1  **A catalogue of resistance gene homologs and a chromosome-scale reference**
2  **sequence support resistance gene mapping in winter wheat**
3
4  Sandip M. Kale[1]*, Albert W. Schulthess[1]*, Sudharsan Padmarasu[1], Philipp H. G. Boeven[2],
5  Johannes Schacht[2], Axel Himmelbach[1], Burkhard Steuernagel[3], Brande B. H. Wulff[3,4], Jochen
6  C. Reif[1], Nils Stein[1,5#], Martin Mascher[1,6#]
7
8  [1]Leibniz Institute of Plant Genetics and Crop Plant Research (IPK) Gatersleben, Seeland,
9  Germany
10  [2]Limagrain GmbH, Peine-Rosenthal, Germany
11  [3]John Innes Centre, Norwich Research Park, Norwich, UK
12  [4]Center for Desert Agriculture, Biological and Environmental Science and Engineering
13  Division (BESE), King Abdullah University of Science and Technology (KAUST), Thuwal, Saudi
14  Arabia.
15  [5]Center for Integrated Breeding Research (CiBreed), Georg-August-University Göttingen,
16  Göttingen, Germany
17  [6]German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Leipzig,
18  Germany
19
20  *These authors contributed equally.
21  #Correspondence should be addressed to Nils Stein (stein@ipk-gatersleben.de) or Martin
22  Mascher (mascher@ipk-gatersleben.de)
23
24
25
26
27  **Abstract**
28
29  A resistance gene atlas is an integral component of the breeder's arsenal in the fight against
30  evolving pathogens. Thanks to high-throughput sequencing, catalogues of resistance genes
31  can be assembled even in crop species with large and polyploid genomes. Here, we report
32  on capture sequencing and assembly of resistance gene homologs in a diversity panel of 907
33  winter wheat genotypes comprising *ex situ* genebank accessions and current elite cultivars.
34  In addition, we use accurate long-read sequencing and chromosome conformation capture
35  sequencing to construct a chromosome-scale genome sequence assembly of cv. Attraktion,
36  an elite variety representative of European winter wheat. We illustrate the value of our
37  resource for breeders and geneticists by (i) comparing the resistance gene complements in
38  plant genetic resources and elite varieties and (ii) conducting genome-wide associations
39  scans (GWAS) for the fungal diseases yellow rust and leaf rust using reference-based and
40  reference-free GWAS approaches. The gene content under GWAS peaks was scrutinized in
41  the assembly of cv. Attraktion.
42

## Introduction

Maintaining plant health in the face of evolving pathogen populations is a perennial goal of breeders. Key to this endeavor is the discovery and deployment of disease resistance (R) genes. Hafeez et al. (2021) put forward the concept of an R gene atlas and illustrated its potential for crop improvement in one of our most widely grown crops, wheat. An important component of populating the wheat R gene atlas is genotyping diversity panels, or more broadly, knowledge of as large a fraction of the resistance gene complement of as many genotypes as possible. One approach to this aim, resistance gene enrichment sequencing [RenSeq, Jupe et al. (2013)], was developed with large-crop genomes in mind. To reduce genomic complexity of sequencing libraries, and hence the required sequence effort, capture probes are designed to target R gene homologs from the nucleotide-binding and leucine-rich repeat (NB-LRR) family, or more generally the family of NB-LRR-related genes (NLRs, Ting et al. (2008)). In its original implementation, RenSeq was combined with short-read sequencing on the Illumina platform (Jupe et al., 2013). A combination of RenSeq with long-read sequencing has been used to assemble the full complement of R genes in the model plant *Arabidopsis thaliana* and analyze their evolutionary dynamics (Van de Weyer et al., 2019).

RenSeq data for diversity panels in combination with matching phenotype data has been used for genome-wide associations scans (GWAS) to find genetic markers associated with disease resistance (Arora et al., 2019). In the best case, this method, termed AgRenSeq, can zoom in on individual candidate genes. However, the limits of association mapping such as population structure (Yu et al., 2006) and sensitivity to the genetic architecture of the trait under study (Lopez-Arboleda et al., 2021) also apply to AgRenSeq. Recently, Gaurav et al. (2021) reported the use of whole-genome shotgun sequencing for association mapping of disease resistance in the wheat diploid progenitor *Aegilops tauschii*. An advantage of WGS over RenSeq is its ability to access also non-NLR resistance genes; a potential drawback is the inability to assemble full-length genes from low to medium-coverage (3x-10x) short-read data.

Independent of choice of sequencing strategy, a potential impediment to GWAS as well as a crucial aspect of R gene evolution is structural variation (SV). R genes are subject to ubiquitous presence-absence and copy-number variations (Michelmore and Meyers, 1998; Van de Weyer et al., 2019). Reference-free GWAS approaches have shown that the presence of peaks can be influenced by the choice of reference sequence (Voichek and Weigel, 2020). In principle, the best resource for studying intra-species NLR diversity are high-quality genome assemblies for a representative diversity panel comprising hundreds of accessions, i.e. a pan-genome. Constructing pan-genome infrastructures for all major crops has recently turned from a moon shot into a realistic mid-term research goal (Della Coletta et al., 2021). But the wheat pan-genome is not there yet: chromosome-scale reference genome sequences for ten wheat varieties, most of them recent elite cultivars, have recently been released (Walkowiak et al., 2020), but this small panel is not comprehensive enough to underpin a species-wide resistance gene inventory.

In the present manuscript, we report on a contribution to the wheat R gene atlas. We constructed an R gene inventory for a diversity panel of winter wheat, the predominant type of wheat in Europe. A chromosome-scale reference genome sequence was constructed for one representative winter wheat cultivar. To illustrate the value of this resource for the wheat genetics and breeding community, we (i) compare patterns of R gene diversity between plant genetic resources (PGR) and elite cultivars; (ii) conduct GWAS for the fungal

91  diseases yellow rust and leaf rust; and (iii) analyze structural variants in close proximity to
92  significantly associated markers.
93
94  **Results**
95
96  *R gene capture in a winter wheat diversity panel*
97
98  We conducted RenSeq for a panel of 907 winter wheat genotypes (**Figure 1**, **Table S1**) and
99  the reference genotypes Chinese Spring (The International Wheat Genome Sequencing
100  Consortium (IWGSC), 2018) and Julius (Walkowiak et al., 2020). Of these, 779 are part of a
101  previously described core set enriched for disease resistant genotypes (Schulthess et al.,
102  2021) comprising 587 PGRs and 192 European elite cultivars. The remaining 128 genotypes
103  are recent German elite breeding lines. We used the Triticeae RenSeq Baits V3 (Tv3) probe
104  set comprising 217,827 oligonucleotide baits (Zhang et al., 2021). Alignment of this bait set
105  to the Chinese Spring reference genome (RefSeq v1.0, The International Wheat Genome
106  Sequencing Consortium (IWGSC) (2018)) indicated that 18 Mb of annotated NBS-LRR gene
107  sequence are targeted. On average, sequences originating from the predicted target were
108  enriched 220-fold.
109  RenSeq reads of individual genotypes were assembled *de novo*, yielding 67,731 to 4,583,579
110  contigs per accession (**Table S2**). Of these, 417 to 2,304 per accession (mean: 1,690)
111  contained full-length NLRs. Coiled-coil NLRs were the most abundant class of NLRs (**Figure
112  2**). Almost all (1,911/1,937) NLRs assembled from the Chinese Spring RenSeq data were
113  aligned to RefSeq v1.0 with 95 % alignment coverage and 95 % alignment identity. A total of
114  1,486 (77 %) Chinese Spring *de novo* assembled NLRs overlapped with RefSeq v1.0 gene
115  models, indicating the absence of resistance gene homologs in the reference annotation,
116  possibly because of lack of expression or pseudogenization. In other genotypes, on average
117  77 % of assembled NLRs were mapped to RefSeq v1.0, consistent with pervasive presence-
118  absence variation (PAV) in resistance genes.
119
120  *Diversity of NLRs in winter wheat genepools*
121
122  To understand the extent of PAV in NLRs in our panel, we performed similarity-based
123  clustering of the assembled NLR from all genotypes. A total of 1,469,694 (96 %) of NLRs were
124  clustered in 39,073 orthogroups, the remainder were singletons without close matches to
125  other NLRs. Fewer than 1 % of orthogroups contained two or more NLRs from the same
126  accession, pointing to a potential collapse of highly similar, recently duplicated NLRs. Most (>
127  85 %) of orthogroups had NLRs from at least 20 different accessions. However, very few
128  orthogroups (472, 1.21%) had members from more than 500 accessions (**Figure S1**). This is at
129  odds with patterns of NLR diversity in the model plant *Arabidopsis thaliana* (Van de Weyer
130  et al., 2019), where the "core-NLRome" comprising genes present in almost all genotypes is
131  substantial. The likely explanation is random sequence dropout due to competition between
132  capture probes and/or low sequencing depth. For example, at a 1 % dropout rate (i.e., a 99
133  % change of being captured and sequenced at sufficient depth), a gene present in 900
134  genotypes has a negligible chance ($0.99^{900}$ = 0.01 %) of being present in all their assemblies.
135  A saturation analysis indicated that a near-complete set of NLR orthogroups assembled in
136  the whole panel can be captured with a rather small number of accessions: 95 % of
137  orthogroups were captured with only 70 genotypes selected at random from the universe of
138  907 accessions (**Figure 3**). Because of random dropout, these figures are likely

