 65520023 |
| RDmB-h3:SNP_320 | GTGGGTGCGATGCTGAGAACATTACTGGATGCTGGCCATTAACTGATTCATCTCGCTTGGCTGAAAGGATCTACCGT | T | G | 653323474 | 653324244 | 644607573 | 653621370 | 653621370 | 655200478 | 65520057 |
| RDmB-h3:SNP_321 | AGAAAACATTCTGGATGCTGGCGATTAACTGATTCATCTCGCTTGGTAAAGGCGATTCTTACCGCTGATGAGAAGTCAGAACAGTC | T | C | 653324343 | 653324443 | 644607592 | 653621289 | 653621389 | 655200497 | 65520059 |
| RDmB-h3:SNP_322 | TGGATCAGCAGAACACACAAAGCTGTCCTACTGGGAAGGTCCTTAAATGTCGACCAAGGAAAGTCCTTAAATGGGGCATACGCCCTAACATGTCAG | A | C | 653326547 | 653326647 | 644609796 | 653623493 | 653623593 | 655202701 | 65520280 |
| RDmB-h3:SNP_323 | GGGTAGAAGGTCCTTGGCTCATGGTCTGGCTGAGGCTGTTGCTCTGGCTGAGGACACTAGCTGTTCTCTGCTGCTGAGGCG | T | C | 653328297 | 653328392 | 644611558 | 653625238 | 653625238 | 655204446 | 65520454 |
| RDmB-h3:SNP_324 | AGAAAGCAGCGGAAGTCATGCTTCTGTTGCTGAGGACACTAGCTGTTCTCTGCTGCTTGGCTGAGGACAGAGCG | T | C | 653328308 | 653328408 | 644611574 | 653625254 | 653625254 | 655204462 | 655204546 |
| RDmB-h3:SNP_325 | TGAGTGGACGGGACAAGGAGACAGGGACATCAGACGCCAACCTGGCGAAGACTCTCCACTCTGACTGGAATCTGTCAGTCAGTC | C | A | 653328380 | 653328480 | 644611646 | 653625326 | 653625426 | 655204534 | 65520463 |
| RDmB-h3:SNP_326 | TCGTTGGCTGGCTGGTGGCTGCTCTGGCTCATCTGGCCATACATTCGGGACTTGGCTGTTGGCTGCTGGCTGGCTGGAGTGT | A | G | 653329764 | 653329864 | 644613130 | 653626710 | 653626810 | 655205918 | 655206018 |
| RDmB-h3:SNP_327 | GGCGTTGGCATGGCTTCTGGCTCATCATAACTGGGGCTGGCTGGCTGGCTGGCTGGCTGGCTGGCTGGAGTGGAGAAGA | G | A | 653329773 | 653329873 | 644613039 | 653626719 | 653626819 | 655205927 | 65520602 |
| RDmB-h3:SNP_328 | CGGTGGCATGAGTACAACACTGAGTGGTTCTGAAACAGACAGGACAAAGTCATGACTGTTACGAGAACAGATAAGGAGGAGCT | G | A | 653330643 | 653330743 | 644613090 | 653627589 | 653627689 | 655206797 | 65520689 |
| RDmB-h3:SNP_329 | CAGCTCATGAGGGGGCAATCTGGCTGAGTGGTTCTCCAGTGTCTGCTTACCTGGCTGATGAGGAACTGATGAAAGGAGGAAACCAACT | T | C | 653338187 | 653338287 | 644621453 | 653635133 | 653635233 | 655214341 | 655214444 |
| RDmB-h3:SNP_330 | CTGGGGCGCAATCGAGCTGGTCTCCAGTGTCTGACTTCTGATGAAATAGGAACGTAAGTGGACCAAAACATATTGACCC | T | C | 653338195 | 653338295 | 644621461 | 653635141 | 653635241 | 655214349 | 655214439 |
| RDmB-h3:SNP_331 | CGGGTGTAGTGTAGTGAAGGGCTGGATGATTGGATTGGCTTCTGTCGAGGCTGGCGAGTGGCTGGATGAGGCTGGAGAG | G | C | 653340048 | 653340184 | 644623350 | 653637030 | 653621638 | 655216338 | 655216338 |
| RDmB-h3:SNP_332 | ATGTAAGGCTGCCATGGATTAGGCTTCTGTCGAGGTGGCCGGCAGCCTGTCGCGATGTTGTCGAGGCGAGGGCTCGGATCC | C | A | 653340097 | 653340197 | 644623363 | 653637043 | 653637143 | 655216251 | 655216353 |
| RDmB-h3:SNP_333 | CGGGATGTTGCTGGCTGGAGCTGGCTGCTGGCTGAGCTGGCTGAGCTGGCTGAGCTGGCTGAGCTGGCTGGCTGGAGTTC | T | C | 653340159 | 653340259 | 3114195309 | 1314195309 | 15029270 | 15029270 | 324938575 |
| RDmB-h3:SNP_334 | AGCGGGATCATGTCCTTACCTGGGGAGCAGCTTACGGTGTGGTAGTCTGGGGCTGGCTGAGCTGGCTGAGCTGGCTGAGGACT | A | G | 653340902 | 653341002 | 644624168 | 653637848 | 653637948 | 655217156 | 655217156 |
| RDmB-h3:SNP_335 | CCACCTCTAATCTGGGGAGCAGCTTACGGTGTGGTAGCTTGGGGAGGCTGGCTGAGCTGGCTGAGCTGGCTGAGCTGG | T | G | 653341019 | 653341019 | 644624285 | 653637865 | 653637865 | 655217173 | 655217173 |
| RDmB-h3:SNP_336 | AGCAGATCTGATCAGCTTAAGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGAGGGATCTCT | C | T | 653341020 | 653341120 | 644624286 | 653637966 | 653638066 | 65521727 | 65521727 |
| RDmB-h3:SNP_337 | ATCTGATGACCTGCTGAGCTGGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGAGGATCT | C | A | 653341025 | 653341291 | 644624291 | 653638071 | 653638071 | 65521727 | 65521727 |
| RDmB-h3:SNP_338 | CTTTTAAATTATGATGATGTTGGTAAAGAGATGAGTACCATACAGGAGGAAAGGCTTACCTGGCTGAGCTGGCTGAGGCTGGCTGG | G | T | 653342402 | 653342502 | 644625768 | 653639348 | 653639448 | 655218556 | 655218556 |
| RDmB-h3:SNP_339 | TATTCCTCTTCTGAGCTAAGGACACAGCTCTGCCAACTACTCTGAGTCTGCTGAGGAGACTCTGTCACCGCTGATCC | G | A | 653344415 | 653344515 | 644627680 | 653641360 | 653641460 | 655220658 | 655220665 |
| RDmB-h3:SNP_340 | TCCTGGGAAAGAACATTCTGGATGCTGAGCTTACCTGGCTCTGAAAGACCCCTGGAGGATGTTGGTGGAGAAGAAATTCTCTGGAGGCA | A | C | 653348853 | 653348953 | 644632118 | 653645728 | 653645789 | 655225006 | 655225102 |
| RDmB-h3:SNP_341 | CTGCACTTAACTGATCTGCTTACAGACCCCTGAAAGATGTTGGTGTGAGAAGAAATTACCTGGCTGGCTGGAGGCTGGCTGG | G | T | 653348876 | 653348976 | 644632241 | 653645821 | 653645821 | 655225129 | 655225129 |
| RDmB-h3:SNP_342 | GTGGTGTGAAACTGGTGTAGTGTAGGAGTACAGCCTGGCAGAACAGAGCAGCTGGCTGAAATGTCGAGTACACTGGTGAAGCTATCA | C | G | 653350456 | 653350536 | 644633701 | 653647481 | 653647481 | 655226689 | 655226689 |
| RDmB-h3:SNP_343 | TGGTAGGAAAAGAACATTCTGGCTGGCTGGACTAACATTAACTAAGGAAATTAAATTGTCGCTGGCTGAGTCAACCTGGCTGG | T | C | 653355899 | 653359099 | 644642464 | 653655944 | 653656044 | 6552352512 | 6552352512 |
| RDmB-h3:SNP_344 | TTCTGTTTCTGAGCAGAACATTCTGATCATCTGCTTCTGAGGAGCTTGGCTGAGCTGGCTGAGCTGGCTGAGCTGGCTGAGG | G | A | 653360235 | 653360335 | 644643500 | 653657180 | 653657280 | 65523648 | 65523648 |
| RDmB-h3:SNP_345 | CAAACAAACGGTAACTGGGAGACTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | T | C | 653362267 | 653362370 | 644645535 | 653659215 | 653659315 | 655238423 | 655238423 |
| RDmB-h3:SNP_346 | GAATCAATATTCTCATGATGAGCTGGATCATAAACCCCAATTCTGGTCTCAAGAAACACCGGAAAAAGAAGATTACATGCAATAGCTTCCAC | T | C | 653362338 | 653362438 | 644645703 | 653659283 | 653659383 | 655238591 | 655238591 |
| RDmB-h3:SNP_347 | ATTATCATGATGAGCTGGATCATAAACCCCAATTCTGGTCTGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | G | C | 653362347 | 653364247 | 644645712 | 653659292 | 653659392 | 655238500 | 655238500 |
| RDmB-h3:SNP_348 | ACCTGCAACTGGTGTAGTGTAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | A | C | 653375163 | 653375263 | 644658584 | 644658684 | 644658720 | 655215155 | 655215155 |
| RDmB-h3:SNP_349 | GTAGTACCTGGAGCAGGGAGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGG | T | T | 653375175 | 653375279 | 644658600 | 644658700 | 644658722 | 655215168 | 655215168 |
| RDmB-h3:SNP_350 | CGATCAGTGGCAATCCGGTACAGACCTGGCTGACATGGTAAAGAACAGGAGCAAGCCGGCGAGTGTGCTGAGCAATCTGGAGGAGCTGG | T | C | 653375650 | 653375750 | 644659071 | 644659171 | 653672596 | 655215193 | 655220203 |
| RDmB-h3:SNP_351 | TTCTCCAAAGAAAGAACACTGATCGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGG | A | C | 653398930 | 653399030 | 644682253 | 644682353 | 644682353 | 655275221 | 655275221 |
| RDmB-h3:SNP_352 | AAAGAACACTGATCGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGG | G | C | 653398940 | 653399040 | 644682463 | 644682463 | 644682463 | 655275231 | 655275231 |