3

139    overestimates, i.e. an even smaller panel may suffice to reach the 95 % threshold. Still, the
140    analysis of orthogroups allows comparison of relative diversity between gene pools. When
141    considering PGR and elite accessions separately, near-saturation can be achieved with 80
142    and 150 accessions, respectively, indicating that, not unexpectedly, NLR diversity is higher in
143    PGRs. However, elite lines were more resistant against yellow rust compared to PGRs and
144    contained a higher number of NLRs that preferentially occur in resistant genotypes,
145    supporting the notion that breeder's efforts to stack resistance genes have been successful
146    (**Figure 4a**). A potential caveat, though, is that NLRs private to elite varieties were localized
147    to regions previously reported as harboring alien introgressions. By contrast, PGR-specific
148    NLRs were distributed more uniformly across the chromosomes (**Figure 4b**). Hence, most
149    NLRs occurring in highly resistant elite varieties may not confer resistance on their own, but
150    have only hitchhiked along one or a few functional resistance genes targeted by breeders.
151
152    *A chromosome-scale assembly of cv. Attraktion*
153
154    Several recent studies suggest that the choice of reference genome impacts the
155    contextualization or even the very presence of GWAS peaks (Arora et al., 2019; Voichek and
156    Weigel, 2020). It is likely that a better reference genome than the assembly of Chinese
157    Spring – indeed a spring-sown landrace from China (Sears and Miller, 1985) – can be
158    selected for mapping resistance genes in winter types. We chose cv. Attraktion because our
159    prior analysis of shallow-coverage whole-genome shotgun data had shown that this cultivar
160    carries large alien introgressions, some of them co-incident with GWAS peaks for resistance
161    to yellow and leaf rust (Schulthess et al., 2021).
162    We sequenced the Attraktion genome to 22-fold coverage with HiFi reads with an average
163    length of 17.8 kb of circular consensus reads. In addition, chromosome conformation
164    capture sequencing (Hi-C) was performed, resulting in 994 million read pairs. Genome
165    assembly following a previously described approach (Mascher et al., 2021; Sato et al., 2021)
166    combining primary contig assembly with Hifiasm (Cheng et al., 2021) and pseudomolecule
167    construction with the TRITEX pipeline (Monat et al., 2019) yielded a set of 1,553 contigs
168    (14.25 Gb) assigned to chromosomal locations. A further 3,442 contigs (434 Mb) remained
169    unplaced. A BUSCO analysis (Simao et al., 2015) indicated that 98.2 % of conserved single-
170    copy genes were present in the assembly (**Table 1**). The inspection of Hi-C contact matrices
171    and alignment to the Chinese Spring RefSeq v2.1 (Zhu et al., 2021b) supported the structural
172    integrity of the pseudomolecules (**Figures S2, S3**).
173    Regions of high divergence between Attraktion and Chinese Spring indicative of the
174    presence of alien introgressions were found on four chromosomes: 4A, 2B, 5B and 2D
175    (**Figure 5**). Interestingly, a 55 Mb introgression on the long arm of chromosome 2B in
176    Attraktion overlapped with a much larger 427 Mb introgression from *T. timopheevi* in
177    LongReach Lancer. Attraktion and Lancer have the same haplotype in the overlapping
178    region, pointing to shared ancestry. Most likely, breeders had decreased the size of this
179    introgression in Attraktion in an attempt to reduce linkage drag.
180
181    *Different GWAS approaches identify a yellow rust resistance locus on chromosome 6A*
182
183    To illustrate the value of our resource for genetic mapping of disease resistance, we
184    conducted GWAS for yellow rust (*Puccinia striiformis* f.sp. *tritici*) resistance in our panel. The
185    degree of yellow rust infection was scored in multi-environment field trials, relying on
186    natural and artificial infection. Details are described elsewhere (Schulthess et al., 2021). We

4

187   followed three different approaches to obtain matrices of bi-allelic markers for use in GWAS.
188   First, we aligned RenSeq reads to a reference genome sequence assembly, called single-
189   nucleotide polymorphisms (SNPs), and used genotype calls at SNP sites as markers. This is
190   the most commonly applied approach for marker discovery, which, however, can capture
191   structural variants only if they are in close linkage disequilibrium (LD) with SNPs. Second, we
192   conducted kmerGWAS (Voichek and Weigel, 2020) which queries the presence-absence
193   state of short oligonucleotides of a fixed length (*k*-mers, *k*=31) as proxies for structural
194   variants. Third, we used SNP sites discovered from the alignment to the reference assembly,
195   but instead of allelic status, we used presence-absence states of genotype calls as markers,
196   similar to what Gabur et al. (2018) did with SNP chip data of rapeseed. We refer to this
197   method as paGWAS. Two different reference genome sequences, Chinese Spring RefSeq
198   V2.1 and our Attraktion assembly, were used to position markers. Note that kmerGWAS is a
199   reference-free approach; associated *k*-mers were aligned *post hoc* to the genome assemblies
200   to place them.
201   Manhattan plots for all three methods and the two references are shown in **Figure 6**. The
202   most prominent feature is a peak on the long arm of chromosome 6A, for which significantly
203   associated markers were reported in GWAS scans. However, it is less prominent in paGWAS
204   and kmerGWAS against the Chinese Spring reference, possibly reflecting the absence of the
205   resistant haplotype in that genotype. SNP GWAS against Chinese Spring did result in a
206   pronounced peak, likely because of SNPs in linkage disequilibrium with the causal variant.
207   Further peaks were observed on other chromosomes, but were not common between all
208   methods.
209   To the best of our knowledge, a resistance gene against yellow rust has not been reported
210   on chromosome 6A in the region pinpointed by our GWAS. We scrutinized the region under
211   the peak in the genome assembly of cv. Attraktion, which scored highly in our resistance
212   trials and carried the resistant haplotype at the 6A peak. The significantly associated markers
213   spanned an interval of 946 kb in the Attraktion genome (**Figure 7a**), containing 121 gene
214   models annotated *ab initio* (Stanke et al., 2006), many which are actually derived from
215   transposable elements. Seven of these genes were NLRs. One of them, spanning a 4.9 kb at
216   around sequence coordinate 612.5 Mb on the 6A pseudomolecule of Attraktion, was
217   identical to the representative contig of the orthogroup "cluster49707" identified from the
218   RenSeq *de novo* assemblies. This representative contig harbored 46.9 % of the significantly
219   associated 31-mers. Among the 10 genome sequence assemblies reported by Walkowiak et
220   al. (2020), that of SY Mattis had the same haplotype as Attraktion (**Figure 7a**). The other nine
221   genomes as well as the Chinese Spring reference lacked several genes present in the
222   Attraktion haplotype, including cluster49707.
223   Interestingly, no significant marker-trait associations were detected in the peak region when
224   only PGRs were included in the association scan (**Figure S4**). The high synteny between
225   Chinese Spring and Attraktion in the vicinity of the peak rules out the presence of an alien
226   introgression. We speculate that the resistant haplotype may have segregated at low
227   frequency – below the detection threshold of GWAS – in landraces and old varieties. Recent
228   shifts in European pathogen populations (Hovmøller et al., 2016) may have favored the
229   resistant haplotype in European winter wheats (**Figure 7b**). Future work should focus on the
230   identification and validation of the causal gene conferring yellow rust resistance and on
231   singling out the yellow rust isolates that is recognizes.
232
233   *GWAS for leaf rust detects known and novel loci*
234

235 The second trait for which we did GWAS is leaf rust (*Puccinia triticina* f.sp. *tritici*) resistance.
236 The degree of natural infection with leaf rust was scored under field conditions (**Tables S3**
237 **and S4**). SNP, paGWAS and kmerGWAS gave partially overlapping results (**Figure 8a**,
238 **Supplementary Figure S5**). Common to all approaches was a peak towards the distal end of
239 the long arm of chromosome 4A. This association had been reported before by Liu et al.
240 (2020a), who analyzed 133 genotypes and 1,574 of their hybrid offspring by exome
241 sequencing. The GWAS peak was co-located with a 26 Mb region of high sequence
242 divergence between Chinese Spring and Attraktion (**Figure 5**), which we attribute to an alien
243 introgression, possibly originating from *T. dicoccoides* (Przewieslik-Allen et al., 2021).
244 Because of suppressed recombination, the introgression is inherited as one large linkage
245 block (**Figure 8b**).
246 Significant associations on the long arm of chromosome 5B were detected by multiple GWAS
247 approaches. Cultivar Attraktion had high resistance scores and had a haplotype associated
248 with resistance in the peak region, which extended from about 692.4 Mb to 694.1 Mb in the
249 Attraktion genome assembly, spanning 235 gene models. Of these, thirteen were NLRs.
250 Significantly associated *k*-mers mapped to 15 orthogroups of NLRs *de novo* assembled from
251 our RenSeq data, which corresponded to three gene models in the Attraktion assembly. One
252 of them was highly similar to the cloned *Lr21* resistance gene (Huang et al., 2003).
253 Interestingly, none of the wheat pangenome assemblies (Walkowiak et al., 2020) harbored
254 this gene, illustrating the need for bespoke genome sequence assemblies to capture the
255 gene content of resistance gene clusters. Further work is needed to prove the causal link
256 between resistance to leaf rust and presence of the *Lr21*-related gene or other NLRs under
257 the GWAS peak.
258
259 **Discussion**
260
261 We have reported RenSeq assemblies as a component of the burgeoning wheat R gene atlas.
262 Due to short-comings of our short-read capture sequencing approach, we were unable to
263 construct a comprehensive NLR-ome as was done with long-reads in *A. thaliana* (Van de
264 Weyer et al., 2019). However, our data set did reveal a contrasting repertoire of R gene
265 homologs in elite varieties and PGRs, documenting breeders' efforts at enhancing genetic
266 resistance by selecting haplotypes bearing R genes.
267 The assembly of one recent elite cultivar, Attraktion, proved instrumental in the analysis of
268 the gene content surrounding two GWAS peaks on chromosomes 6A and 5B for yellow rust
269 and leaf rust, respectively. But no single genotype can capture all resistance genes and
270 sequence assemblies of other genotypes will be required to zoom in on candidate genes
271 against other diseases or other isolates of yellow rust. Fortunately, the cost for wheat
272 whole-genome assembly has decreased substantially in recent years. The Attraktion
273 assembly was completed within three months after selection of that genotype and cost
274 approximately EUR 40,000. This shows that whole-genome assembly still entails large
275 expenses, which, however, may constitute a worthwhile investment if candidate genes
276 cannot be pinned down by cheaper alternatives. As long as capture and selective sequencing
277 of large (> 500 kb) genomic regions has not become routine (López-Girona et al., 2020),
278 whole-genome assembly is a viable alternative even if only a single locus is of interest.
279 Reference genome sequences of diversity panels large enough for GWAS (i.e. 100-1000
280 genotypes) would render both RenSeq and WGS superfluous: the comprehensiveness and
281 context afforded by genome assembly cannot be matched by short-read approaches.
282 However, the large size of the wheat genome (15 Gb) makes the assembly of hundreds or

283 thousands of genotypes cost-prohibitive at the time of writing. Consequently, the resources
284 reported in this article will retain their usefulness in the medium term.
285 The three different GWAS approaches we took, SNP GWAS, paGWAS and kmerGWAS, were
286 only partially concordant, highlighting the potential benefits that may be reaped from an
287 integration of the GWAS and pan-genomics toolkits. Computational frameworks to construct
288 and analyze pan-genome graphs are under active development (Hickey et al., 2020; Li et al.,
289 2020). Reduced-representation approaches focusing on the single-copy or repeat-depleted
290 part of the genome have been applied in soybean (Liu et al., 2020b) and barley (Jayakodi et
291 al., 2020). It is unlikely, however, that single-copy sequence can represent copy-number
292 variation in rapidly evolving resistance genes. Future algorithmic work should focus on the
293 graph-based representation of pan-genomes for complex plant genomics, graph-based read
294 mapping and GWAS with multi-allelic structural variants captured in pangenome graphs.
295
296 **Methods**
297
298 *DNA extraction, library preparation and sequencing for RenSeq*
299
300 A total of 779 genotypes from a trait-customized core collection of winter wheat (Schulthess
301 et al., 2021) along with 128 advanced elite lines were used for resistance gene enrichment
302 sequencing (RenSeq). DNA was extracted from a single leaf of about 10 cm length harvested
303 from a 10 days old seedling using the DNeasy 96 Plant Kit (QIAGEN, Hilden, Germany) as per
304 the manufacturer's instructions. DNA quality and quantity were determined using a 0.8%
305 agarose gel and Qubit fluorometer (Life Technologies, CA, USA). The RenSeq libraries were
306 prepared using the protocol of Steuernagel et al. (2017) with minor technical modifications.
307 Briefly, 1 µg DNA from each genotype was fragmented to ~500 bp size using the Covaris S2
308 (Covaris, MA, USA). The fragmented DNA was purified using 0.6X AMPure® XP beads
309 (Beckman Coulter, IN, USA) according to the manufacturer's instructions. Paired end libraries
310 for Illumina sequencing were constructed using NEBNext® Ultra™ II DNA Library Prep Kit for
311 Illumina® (New England Biolabs Inc, MA, USA) as per the manufacturer's instructions, except
312 AMPure® XP beads were used for all the purification and size selection steps. For PCR
313 amplification, 10 µl of adapter ligated DNA from each genotype was used along with 25 µl 2x
314 KAPA HiFi HotStart ReadyMix (Kapa Biosystems, MA, USA), 1 µl Index and Universal PCR
315 Primer and 13 µl water. The library from each genotype was indexed using Unique Dual
316 Index Primer Pairs (NEBNext Multiplex Oligos for Illumina) in order to perform multiplexed
317 sequencing.
318 The enrichment of NBS-LRR DNA fragments was achieved through hybridization of PCR
319 amplified genomic DNA libraries prepared above with 200K Triticeae NLR bait libraries (Tv3,
320 Zhang et al. (2021)) available at
321 https://github.com/steuernb/MutantHunter/blob/master/Triticea_RenSeq_Baits_V3.fasta.g
322 z. The libraries were quantified using Qubit fluorometer (Life Technologies, CA, USA) and
323 average fragment size was determined using the 4200 Tape Station (Agilent Technologies,
324 CA, USA). The libraries from eight genotypes were pooled in an equimolar manner and
325 hybridized with the bait library (myBaits-11; Arbor Biosciences, Ann Arbor, MI, USA). The
326 hybridization reaction was carried out at 65 °C for 18 hours and the hybridized fragments
327 were captured using MyOne Streptavidin C1 magnetic beads (ThermoFisher Scientific). The
328 hybridization and capture of NBS-LRR fragments was performed according to MYbaits v4.0
329 protocol. Finally, PCR amplification of captured fragments was carried out using 2x KAPA HiFi
330 HotStart ReadyMix and standard Illumina P5 and P7 primers. Twelve capture libraries (96

7

331   genotypes) were pooled in equimolar amounts, quantified using qPCR and sequenced
332   (paired-end, 2 x 250 cycles) on the NovaSeq 6000 (Illumina).
333
334   *DNA extraction, library preparation and sequencing for PacBio HiFi sequencing*
335
336   High molecular weight (HMW) DNA for PacBio circular consensus sequencing (CCS) was
337   prepared from 100 one-week old seedlings of cultivar Attraktion following the protocol from
338   Dvorak et al. (1988). Briefly, nuclei were extracted from ground leaves in a sucrose-based
339   homogenization buffer. The protein contamination was removed by proteinase-K treatment
340   and phenol:chloroform extraction. The HMW DNA was then spooled out of solution during
341   sodium acetate and ethanol precipitation. Size profile of the extracted DNA was checked
342   using Femtopulse system genomic DNA 165 kb kit (Agilent Technologies, CA, USA). Eight HiFi
343   SMRTbell® libraries were prepared using the SMRTbell$^{TM}$ Express Template Prep Kit 2.0
344   according to the manufacturer's instructions (Pacific Biosciences protocol: PN 101-853-100
345   Version 03, January 2020). In short, the protocol involves fragmenting the HMW DNA to a
346   mean fragment length of 20 kb using Megaruptor 3 (Diagenode), followed by DNA damage
347   repair, end repair/A-tailing and adapter ligation. Linear DNA fragments were removed by
348   nuclease treatment of the SMRTbell libraries. Size selection of libraries was carried out using
349   the Sage ELF system and fractions with 15-20 kb mean insert sizes were used for sequencing.
350   Polymerase/insert complex formation and clean up was performed using Sequel II$^{TM}$ binding
351   kit 2.0 based on manufacturer's instructions. Sequencing was performed on 16 8M SMART
352   cells using sequencing chemistry V2.0 and with a 2-hours pre-extension and 30-hours movie
353   time setting. CCS reads were obtained with PacBio CCS software
354   (https://github.com/PacificBiosciences/ccs).
355
356   Chromosome conformation capture sequencing (Hi-C) libraries were prepared from 1.5 g of
357   leaf material from one-week-old seedlings of cv. "Attraktion" as per the protocol of
358   Padmarasu et al. (2019) with a few modifications. The modifications include the use of nuclei
359   isolation protocol (Dvorak et al., 1988) and Ampure bead-based size selection instead of
360   SYBR-gold agarose gel-based size selection. The prepared library was quantified using qRT-
361   PCR using known concentration standards and sequenced on two lanes of a NovaSeq 6000
362   SP flow-cell using 200 cycles (2 x 100 bp paired-end mode).
363   *Data processing and enrichment efficiency calculation*

364   The adapter and low quality bases from raw RenSeq reads were removed using cutadapt
365   v1.16 (Martin, 2011) with a minimum read length of 30 bp after trimming. The quality check
366   for adapter and quality trimming was carried out using FastQC v0.11.7
367   (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/). Trimmed reads were
368   aligned against the reference genome assembly of cv. Chinese Spring (RefSeq v1.0, The
369   International Wheat Genome Sequencing Consortium (IWGSC) (2018)) using BWA-MEM
370   v0.7.17 (Li, 2013) with default parameters. The output was converted to binary alignment
371   map (BAM) format using SAMtools v1.9 (Li et al., 2009) and then the sorting was carried out
372   using NovoSort (V3.06.05). Sequences from the bait library were aligned to the RefSeq v1.0
373   using BLASTn v2.9.0 program (Altschul et al., 1990). Alignments with 95% identity and 70%
374   query coverage were retained and alignments separated by 120 bp or less were merged
375   using bedtools v2.29.2 (Quinlan and Hall, 2010). Finally, the total size of regions >= 100 bp in
376   the reference genome covered by alignments to the baits was calculated and considered as

377 the size of our capture target. The sorted BAM file of each genotype was used to calculate
378 the number of reads mapped on target and on the whole genome using SAMtools. The
379 enrichment factor (EF) was then determined as (N/M)/(T/G), where N is the number of reads
380 mapped on target, M indicates the total number of mapped reads, T denotes the size of the
381 targeted region and G is the size of the genome.
382
383 *De novo RenSeq assembly and NBS-LRR identification*
384
385 Only genotypes with at least 1 million reads and an enrichment factor >=100 were
386 considered. The quality trimmed data from each genotype was assembled *de novo* with CLC
387 Assembly cell (https://digitalinsights.qiagen.com/products-overview/discovery-insights-
388 portfolio/analysis-and-visualization/qiagen-clc-assembly-cell/) using the parameters -w=64 -
389 p fb ss 200 900. The contigs from each genotype were annotated with AUGUSTUS v3.3.145
390 (Stanke et al., 2006) using wheat gene models as training datasets, and contigs harboring
391 complete genes were identified. Amino acid (AA), coding sequence (CDS) and transcript
392 sequence for each complete gene were extracted using getAnnoFasta.pl script from the
393 AUGUSTUS package. AA sequences of gene models were used to predict protein domains
394 using the pfam_scan.pl script from PfamScan (Chojnacki et al., 2017), which searches FASTA
395 sequences against the Pfam HMM database (Mistry et al., 2020). The script was run with
396 sequence e-value cutoff of $10^{-5}$ and domain e-value cutoff 0.2 keeping other parameters to
397 default. Genes containing at least one NB-ARC (NBS) domain (pfam ID PF00931.23) were
398 considered as NLRs and used for downstream analysis. There is no standard tool available to
399 predict coiled coil (CC) domain, therefore all NLRs with the "Rx_N" (PF18052) domain which
400 is predicted as coiled coil were classified as coiled coil (CC). The sequences were further
401 classified as NBS (only NBS domain), NLs (NBS + LRRs), CNs (Rx_N + NBS), CNLs (Rx_N + NBS +
402 LRRs), XN (Integrated domain (ID) + NBS), XNL (ID + NBS + LRR), XRN (ID + Rx + N), XCNL (ID +
403 Rx_N + NBS + LRR) based on domain composition. A bash script was used to retrieve gene
404 structure information such as gene length, number of exons and introns from GFF files. The
405 CDS sequences of NLRs from each genotype were aligned against RefSeq v1.0 using GMAP
406 (Wu and Watanabe, 2005). The alignments were filtered with 70% query coverage and 95%
407 identity cutoff.
408
409 *Clustering and saturation analysis*
410
411 The AA sequences of all the NLRs identified from all genotypes were clustered using the
412 easy-linclust workflow from MMseqs2 software suite (Steinegger and Söding, 2017). The
413 program was run with e-value 1e-15, --min-seq-id 0.95, --seq-id-mode 2 (longer sequences),
414 --cluster-mode 2 (coverage of query), --kmer-per-seq 100, keeping other parameters to
415 default). The saturation analysis was carried out to determine the number of genotypes
416 needed to capture 95% OGs. This was done by making random selection of the genotypes
417 and counting the number of OGs present in these selections. The analysis was carried out
418 separately for all the genotypes, only plant genetic resources and only using elite lines. The
419 process was repeated 100 times starting with two and ending with a maximum number of
420 genotypes for each category.
421
422 *Identification of elite and PGR-specific OGs*
423