| RDmB-h3:SNP_353 | AAACACTGTCATCGAGACTCGCTGAGCACGGGTGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGG | T | C | 653398944 | 653399044 | 644682367 | 644682467 | 644682467 | 655275235 | 655275235 |
| RDmB-h3:SNP_354 | AAACAGGAGAACGAGTGAAGTACGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGG | A | G | 653399885 | 653399885 | 644683308 | 644683604 | 644683604 | 655276176 | 655276176 |
| RDmB-h3:SNP_355 | GAGAGTGAAGTCTGGCGAGGGCTTACCTGGCTGAGTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGG | C | T | 653399895 | 653399995 | 644683418 | 644683681 | 644683681 | 655276285 | 655276285 |
| RDmB-h3:SNP_356 | AGGCCCTGGCGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | G | C | 653400740 | 653400840 | 644684163 | 644684263 | 644684263 | 655277103 | 655277103 |
| RDmB-h3:SNP_357 | GTATTTTCTCTGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | C | A | 653400840 | 653400940 | 644684263 | 644684363 | 644684363 | 655277223 | 655277223 |
| RDmB-h3:SNP_358 | CCCATGACTGATGTTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | G | C | 653400852 | 653400952 | 644684275 | 644684375 | 644684375 | 655277143 | 655277143 |
| RDmB-h3:SNP_359 | AAACACTGTCATCGAGACTCGCTGAGCACGGGTGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | A | T | 653398944 | 653399044 | 644682367 | 644682467 | 644682467 | 655275335 | 655275335 |
| RDmB-h3:SNP_360 | TTCCGGGATGGCTGCTCTGGCTGACTGATCTGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | G | T | 653400865 | 653400965 | 644684288 | 644684388 | 644684388 | 655277225 | 655277225 |
| RDmB-h3:SNP_361 | GGATGGCTGCTGCTCTGGCTGACTGATCTGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGGAGGAGCTGGCTGG | C | G | 653400871 | 653400971 | 644684294 | 644684394 | 644684394 | 655277226 | 655277226 |
| RDmB-h3:SNP_362 | ATAAGGACTAACAACTAACTGAGAACAAAAAACTTACAGTCAGGAGAACACAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | T | C | 653409188 | 653409288 | 644692640 | 644692740 | 644692740 | 655277231 | 655277231 |
| RDmB-h3:SNP_363 | TTGATGACTGCTGCTGCTGACTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | G | A | 653411853 | 653411953 | 644695329 | 644695429 | 644695429 | 655278808 | 655278808 |
| RDmB-h3:SNP_364 | TTGATGACTGCTGCTGCTGACTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | T | C | 653411859 | 653411959 | 644695335 | 644695435 | 644695435 | 655278804 | 655278804 |
| RDmB-h3:SNP_365 | GGACAGGAGATGAGTACCATGACGAGAACGGGGAGTACATGTCAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGGAGGAGCTGG | A | C | 653412403 | 653412503 | 644695879 | 644695979 | 644695979 | 655288563 | 655288563 |
| RDmB-h3:SNP_366 | GGAGGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGGAGGAGCTGG | T | C | 653414648 | 653414784 | 644698160 | 644698260 | 644698260 | 655290944 | 655290944 |
| RDmB-h3:SNP_367 | TTGATGACTGCTGCTGCTGACTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | G | T | 653412501 | 653412601 | 644704977 | 644705077 | 644705077 | 655297657 | 655297657 |
| RDmB-h3:SNP_368 | ACAGGGCTGATTCCTGGCTGACTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | A | G | 653414823 | 653414823 | 644698199 | 644698299 | 644698299 | 655290883 | 655290883 |
| RDmB-h3:SNP_369 | GGGGCCAGATCTGGCTGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | C | A | 653414813 | 653414913 | 644698289 | 644698389 | 644698389 | 655291073 | 655291073 |