424    The genotype information of NLRs from each OG (variable *OG*). was used to classify the OG
425    as specific to elite lines, specific to plant genetic resources (PGR) or common (variable *Type*).
426    Further, the resistant OGs were identified based on the following linear model:

$$Phenotype = Type + OG$$

427    The OGs with negative effect and P value <0.01 were considered as resistant OGs. The
428    number of resistant OGs from respective accessions were counted and correlation of OG
429    count with disease susceptibility was studied.
430
431    *Chromosome scale genome assembly of "Attraktion" cultivar*
432
433    The HiFi reads were assembled using hifiasm v0.14 (Cheng et al., 2021) to generate a
434    primary contig assembly. The pseudomolecule construction was carried out using the TRITEX
435    pipeline (Monat et al., 2019). For this, the guide map was constructed by aligning the single
436    copy sequences from Julius to the Attraktion contig assembly. The Hi-C data was then used
437    for chimera breaking and contig ordering to generate pseudomolecules.
438    Transposable elements (TE) were annotated using a homology based approach implemented
439    in RepeatMasker v4.0.8 (Smit et al., 2004). A custom library was created by downloading and
440    combining wheat TE sequences from ClariTeRep: https://github.com/jdaron/CLARI-TE) and
441    2825      complete      plant      TE      sequences      ([http://botserv2.uzh.ch/kelldata/trep-](http://botserv2.uzh.ch/kelldata/trep-db/downloads/trep-db_complete_Rel-16.fasta.gz)
442    [db/downloads/trep-db_complete_Rel-16.fasta.gz](http://botserv2.uzh.ch/kelldata/trep-db/downloads/trep-db_complete_Rel-16.fasta.gz)).
443    Gene annotation was carried out using AUGUSTUS v3.3.1. Initially, CDS sequences of high
444    confidence (HC) genes from RefSeq v2.1 (Zhu et al., 2021a) were aligned to the Attraktion
445    assembly using GMAP-GSNAP and Alignment was filtered with 70% coverage and 95%
446    identity and top hits for each gene were extracted. The outputs of RepeatMasker and
447    GMAP-GSNAP were combined and a GFF file was created. This GFF file along with wheat
448    gene models served as a training dataset for AUGUSTUS.
449
450    The assembly completeness was assessed with 1,614 Benchmarking Universal Single Copy
451    Orthologs (BUSCO v5.1.2) (Simao et al., 2015) genes from plants using "genome mode". The
452    assembly quality was also evaluated both as genome and gene level. For genome-wide
453    comparison, single copy sequences from RefSeq v2.1 were aligned to Attraktion assembly
454    using minimap2 v2.17 (Li, 2018). For gene level comparison, the transcript sequences of HC
455    genes from RefSeq v2.1 were aligned against the transcript sequences of Attraktion using
456    LAST (Kiełbasa et al., 2011). Alignment filtration and synteny analysis was carried out using
457    MCScan (https://github.com/tanghaibao/jcvi/wiki/MCscan-%28Python-version%29).
458
459    The structural variations (SVs) were detected using the SyRI pipeline (Goel et al., 2019) with
460    default parameters. For this, the Refseqv2.1 assembly was aligned to the Attraktion
461    assembly using unimap, a fork of minimap2 optimized for assembly-to-reference
462    comparison. The script "sam2delta.py" from RaGOO (Alonge et al., 2019) was used for SAM
463    to mummer-delta format conversion. The delta file was filtered using delta-filter utility from
464    MuMmer v4.0 (Marçais et al., 2018) to filter out smaller alignments (2000 bp) and the file
465    was converted to TSV format using show-coords utility from MuMmer v4.0. The TSV file
466    served as input for SyRI. The SVs from SyRI were reclassified into presence-absence
467    variations (PAVs), inversions and translocations as follows: The CPL, DEL, DUP/INVDP (loss)
468    variants, and the Attraktion sequences in NOTAL and TDM were converted as Absence SVs
469    (relative to Attraktion). The CPG, INS, DUP/INVDP (gain) variants, and the query sequences in
470    NOTAL and TDM were converted as Presence SVs (relative to Attraktion). The INV variants

471    were regarded as inversions while the TRANS and INVTR were both regarded as
472    translocation SVs.
473
474    *Phenotypic records and analyses*
475
476    The experimental setup and quality assessment of yellow rust data were already presented
477    in detail elsewhere (Schulthess et al., 2021). Briefly, five yellow rust artificially inoculated
478    experiments plus two experiments relying on natural infections were conducted at five
479    different German locations during harvest years 2019 and 2020. In three out of these seven
480    field experiments, the presence of natural leaf rust infection was also recorded. Further
481    details on these experiments can be found in **Table S3**. Yellow and leaf rust infection severity
482    were expressed in a 1 (no symptoms) to 9 (severe infection) scoring scale according to the
483    Protocols       of       the       German       Federal       Plant       Variety       Office
484    (http://www.bundessortenamt.de/internet30/fileadmin/Files/PDF/Richtlinie_LW2000.pdf).
485    Outlier correction within and heritability within and across leaf rust experiments were
486    assessed in the same way as for yellow rust (Schulthess et al., 2021). The best linear
487    unbiased estimations (BLUEs) across experiments of yellow and leaf rust were used as
488    phenotypes for downstream analyses (**Table S4**).
489
490    *Reference-based GWAS*
491
492    The alignment records against the Chinese Spring reference in BAM format (see above) were
493    used for variant identification. The variant calling was performed using the mpileup and call
494    functions from SAMtools v1.9 and BCFtools (v1.8) (Li, 2011). The software was run with the -
495    DV parameter for SAMtools mpileup and minimum read quality (-q) cutoff of 20. The bi-
496    allelic SNPs were further filtered with minimum QUAL >=40; minimum read depth for
497    homozygous call >=2 and minimum read depth for heterozygous calls >=4 using a custom
498    AWK script. The commands were run in parallel wherever applicable to reduce
499    computational time using GNU parallel (Tange, 2011).
500    The variant calling was also performed by aligning adapter and quality trimmed reads from
501    each genotype against the genome assembly of Attraktion. Variant calling and SNP filtration
502    were performed as described above except that minimap2 v.2.17 was used for mapping
503    reads against the Attraktion assembly.
504    The GWAS for yellow rust and leaf rust was carried out using a univariate linear mixed model
505    from GEMMA (v0.98) software (Zhou and Stephens, 2012) with -lmm 4 -miss 0.2 -maf 0.01
506    parameters, while keeping other parameters to default settings. Relatedness and population
507    structure were accounted for using a kinship matrix in the form $2*(1-RD)$, where RD are the
508    Rogers' distances between genotypes computed from 17,840 high-quality GBS SNPs
509    (Schulthess et al., 2021). GWAS was also carried out separately using PGRs and elite
510    accessions to identify novel sources of resistance from PGRs, avoiding over-correction by the
511    kinship matrix.
512    For paGWAS, SNPs with more than 20% missing data were scored as presence-absence
513    variants. The presence-absence status of genotype calls was converted to
514    reference/alternate allele calls. GWAS with these data were done with GEMMA as described
515    above.
516
517    *Reference-free GWAS*
518

11

519 The reference free GWAS was carried out using the kmersGWAS pipeline (Voichek and
520 Weigel, 2020). Briefly, 31 bp *k*-mers that were supported by at least five reads were
521 extracted using kmctools v3.1.1 (Kokot et al., 2017). The *k*-mers from all the genotypes were
522 combined and a non-redundant *k*-mer presence-absence genotype matrix was generated.
523 The pipeline was run with 100 permutations, the 5 million top *k*-mers and minor allele
524 frequency 0.01, while setting other parameters to default. The significance threshold was
525 determined by selecting 5th top P-value from the 100 top P values obtained from 100
526 permutations. To position the *k*-mers in the genome, they were mapped against the
527 reference assemblies of Chinese Spring and Attraktion and positions of uniquely mapped *k*-
528 mers were retrieved. The positional information along with p-values were used for
529 generation of Manhattan plots using the qqman (Turner) R package.
530 The significant *k*-mers were also aligned to the NBS-LRR transcript database generated
531 above and the number of *k*-mers aligned with 100% identity and 100% coverage to
532 representative transcript sequences of each cluster were counted. The results were
533 manually inspected and the candidate clusters with large proportions of significant *k*-mers
534 for yellow rust and leaf rust were identified.
535
536 *Candidate gene identification*
537
538 The results from various methods mentioned above were compared and consensus regions
539 for yellow rust resistance were determined. The MCscan
540 (https://github.com/tanghaibao/jcvi/wiki/MCscan-(Python-version)) software with default
541 parameters was used to study local synteny between different wheat genome assemblies.
542 The CDS sequences of candidate clusters identified based on the kmerGWAS method were
543 aligned separately to each of the consensus regions identified for yellow rust using GMAP
544 (Wu and Watanabe, 2005) and candidate genes were identified. The AA sequences of the
545 NBS domain of candidate genes were extracted and aligned with NBS domains of cloned R
546 genes from wheat using MAFFT v7.305 (Katoh et al., 2002). The spurious sequences or
547 poorly aligned regions were removed using trimAl v1.2 (Capella-Gutiérrez et al., 2009). The
548 phylogenetic analysis was carried out using IQ-TREE v1.6.12 (Nguyen et al., 2015). The
549 phylogenetic tree was visualized using ggtree (Yu et al., 2017).
550
551 **Data availability**

568
569    **Author contributions**
570    JCR, NS and MM designed research. AWS, PHB and JS produced phenotypic data. SMK and
571    AH performed capture sequencing. AWS selected plant material and analyzed phenotypic
572    data. SMK analyzed molecular data and performed association analyses. BS and BBHW
573    provided the bait library. SP and AH generated PacBio HiFi and Hi-C data. SMK and MM
574    constructed the reference assembly of cv. Attraktion. SMK, AWS and MM wrote the paper
575    with input from all co-authors.
576

577 **Tables**

578

579 **Table 1: Statistics of the genome sequence assembly of cv. Attraktion**

580

| | |
|---|---|
| Assembly size | 14.7 Gb |
| Number of contigs | 4,953 |
| Contig N50 | 17.3 Mb |
| Contig N90 | 4.1 Mb |
| Pseudomolecule size | 14.3 Gb |
| Number of contigs in pseudomolecules | 1,553 |
| Complete BUSCOs | 1,584 (98.2%) |

581

582

14

## Supplementary items

**Table S1:** Passport data and ENA accession numbers for 907 winter wheat accessions.

**Table S2:** RenSeq assembly statistics.

**Table S3:** Experimental setup and data quality assessment of leaf rust data.

**Table S4:** BLUEs of yellow and leaf rust used for association analyses.

**Figure S1: Distribution of NBS-LRR orthogroup size in the hexaploid wheat collection.** No orthogroups common to all the accessions were found because of random sequence drop out.

**Figure S2: Intrachromosomal Hi-C contact matrices for pseudomolecules of cv. Attraktion.**

**Figure S3: Whole-chromosome alignments of the pseudomolecules of cv. Attraktion to Chinese Spring RefSeq V2.1 assembly (Zhu et al. 2021).**

**Figure S4: GWAS using SNPs identified relative to the reference assembly of Chinese Spring (RefSeq V1.0).** GWAS was done in the panels: elite **(a)**, PGR **(b)**, and Elite and PGR combined **(c)**. No significant marker-trait associations were detected in the PGR panel at the 6A locus, indicating that the resistant haplotype is rare in PGR.

**Figure S5**: **Association scans for leaf rust resistance using different marker systems. (a)** SNPs identified relative to Chinese Spring RefSeq V1.0. **(b)** presence-absence GWAS using SNPs identified relative Chinese Spring RefSeq V1.0 and scored as presence-absence markers. **(c)** $k$-mers mapped against Chinese Spring RefSeq V1.0. Panels **(d)**, **(e)** and **(f)** show the results of SNP-based, presence-absence and $k$-mer based GWAS when the reference sequence of cv. Attraktion was used for SNP identification or $k$-mer mapping. The blue horizontal lines indicate threshold above which associations are statistically significant.