| RDmB-h3:SNP_370 | TTGATGACTGCTGCTGCTGACTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | T | C | 653419064 | 653419164 | 644702540 | 644702640 | 644702640 | 653716093 | 653716093 |
| RDmB-h3:SNP_371 | TTGATGACTGCTGCTGCTGACTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | G | T | 653421456 | 653421556 | 644704932 | 644705032 | 644705032 | 655277162 | 655277162 |
| RDmB-h3:SNP_372 | GTGGATTATTCAGAGGAGCTGGCTGCTGACTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGGAGGAGCTGG | C | A | 653421501 | 653421601 | 644704977 | 644705077 | 644705077 | 655297657 | 655297657 |
| RDmB-h3:SNP_373 | TTGATGACTGCTGCTGCTGACTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGGAGGAGCTGG | G | T | 653421609 | 653421709 | 644705085 | 644705185 | 644705185 | 655297864 | 655297864 |
| RDmB-h3:SNP_374 | ATGATGATGAGTATGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGAGGAGCTGGCTGG | T | A | 653421707 | 653421807 | 234399464 | 234399564 | 234399564 | 244586451 | 244586451 |

| | | | | | | | | | | | |
|--------------------------|---|---|---|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| RD M b-h3:SNP_503 | CA GT CC A T T T C T G T C A G G G T C C T A C G A T C A C T G G A T C T T A A G G A G A C T C T T T G C C C A A T C A G G T G A C C A A C A C T G A T G A A G G G G C A T G C A | A | G | 653532368 | 653532468 | 644816073 | 644816173 | 653829128 | 653829228 | 655408805 | 655408905 |
| RD M b-h3:SNP_504 | T C T T G C T T C A T G T C G T T A G T G A T G A T G C T C A A A C C C T C A C T A C T T T C T C A C A T G A C C T A A C C T A A A G C C C A A T C G G T C A C C G A T T C T | T | C | 653534008 | 653534108 | 644809477 | 644809577 | 653822532 | 653822632 | 655402209 | 655402309 |
| RD M b-h3:SNP_505 | T G C T T C A T G T C G T T A G T G A T G C T C C A A A C C C T C A C T A C T T T C T C A C A T G A C C T A A C C T A A A G C C C A A T C G G T C A C C G A T T C T T C | C | T | 653534011 | 653534111 | 644809480 | 644809580 | 653822535 | 653822635 | 655402212 | 655402312 |
| RD M b-h3:SNP_506 | A A A C C C T C A C T A C T T T C A C A T C C A C A T A G C C T A A A C C T A A A G C C C A A T C G G T C A C C C G A T T C T T C T C C G G C C C A C C G A G G T C A G A T G T C A | C | A | 653534039 | 653534139 | 644809508 | 644809608 | 653822563 | 653822663 | 655402240 | 655402340 |
| RD M b-h3:SNP_507 | C C C T C A C T A C T T T C A C A T C C A C A T A G C C T A A A C C T A A A G C C C A A T C G G T C A C C C G A T T C T T C T C A T C C G G C C C A C C G A G G T C A G A T G T C A T A G | A | C | 653534042 | 653534142 | 644809511 | 644809611 | 653822566 | 653822666 | 655402243 | 655402343 |
| RD M b-h3:SNP_508 | C A C A A T G G T A T C T C G C T A G C C C T T C T G A G A C C C G C C C A A A C T A A A T G C A T G A A C C C T T A A A G G A T C A A G G A C T G A C T A A A C A C T T G T A A A G G | T | G | 653535699 | 653535799 | 644819406 | 644819506 | 653823459 | 653823559 | 655412137 | 655412237 |
| RD M b-h3:SNP_509 | A T T A A T G C A T A G A C C T T A A A G G A T C A A G G A C T G A C T A A A C A T T G T A A A G G A T C A G G T A C G T A C A T C C C A A T A A A C C C A A T A A T A A T G A A T | C | T | 653535739 | 653535839 | 644819446 | 644819546 | 653823499 | 653823599 | 655412177 | 655412277 |