## References

Alonge, M., Soyk, S., Ramakrishnan, S., Wang, X., Goodwin, S., Sedlazeck, F.J., Lippman, Z.B. and Schatz, M.C. (2019) RaGOO: fast and accurate reference-guided scaffolding of draft genomes. *Genome Biol.* **20**, 224.

Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J. (1990) Basic local alignment search tool. *J Mol Biol* **215**, 403-410.

Arend, D., Junker, A., Scholz, U., Schüler, D., Wylie, J. and Lange, M. (2016) PGP repository: a plant phenomics and genomics data publication infrastructure. *Database* **2016**.

Arend, D., Lange, M., Chen, J., Colmsee, C., Flemming, S., Hecht, D. and Scholz, U. (2014) e! DAL-a framework to store, share and publish research data. *BMC bioinformatics* **15**, 214.

Arora, S., Steuernagel, B., Gaurav, K., Chandramohan, S., Long, Y., Matny, O., Johnson, R., Enk, J., Periyannan, S., Singh, N., Asyraf Md Hatta, M., Athiyannan, N., Cheema, J., Yu, G., Kangara, N., Ghosh, S., Szabo, L.J., Poland, J., Bariana, H., Jones, J.D.G., Bentley, A.R., Ayliffe, M., Olson, E., Xu, S.S., Steffenson, B.J., Lagudah, E. and Wulff, B.B.H. (2019) Resistance gene cloning from a wild crop relative by sequence capture and association genetics. *Nat Biotechnol* **37**, 139-143.

Capella-Gutiérrez, S., Silla-Martínez, J.M. and Gabaldón, T. (2009) trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972-1973.

Cheng, H., Concepcion, G.T., Feng, X., Zhang, H. and Li, H. (2021) Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nature Methods* **18**, 170-175.

Chojnacki, S., Cowley, A., Lee, J., Foix, A. and Lopez, R. (2017) Programmatic access to bioinformatics tools from EMBL-EBI update: 2017. *Nucleic Acids Research* **45**, W550-W553.

Della Coletta, R., Qiu, Y., Ou, S., Hufford, M.B. and Hirsch, C.N. (2021) How the pan-genome is changing crop genomics and improvement. *Genome Biology* **22**, 3.

Dvorak, J., McGuire, P.E. and Cassidy, B. (1988) Apparent sources of the A genomes of wheats inferred from polymorphism in abundance and restriction fragment length of repeated nucleotide sequences. *Genome* **30**, 680-689.

Gabur, I., Chawla, H.S., Liu, X., Kumar, V., Faure, S., von Tiedemann, A., Jestin, C., Dryzska, E., Volkmann, S., Breuer, F., Delourme, R., Snowdon, R. and Obermeier, C. (2018) Finding invisible quantitative trait loci with missing data. *Plant Biotechnol J* **16**, 2102-2112.

Gaurav, K., Arora, S., Silva, P., Sánchez-Martín, J., Horsnell, R., Gao, L., Brar, G.S., Widrig, V., Raupp, J., Singh, N., Wu, S., Kale, S.M., Chinoy, C., Nicholson, P., Quiroz-Chávez, J., Simmonds, J., Hayta, S., Smedley, M.A., Harwood, W., Pearce, S., Gilbert, D., Kangara, N., Gardener, C., Forner-Martínez, M., Liu, J., Yu, G., Boden, S., Pascucci, A., Ghosh, S., Hafeez, A.N., O'Hara, T., Waites, J., Cheema, J., Steuernagel, B., Patpour, M., Justesen, A.F., Liu, S., Rudd, J.C., Avni, R., Sharon, A., Steiner, B., Kirana, R.P., Buerstmayr, H., Mehrabi, A.A., Nasyrova, F.Y., Chayut, N., Matny, O., Steffenson, B.J., Sandhu, N., Chhuneja, P., Lagudah, E., Elkot, A.F., Tyrrell, S., Bian, X., Davey, R.P., Simonsen, M., Schauser, L., Tiwari, V.K., Kutcher, H.R., Hucl, P., Li, A., Liu, D.-C., Mao, L., Xu, S., Brown-Guedira, G., Faris, J., Dvorak, J., Luo, M.-C., Krasileva, K., Lux, T., Artmeier, S., Mayer, K.F.X., Uauy, C., Mascher, M., Bentley, A.R., Keller, B., Poland, J. and Wulff, B.B.H. (2021) Evolution of the bread wheat D-subgenome and enriching it with diversity from <em>Aegilops tauschii</em>. *bioRxiv*, 2021.2001.2031.428788.

666    Goel, M., Sun, H., Jiao, W.-B. and Schneeberger, K. (2019) SyRI: finding genomic
667          rearrangements and local sequence differences from whole-genome assemblies.
668          *Genome Biology* **20**, 277.
669    Hafeez, A.N., Arora, S., Ghosh, S., Gilbert, D., Bowden, R.L. and Wulff, B.B.H. (2021) Creation
670          and judicious application of a wheat resistance gene atlas. *Mol Plant* **14**, 1053-1070.
671    Hickey, G., Heller, D., Monlong, J., Sibbesen, J.A., Sirén, J., Eizenga, J., Dawson, E.T., Garrison,
672          E., Novak, A.M. and Paten, B. (2020) Genotyping structural variants in pangenome
673          graphs using the vg toolkit. *Genome Biology* **21**, 35.
674    Hovmøller, M.S., Walter, S., Bayles, R.A., Hubbard, A., Flath, K., Sommerfeldt, N., Leconte,
675          M., Czembor, P., Rodriguez-Algaba, J., Thach, T., Hansen, J.G., Lassen, P., Justesen,
676          A.F., Ali, S. and de Vallavieille-Pope, C. (2016) Replacement of the European wheat
677          yellow rust population by new races from the centre of diversity in the near-
678          Himalayan region. *Plant Pathology* **65**, 402-411.
679    Huang, L., Brooks, S.A., Li, W., Fellers, J.P., Trick, H.N. and Gill, B.S. (2003) Map-based cloning
680          of leaf rust resistance gene Lr21 from the large and polyploid genome of bread
681          wheat. *Genetics* **164**, 655-664.
682    Jayakodi, M., Padmarasu, S., Haberer, G., Bonthala, V.S., Gundlach, H., Monat, C., Lux, T.,
683          Kamal, N., Lang, D., Himmelbach, A., Ens, J., Zhang, X.-Q., Angessa, T.T., Zhou, G., Tan,
684          C., Hill, C., Wang, P., Schreiber, M., Boston, L.B., Plott, C., Jenkins, J., Guo, Y., Fiebig,
685          A., Budak, H., Xu, D., Zhang, J., Wang, C., Grimwood, J., Schmutz, J., Guo, G., Zhang,
686          G., Mochida, K., Hirayama, T., Sato, K., Chalmers, K.J., Langridge, P., Waugh, R.,
687          Pozniak, C.J., Scholz, U., Mayer, K.F.X., Spannagl, M., Li, C., Mascher, M. and Stein, N.
688          (2020) The barley pan-genome reveals the hidden legacy of mutation breeding.
689          *Nature* **588**, 284-289.
690    Jupe, F., Witek, K., Verweij, W., Sliwka, J., Pritchard, L., Etherington, G.J., Maclean, D., Cock,
691          P.J., Leggett, R.M., Bryan, G.J., Cardle, L., Hein, I. and Jones, J.D. (2013) Resistance
692          gene enrichment sequencing (RenSeq) enables reannotation of the NB-LRR gene
693          family from sequenced plant genomes and rapid mapping of resistance loci in
694          segregating populations. *Plant J* **76**, 530-544.
695    Katoh, K., Misawa, K., Kuma, K.-I. and Miyata, T. (2002) MAFFT: a novel method for rapid
696          multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**,
697          3059-3066.
698    Kiełbasa, S.M., Wan, R., Sato, K., Horton, P. and Frith, M.C. (2011) Adaptive seeds tame
699          genomic sequence comparison. *Genome Res.* **21**, 487-493.
700    Kokot, M., Dlugosz, M. and Deorowicz, S. (2017) KMC 3: counting and manipulating k-mer
701          statistics. *Bioinformatics* **33**, 2759-2761.
702    Li, H. (2011) A statistical framework for SNP calling, mutation discovery, association mapping
703          and population genetical parameter estimation from sequencing data. *Bioinformatics*
704          **27**, 2987-2993.
705    Li, H. (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-
706          MEM. *arXiv preprint arXiv:1303.3997*.
707    Li, H. (2018) Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**,
708          3094-3100.
709    Li, H., Feng, X. and Chu, C. (2020) The design and construction of reference pangenome
710          graphs with minigraph. *Genome Biology* **21**, 265.
711    Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.
712          and Durbin, R. (2009) The sequence alignment/map format and SAMtools.
713          *Bioinformatics* **25**, 2078-2079.