| RD M b-h3:SNP_510 | A T G C A T A G A A C C T T A A A G A T C A A G G A C T G A C T A A A C A T T G T A A A G G A G A T C C A G T C A G G T A C A T C C C A A T A T A A A C C A A T A A T A A T G A A T | T | G | 653535743 | 653535843 | 644819450 | 644819550 | 653823503 | 653823603 | 655412181 | 655412281 |
| RD M b-h3:SNP_511 | T C C T C T G A G A G A T A A A C C C A C T A C T G C A T G G A A G A T A T T G G C A G G C T T A G T G C A G G C T G A C A C C A A C A A C A G A T T C A T T G A A G G I T T A G | G | T | 653536277 | 653536377 | 644819984 | 644820084 | 653823037 | 653823137 | 655412715 | 655412815 |
| RD M b-h3:SNP_512 | T A A C C C A C T T A G C A T G G A A G A T A T T G G C A G G C C T A G T G G A A C A T G A G C A G C T G G A C A T A C A A A C A G A A T T C A T T G A A G G I T T A G A G G I T T A G G C A C A T | A | T | 653536289 | 653536389 | 644819996 | 644820096 | 653823049 | 653823149 | 655412727 | 655412827 |
| RD M b-h3:SNP_513 | A A A T C C C C A A C A T T A A G A A T A T C A T A T T T T C T G C A G T G A G G G A G C G A A T T C A A A C C T C C A A T A A T A T T G A T G G A A T T T C T C A A T G A T T I T T A | C | G | 653539349 | 653539449 | 644820306 | 644821356 | 653826109 | 653826209 | 655415787 | 655415887 |
| RD M b-h3:SNP_514 | A C G T C C T C A C A A T G C C G A C A G G T A T A G C A G T T G A T T T A T C A G C C A T T A C A A G G A T A T C A G T G G G G T C A C T A C T C A A G T C T C A C G G T A C A A | A | G | 653539659 | 653539759 | 644823366 | 644823466 | 653836419 | 653836519 | 655416097 | 655416197 |
| RD M b-h3:SNP_515 | A G C C T C C C C C A A G C T T G T G G G A G A C T C T C T T G A T T G C A C A A A C A T A T T T C C C A A A G G C A G C T T G T T C A T A T G A G G C A G G G G G A A T T T A G C | C | T | 653541264 | 653541364 | 644824971 | 644825071 | 653838024 | 653838124 | 655417702 | 655417802 |
| RD M b-h3:SNP_516 | T G G G C T C T C T G A G A G A T A A A C C C C T A G T C C T G C T G G G G T C A T C A T G A T C A T A A T C A T G T G A T G G T C A T C A G G G C T C T G G C T G T G T T T | T | C | 653544423 | 653544523 | 644828131 | 644828133 | 653841183 | 653841283 | 655420862 | 655420962 |
| RD M b-h3:SNP_517 | T C A T G T T G G C T T A A T T C A T G T C G G T A C T T T G A T A T T G A T T G T G G G A G G A C C G G T G C C T T A G T G A A C A A C C T C C A A C T T A T G A T | A | T | 653548424 | 653548524 | 644832132 | 644832232 | 653845184 | 653845284 | 655424863 | 655424963 |
| RD M b-h3:SNP_518 | T C G G T A C T T G A T A T T G A T A T T G T G G G A G G A C C G G T G C T T A G T G A C C A A C C T C C A A C T T A T G A T T A T C T C T C G A A G C A T C C G | C | T | 653548445 | 653548545 | 644832153 | 644832253 | 653845205 | 653845305 | 655424884 | 655424984 |
| RD M b-h3:SNP_519 | C A A C A T A C A T A T C A T A A A A A T C G A A T A C A T A C A G G T A C A T G A C T C A T C A C T G A C A T T C T C T C G G A C C A A G A A G A A G A A G A T C A C T C A C | C | T | 653548761 | 653548861 | 644832469 | 644832569 | 653845521 | 653845621 | 655425200 | 655425300 |
| RD M b-h3:SNP_520 | C T T G G G G T C G A A A C T T G G A A G T T A A A G G C T C C A A A G G C T C C A G G T C G G T T C C T C G G A C T C C C C C G C C A G C C C G A G G A C C T G G | T | C | 653573984 | 653574084 | 644857883 | 644857983 | 653870627 | 31388458 | 31388358 | |