714    Liu, F., Zhao, Y., Beier, S., Jiang, Y., Thorwarth, P., CF, H.L., Ganal, M., Himmelbach, A., Reif,
715            J.C. and Schulthess, A.W. (2020a) Exome association analysis sheds light onto leaf
716            rust (Puccinia triticina) resistance genes currently used in wheat breeding (Triticum
717            aestivum L.). *Plant Biotechnol J* **18**, 1396-1408.
718    Liu, Y., Du, H., Li, P., Shen, Y., Peng, H., Liu, S., Zhou, G.-A., Zhang, H., Liu, Z., Shi, M., Huang,
719            X., Li, Y., Zhang, M., Wang, Z., Zhu, B., Han, B., Liang, C. and Tian, Z. (2020b) Pan-
720            Genome of Wild and Cultivated Soybeans. *Cell* **182**, 162-176.e113.
721    Lopez-Arboleda, W.A., Reinert, S., Nordborg, M. and Korte, A. (2021) Global Genetic
722            Heterogeneity in Adaptive Traits. *Molecular Biology and Evolution*.
723    López-Girona, E., Davy, M.W., Albert, N.W., Hilario, E., Smart, M.E.M., Kirk, C., Thomson, S.J.
724            and Chagné, D. (2020) CRISPR-Cas9 enrichment and long read sequencing for fine
725            mapping in plants. *Plant Methods* **16**, 121.
726    Marçais, G., Delcher, A.L., Phillippy, A.M., Coston, R., Salzberg, S.L. and Zimin, A. (2018)
727            MUMmer4: A fast and versatile genome alignment system. *PLoS Comput. Biol.* **14**,
728            e1005944.
729    Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing
730            reads. *EMBnet. Journal* **17**, pp. 10-12.
731    Mascher, M., Wicker, T., Jenkins, J., Plott, C., Lux, T., Koh, C.S., Ens, J., Gundlach, H., Boston,
732            L.B., Tulpová, Z., Holden, S., Hernández-Pinzón, I., Scholz, U., Mayer, K.F.X., Spannagl,
733            M., Pozniak, C.J., Sharpe, A.G., Šimková, H., Moscou, M.J., Grimwood, J., Schmutz, J.
734            and Stein, N. (2021) Long-read sequence assembly: a technical evaluation in barley.
735            *Plant Cell*.
736    Michelmore, R.W. and Meyers, B.C. (1998) Clusters of resistance genes in plants evolve by
737            divergent selection and a birth-and-death process. *Genome Res* **8**, 1113-1130.
738    Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, Gustavo A., Sonnhammer,
739            E.L.L., Tosatto, S.C.E., Paladin, L., Raj, S., Richardson, L.J., Finn, R.D. and Bateman, A.
740            (2020) Pfam: The protein families database in 2021. *Nucleic Acids Research* **49**, D412-
741            D419.
742    Monat, C., Padmarasu, S., Lux, T., Wicker, T., Gundlach, H., Himmelbach, A., Ens, J., Li, C.,
743            Muehlbauer, G.J., Schulman, A.H., Waugh, R., Braumann, I., Pozniak, C., Scholz, U.,
744            Mayer, K.F.X., Spannagl, M., Stein, N. and Mascher, M. (2019) TRITEX: chromosome-
745            scale sequence assembly of Triticeae genomes with open-source tools. *Genome Biol*
746            **20**, 284.
747    Nguyen, L.-T., Schmidt, H.A., von Haeseler, A. and Minh, B.Q. (2015) IQ-TREE: a fast and
748            effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol.*
749            *Biol. Evol.* **32**, 268-274.
750    Padmarasu, S., Himmelbach, A., Mascher, M. and Stein, N. (2019) In Situ Hi-C for Plants: An
751            Improved Method to Detect Long-Range Chromatin Interactions. *Methods Mol Biol*
752            **1933**, 441-472.
753    Przewieslik-Allen, A.M., Wilkinson, P.A., Burridge, A.J., Winfield, M.O., Dai, X., Beaumont, M.,
754            King, J., Yang, C.-y., Griffiths, S., Wingen, L.U., Horsnell, R., Bentley, A.R., Shewry, P.,
755            Barker, G.L.A. and Edwards, K.J. (2021) The role of gene flow and chromosomal
756            instability in shaping the bread wheat genome. *Nature Plants* **7**, 172-183.
757    Quinlan, A.R. and Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing
758            genomic features. *Bioinformatics* **26**, 841-842.
759    Sato, K., Abe, F., Mascher, M., Haberer, G., Gundlach, H., Spannagl, M., Shirasawa, K. and
760            Isobe, S. (2021) Chromosome-scale genome assembly of the transformation-
761            amenable common wheat cultivar 'Fielder'. *DNA Research* **28**.

762 Schulthess, A.W., Kale, S.M., Liu, F., Zhao, Y., Philipp, N., Rembe, M., Jiang, Y., Beukert, U.,
763      Serfling, A., Himmelbach, A., Fuchs, J., Oppermann, M., Weise, S., Boeven, P.H.G.,
764      Schacht, J., Longin, C.F.H., Kollers, S., Pfeiffer, N., Korzun, V., Lange, M., Scholz, U.,
765      Stein, N., Mascher, M. and Reif, J.C. (2021) GiPS: Genomics-informed parent selection
766      uncovers the breeding value of wheat genetic resources. *bioRxiv*,
767      2021.2012.2015.472759.
768 Sears, E.R. and Miller, T.E. (1985) THE HISTORY OF CHINESE SPRING WHEAT. *Cereal Research*
769      *Communications* **13**, 261-263.
770 Simao, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. and Zdobnov, E.M. (2015)
771      BUSCO: assessing genome assembly and annotation completeness with single-copy
772      orthologs. *Bioinformatics* **31**, 3210-3212.
773 Smit, A., Hubley, R. and Green, P. (2004) RepeatMasker Open-4.0.
774 Stanke, M., Keller, O., Gunduz, I., Hayes, A., Waack, S. and Morgenstern, B. (2006)
775      AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Research* **34**,
776      W435-W439.
777 Steinegger, M. and Söding, J. (2017) MMseqs2 enables sensitive protein sequence searching
778      for the analysis of massive data sets. *Nature Biotechnology* **35**, 1026-1028.
779 Steuernagel, B., Vrána, J., Karafiátová, M., Wulff, B.B.H. and Doležel, J. (2017) Rapid Gene
780      Isolation Using MutChromSeq. In: *Wheat Rust Diseases: Methods and Protocols*
781      (Periyannan, S. ed) pp. 231-243. New York, NY: Springer New York.
782 Tange, O. (2011) Gnu parallel-the command-line power tool. *The USENIX Magazine* **36**, 42-
783      47.
784 The International Wheat Genome Sequencing Consortium (IWGSC) (2018) Shifting the limits
785      in wheat research and breeding using a fully annotated reference genome. *Science*
786      **361**, eaar7191.
787 Ting, J.P.Y., Lovering, R.C., Alnemri, E.S., Bertin, J., Boss, J.M., Davis, B.K., Flavell, R.A.,
788      Girardin, S.E., Godzik, A., Harton, J.A., Hoffman, H.M., Hugot, J.-P., Inohara, N.,
789      Mackenzie, A., Maltais, L.J., Nunez, G., Ogura, Y., Otten, L.A., Philpott, D., Reed, J.C.,
790      Reith, W., Schreiber, S., Steimle, V. and Ward, P.A. (2008) The NLR gene family: a
791      standard nomenclature. *Immunity* **28**, 285-287.
792 Turner, S.D. (2018) qqman: an R package for visualizing GWAS results using Q-Q and
793      manhattan plots. *The Journal of open source software*.
794 Van de Weyer, A.L., Monteiro, F., Furzer, O.J., Nishimura, M.T., Cevik, V., Witek, K., Jones,
795      J.D.G., Dangl, J.L., Weigel, D. and Bemm, F. (2019) A Species-Wide Inventory of NLR
796      Genes and Alleles in Arabidopsis thaliana. *Cell* **178**, 1260-1272 e1214.
797 Voichek, Y. and Weigel, D. (2020) Identifying genetic variants underlying phenotypic
798      variation in plants without complete genomes. *Nat Genet* **52**, 534-540.
799 Walkowiak, S., Gao, L., Monat, C., Haberer, G., Kassa, M.T., Brinton, J., Ramirez-Gonzalez,
800      R.H., Kolodziej, M.C., Delorean, E., Thambugala, D., Klymiuk, V., Byrns, B., Gundlach,
801      H., Bandi, V., Siri, J.N., Nilsen, K., Aquino, C., Himmelbach, A., Copetti, D., Ban, T.,
802      Venturini, L., Bevan, M., Clavijo, B., Koo, D.-H., Ens, J., Wiebe, K., N'Daye, A., Fritz,
803      A.K., Gutwin, C., Fiebig, A., Fosker, C., Fu, B.X., Accinelli, G.G., Gardner, K.A., Fradgley,
804      N., Gutierrez-Gonzalez, J., Halstead-Nussloch, G., Hatakeyama, M., Koh, C.S., Deek, J.,
805      Costamagna, A.C., Fobert, P., Heavens, D., Kanamori, H., Kawaura, K., Kobayashi, F.,
806      Krasileva, K., Kuo, T., McKenzie, N., Murata, K., Nabeka, Y., Paape, T., Padmarasu, S.,
807      Percival-Alwyn, L., Kagale, S., Scholz, U., Sese, J., Juliana, P., Singh, R., Shimizu-
808      Inatsugi, R., Swarbreck, D., Cockram, J., Budak, H., Tameshige, T., Tanaka, T., Tsuji, H.,
809      Wright, J., Wu, J., Steuernagel, B., Small, I., Cloutier, S., Keeble-GagnÃ¨re, G.,

810    Muehlbauer, G., Tibbets, J., Nasuda, S., Melonek, J., Hucl, P.J., Sharpe, A., Clark, M.,
811    Legg, E., Bharti, A., Langridge, P., Hall, A., Uauy, C., Mascher, M., Krattinger, S.G.,
812    Handa, H., Shimizu, K.K., Distelfeld, A., Chalmers, K., Keller, B., Mayer, K.F.X., Poland,
813    J., Stein, N., McCartney, C.A., Spannagl, M., Wicker, T. and Pozniak, C.J. (2020)
814    Multiple wheat genomes reveal global variation in modern breeding. *Nature*
815    **accepted**.
816 Wu, T.D. and Watanabe, C.K. (2005) GMAP: a genomic mapping and alignment program for
817    mRNA and EST sequences. *Bioinformatics* **21**, 1859-1875.
818 Yu, G., Smith, D.K., Zhu, H., Guan, Y. and Lam, T.T.y. (2017) ggtree : an r package for
819    visualization and annotation of phylogenetic trees with their covariates and other
820    associated data. *Methods in Ecology and Evolution* **8**, 28-36.
821 Yu, J., Pressoir, G., Briggs, W.H., Vroh Bi, I., Yamasaki, M., Doebley, J.F., McMullen, M.D.,
822    Gaut, B.S., Nielsen, D.M., Holland, J.B., Kresovich, S. and Buckler, E.S. (2006) A unified
823    mixed-model method for association mapping that accounts for multiple levels of
824    relatedness. *Nat Genet* **38**, 203-208.
825 Zhang, J., Hewitt, T.C., Boshoff, W.H.P., Dundas, I., Upadhyaya, N., Li, J., Patpour, M.,
826    Chandramohan, S., Pretorius, Z.A., Hovmøller, M., Schnippenkoetter, W., Park, R.F.,
827    Mago, R., Periyannan, S., Bhatt, D., Hoxha, S., Chakraborty, S., Luo, M., Dodds, P.,
828    Steuernagel, B., Wulff, B.B.H., Ayliffe, M., McIntosh, R.A., Zhang, P. and Lagudah, E.S.
829    (2021) A recombined Sr26 and Sr61 disease resistance gene stack in wheat encodes
830    unrelated NLR genes. *Nat Commun* **12**, 3378.
831 Zhou, X. and Stephens, M. (2012) Genome-wide efficient mixed-model analysis for
832    association studies. *Nature Genetics* **44**, 821-824.
833 Zhu, T., Wang, L., Rimbert, H., Rodriguez, J.C., Deal, K.R., De Oliveira, R., Choulet, F., Keeble-
834    Gagnère, G., Tibbits, J., Rogers, J., Eversole, K., Appels, R., Gu, Y.Q., Mascher, M.,
835    Dvorak, J. and Luo, M.-C. (2021a) Optical maps refine the bread wheat Triticum
836    aestivum cv. Chinese Spring genome assembly. *Plant J.* **107**, 303-314.
837 Zhu, T., Wang, L., Rimbert, H., Rodriguez, J.C., Deal, K.R., De Oliveira, R., Choulet, F., Keeble-
838    Gagnère, G., Tibbits, J., Rogers, J., Eversole, K., Appels, R., Gu, Y.Q., Mascher, M.,
839    Dvorak, J. and Luo, M.C. (2021b) Optical maps refine the bread wheat Triticum
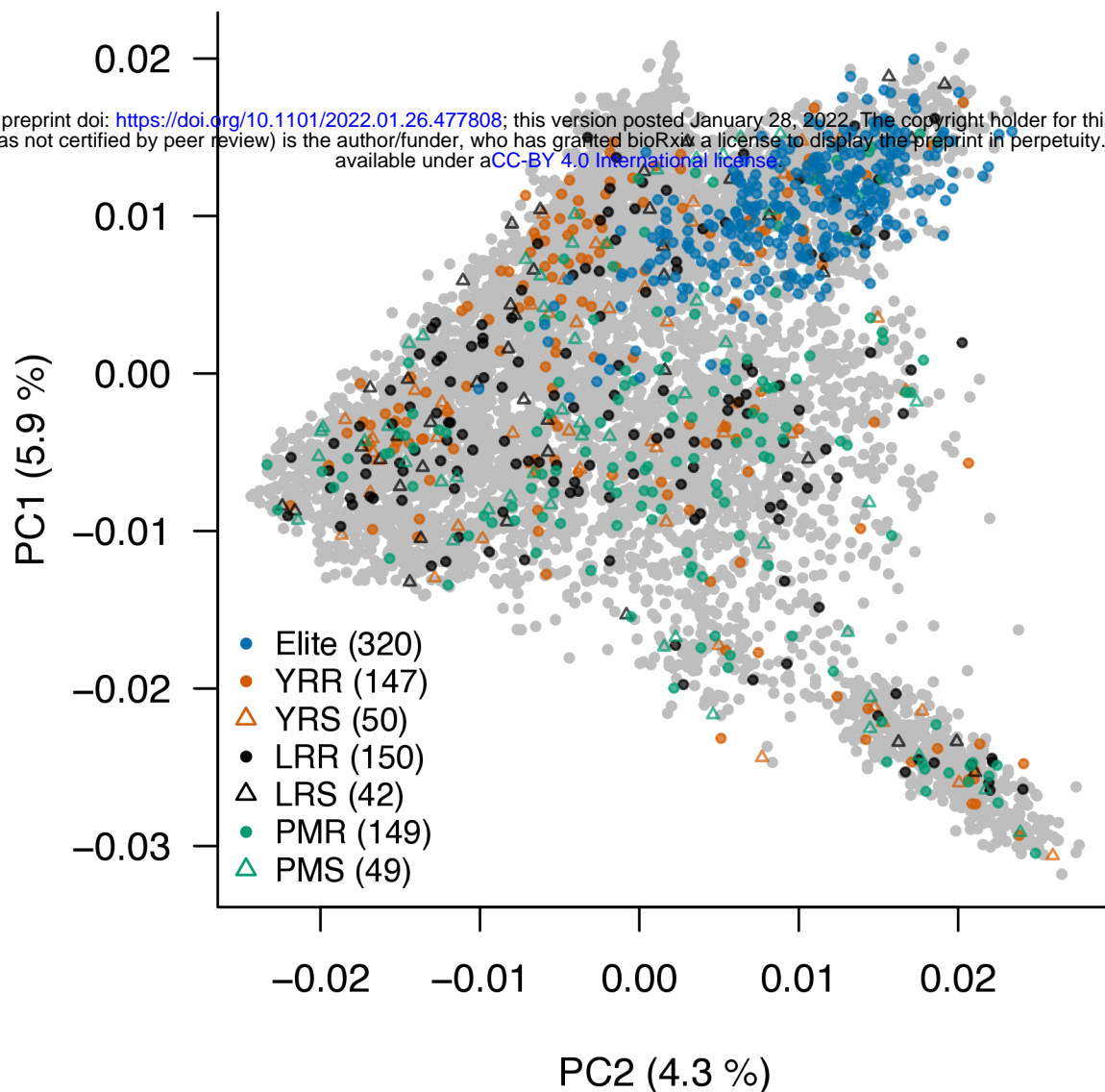840    aestivum cv. Chinese Spring genome assembly. *Plant J* **107**, 303-314.
841

**Figure 1:** Genetic diversity of genotypes selected for RenSeq. The 907 accessions selected for the RenSeq analysis were projected onto the molecular diversity space of the winter wheat collection of the IPK genebank portrayed by the first two principal components (PCs) from a PC analysis on genome-wide SNP markers (Schulthess et al. 2021). Among RenSeq characterizations, 192 European elite cultivars and 128 German elite breeding lines represent the diversity already handled by European breeding. The remaining fraction is composed of 587 plant genetic resources (PGRs) samples from the IPK genebank which was enriched for disease resistant genotypes with minimized population structure (Schulthess et al. 2021). According to selection, PGRs are classified as yellow rust resistant [YRR] or susceptible [YRS], leaf rust resistant [LRR] or susceptible [LRS] and powdery mildew resistant [PMR] or susceptible [PMS].

**Figure 2:** Proportion of NLRs of different classes in individual RenSeq assemblies. CN: Colied-coil-NBS; N: NBS; CNL: Colied-coil-NBS-LRR; NL: NBS-LRR; XCN: integrated domain (ID)-Colied-coil-NBS; XCNL: ID-Colied-coil-NBS-LRR; XNL: ID-NBS-LRR.
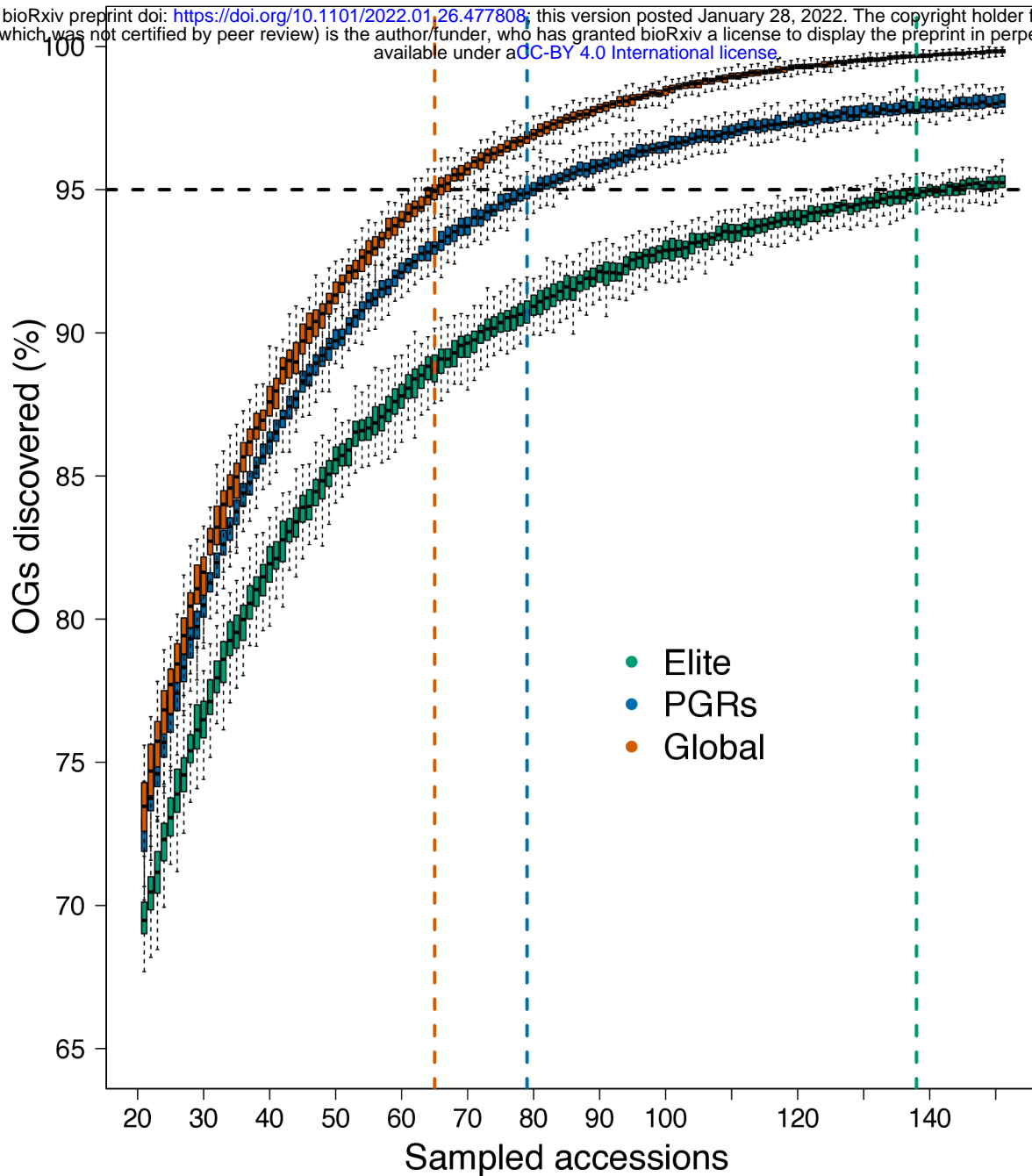
**Figure 3: Saturation analysis.** Fraction of NLR orthogroups recovered from randomly drawn subsets of genotypes. Subsets were selected from the entire population as well as elite varieties and PGRs. Sampling was repeated 100 times for subsets of increasing size. Colored vertical lines indicate the number of accessions required to achieve 95 % representation of the NLR universe.

**Figure 4: Enrichment of resistance-associated NLR clusters in elite lines. (a)** The number of resistance-associated clusters from an accession is plotted against its yellow rust susceptibility score. Elite lines with many resistance-associated clusters were less susceptible to yellow rust. **(b)** Genomic distribution of elite-specific orthogroups (OGs, green), PGR-specific OGs (blue) and OGs present in both elite lines and PGRs (orange) in the three subgenomes of hexploid wheat (A, B, D). The grey boxes mark the positions of alien introgressions.