| RD M b-h3:SNP_521 | T A T T G A T T C A A G C C A T T G A T A A A C C C A G A C A T T C T G A C G A C A T G C T C A T T G G T G A G G T A T T A T C C C A T C T I T G T A G G C A A A G G A T A C T T I G T C A G A G G C C | T | C | 653576912 | 653577012 | 644860811 | 644860911 | 653873455 | 653873555 | 655453296 | 655453396 |
| RD M b-h3:SNP_522 | T C G C A C C G A G G G C A G C C G G A C T T C G A T C T C C G C C G A C T G C T G A C T G G T G A G G T A C T T C G T G C T G A C T T G A T T C T C T A A T T T A | A | G | 653588550 | 653588650 | 644872450 | 644872550 | 653885093 | 653885193 | 655464935 | 655465035 |
| RD M b-h3:SNP_523 | T T G A A T T C G C T T G A C C T T C T C C C C T C T C G G T C A G A T C C A C C A C C C G A G G C C G A G G C C G C C G C G G C G C G G C A | A | G | 653589182 | 653589282 | 644873082 | 644873182 | 653885723 | 653885723 | 655465567 | 655465667 |
| RD M b-h3:SNP_524 | C A A G G C T G C C C C G A G A G C G G G C G A G G G A T T C C G C G C A C T G C A C T G C A G C C C G A G G C C G C A G G C C G C A C C C G C C A T C C C A | A | G | 653589350 | 653589450 | 644873250 | 644873350 | 653885893 | 653885993 | 655465735 | 655465835 |
| RD M b-h3:SNP_525 | T A G C A G C T C A A G C G G T C C T C A A G G A G C A C G C G C G T A C T T G T C G T A G G T C G T C T C T G C T C A A G G C A G C A C T G G A T G C G T A T | T | C | 653591800 | 653591900 | 644875699 | 644875799 | 653888343 | 653888443 | 655468184 | 655468284 |
| RD M b-h3:SNP_526 | A A C T G C A C G C T C C A C A C T T G G T A C T T G T G A C G T G T C A T C T G C C G A A A C C T T G C C G A A A C C G C T T G C C G A A G T G C T C T G A | T | C | 653592459 | 653592559 | 644876358 | 644876458 | 653889002 | 653889102 | 655468843 | 655468943 |
| RD M b-h3:SNP_527 | A C G C A T T G G C T G A T C T G T G C T G G A C C C C C G A G A T T C C A C C G T C A T G T G C C A C G A A C G T A T G G T A C T G C C G T C G C G T | T | C | 653597552 | 653597652 | 288843867 | 288843967 | 653894095 | 653894195 | 299619477 | 299619577 |
| RD M b-h3:SNP_528 | A G T C G A T T C A C T C A G C A T C G G T G A A A G G C A T G C C T T T A C C A G C G A T G A A G A G G C T G A C A G G T C A G T G G A C T C C A G A G C T G A T G C A G C G T | G | T | 653618097 | 653618197 | 644902184 | 644902284 | 653914286 | 653914386 | 655493951 | 655494051 |
| RD M b-h3:SNP_529 | C A G C T G A T G A T C A A A G A T A G G T C C A A G A G A A G G C A T T C C G A G G G G T G T G T T C T G G A G C A T A C A C A C T C C G T C A G T C C A A G A G A T C T T C | G | T | 653619025 | 653619125 | 644903112 | 644903212 | 653915214 | 653915314 | 655494879 | 655494979 |
| RD M b-h3:SNP_530 | G A T C A T G G A G G T T T G G T G C A T T C C C G A C T G G A C G A T T C C G T A T G T C A C A T A T A C T C T G A G G A A A G G C C T C C G G A G G A T G A C A G A | A | G | 653619197 | 653619297 | 644903284 | 644903384 | 653915386 | 653915486 | 655495051 | 655495151 |
| RD M b-h3:SNP_531 | T G T C A C A T A T C T G A G G A A A G G C T T C C C G G A G G T G A C A G A T C G T C C A T G A T C A T A A G G C C T T T C C C T A A C A A G G A C A T T A T A T | T | C | 653619242 | 653619342 | 644903329 | 644903429 | 653915431 | 653915531 | 655495096 | 655495196 |
| RD M b-h3:SNP_532 | A C G C A T G C A G C T T T A G A G G A A G A A G G G A G A T G G C T C T G A T C C C G T C A C C A T T C A C C A A G A T T T G C C G C T T C C C G C A G G A A T T G T G A G A G G T | A | G | 653620165 | 653620265 | 586074136 | 586074036 | 653916354 | 653916454 | 655496019 | 655496119 |
| RD M b-h3:SNP_533 | C C C C T A C C A C G T C G G T G A A G G G C T G G A C A A G C T C C G C C T T C T G A C C G C C T T C C T C T C T C T C T C G C A C G T | C | T | 653624231 | 653624331 | 644907962 | 644908062 | 653919637 | 653919737 | 655500059 | 655500159 |
| RD M b-h3:SNP_534 | T G T C G G C C A A T T G G A T C A G C G G T G A A G G G C T G G A C C C C T A G A G G G A A G T T T A T C T G A A T A G C C G T G T C C C T C A A A A A G G A C A G C G | G | C | 653639329 | 653639429 | 644922838 | 644922938 | 653934450 | 653934550 | 65514588 | 65514688 |
| RD M b-h3:SNP_535 | T C G A G C G C T G T G G G T T A T G T G G T G A C C C C T A G A G G G A A G T T T A T C T G A A T A G C C G T G T C C C T A A A A A G G A C G A C G G G A G G T G T A C C T G A C C | C | T | 653639346 | 653639446 | 644922855 | 644922955 | 653934567 | 653934605 | 655514705 | |
| RD M b-h3:SNP_536 | C G A G C G T G T G G G T T A T G T G G T G A C C C C T A G A G G G A A G T T T A T C T G A A T A G C C G T G T C C C T C A A A A A G G A C G A C G C G G A G G T G T A C C T G A C C T | T | C | 653639347 | 653639447 | 644922856 | 644922956 | 653934468 | 653934568 | 655514606 | 655514706 |
| RD M b-h3:SNP_537 | G T G C T T A G T A A G A A T T T G T A G T A A G G T T T G T A T T C C C T A T G C A A G A C G A A A G T C A A T G A T A T G C A A T G A A A T T T A T A T C T C A C T T A T G G T G C A T T A T C A G T G T A C T T A | A | C | 653641912 | 653642012 | 644925421 | 644925521 | 653937033 | 653937133 | 655517171 | 655517271 |
| RD M b-h3:SNP_538 | T A T T G A T G T A A G G G T T G T A T T C C C T A T G C A A G A C G A A A G T C A A T G A T A T G C A A T G A A A T T T A T A T C C T A C T T A T G G T G C A T T A T C A G T G T A C T T A | A | T | 653641927 | 653642027 | 644925436 | 644925536 | 653937048 | 653937148 | 655517186 | 655517286 |

7. DECLARATION OF AUTHORSHIP

I hereby confirm that this thesis and the work presented in it is entirely my own. Where I have consulted the work of others this is always clearly stated. All statements taken literally from other writings or referred to by analogy are marked and the source is always given. This paper has not yet been submitted to another examination office, either in the same or similar form. I agree that the present work may be verified with an anti-plagiarism software.

Place, Date

Original Signature

ACKNOWLEDGEMENT

This work is dedicated to L. for your continuous support and advice and for listening to my worries.

I would like to thank Dr. Makhoul heartfully for guiding me throughout this project and giving me inspiration to accomplish my goals.

I want to express my gratitude to Dr. Obermeier for the help provided to target my research topic and for his explanations and suggestions during the correction of the report.

I also would like to show my appreciation to Prof. Dr. Snowdon for giving me the opportunity to challenge myself and improve my knowledge and skills with this project.

Last but not least, thanks to Dr. Hao and colleagues for their extremely useful collaboration.