**Figure 5: Structural variants detected by Synteny and Rearrangement Identifier (SyRI, Goel et al. 2019) between the genome assemblies of cv. Attraktion (reference) and Chinese spring RefSeq V2.1 (query).**
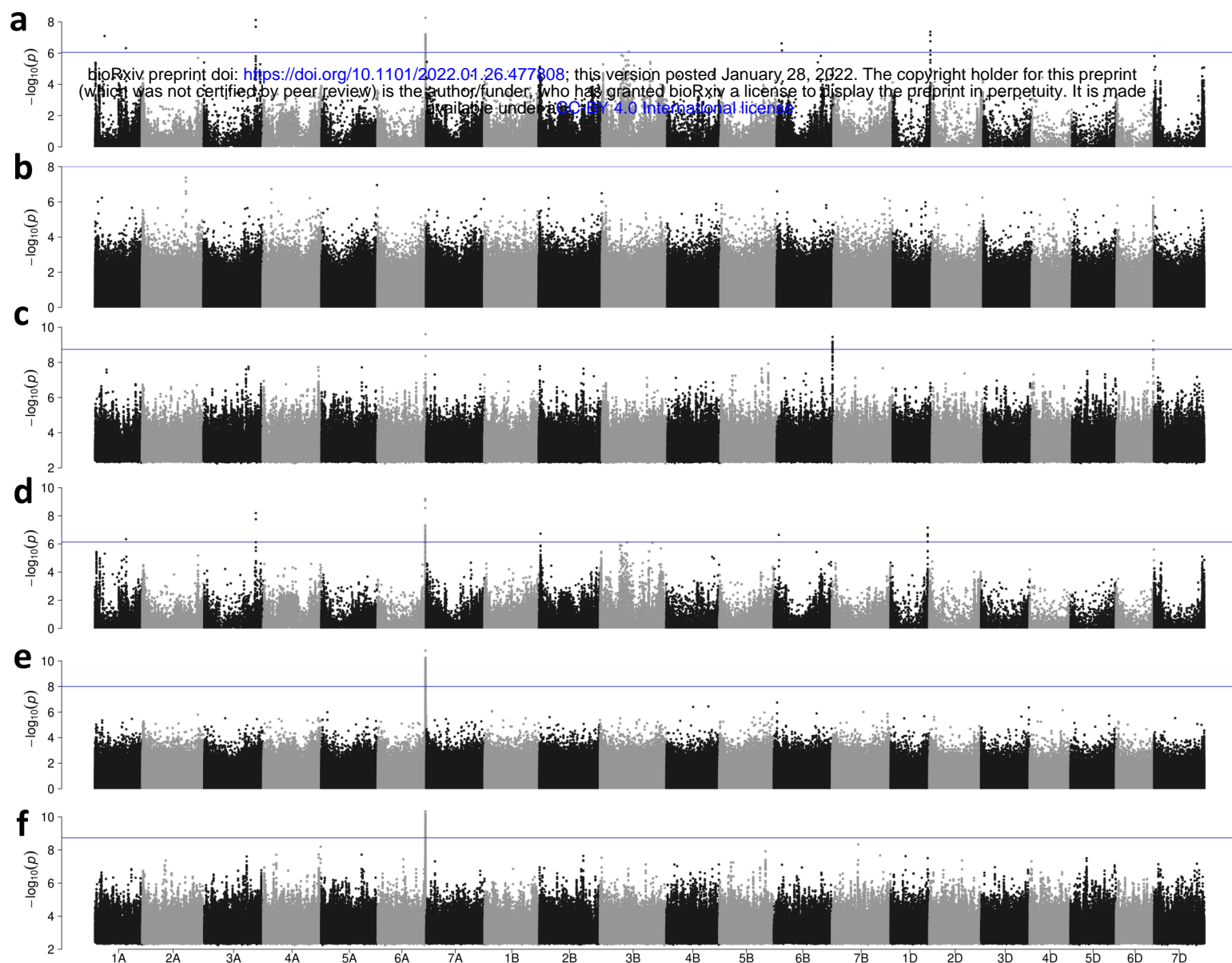
**Figure 6: Association scans for yellow rust resistance using different marker systems. (a)** SNPs identified relative to Chinese Spring RefSeq V1.0. **(b)** presence-absence GWAS using SNPs identified relative to Chinese Spring RefSeq V1.0 and scored as presence-absence markers. **(c)** *k*-mers mapped against Chinese Spring RefSeq V1.0. Panels **(d)**, **(e)** and **(f)** show the results of SNP-based, presence-absence and *k*-mer based GWAS when the reference sequence of cv. Attraktion was used for SNP identification or *k*-mer mapping. The blue horizontal lines indicate the threshold above which associations are statistically significant.
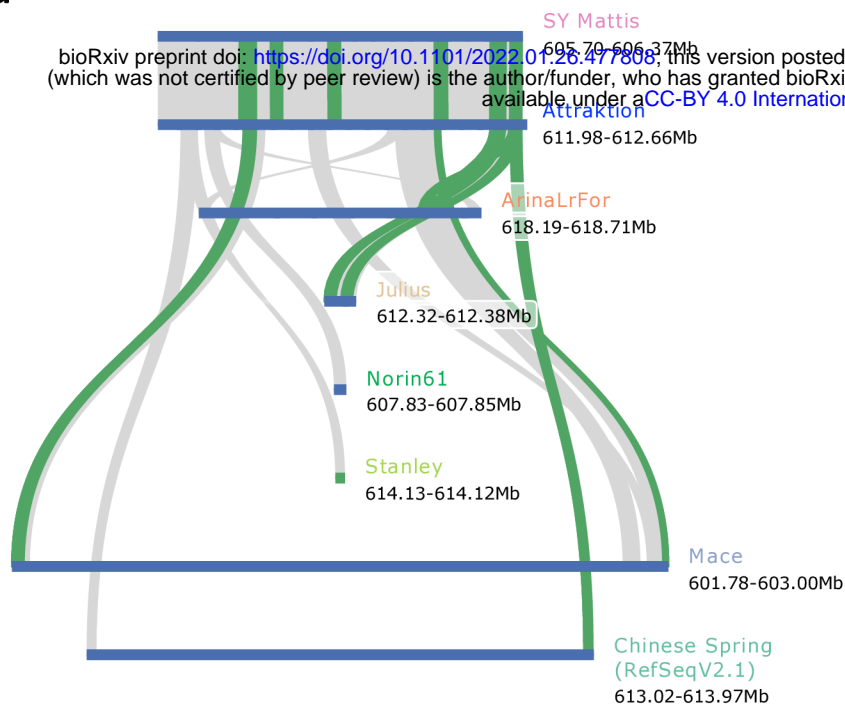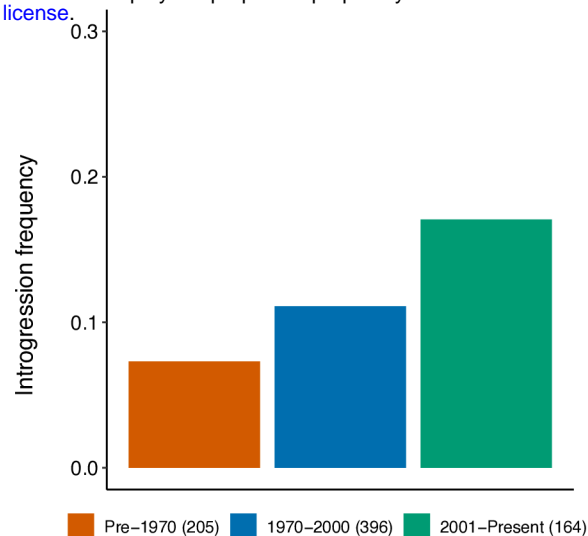
**Figure 7: Tracing the history of a novel yellow rust resistance locus on chromosome 6A. (a)** Gene-based colinearity analysis of the yellow rust resistance locus identified on chromosome 6A in the Attraktion assembly with reference assemblies from the wheat pan-genome (Walkowiak et al. 2020). (**b**) Frequency of resistant haplotypes in accessions from different time period. The numbers in parentheses indicate size of each group.
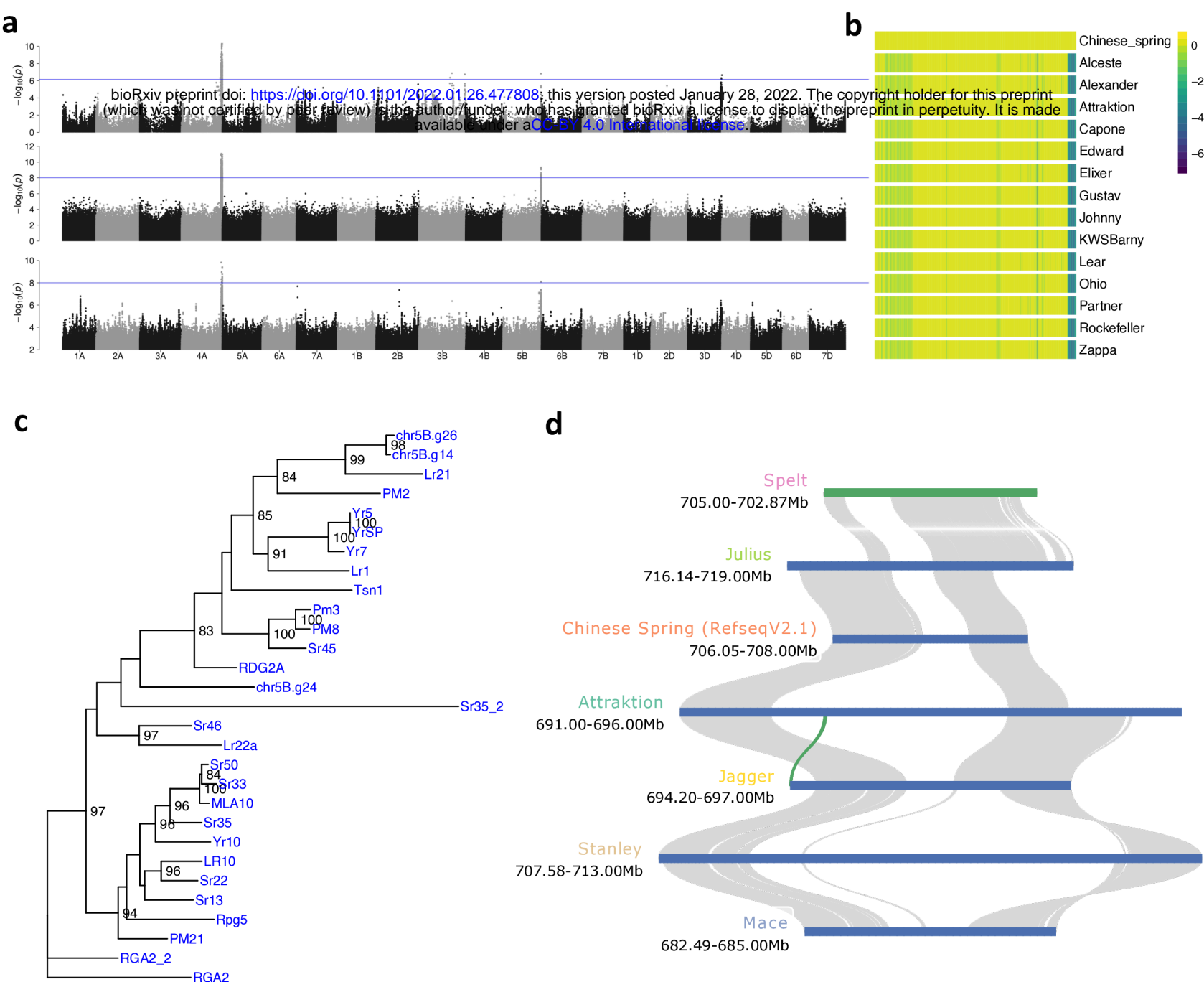
**Figure 8: Identification of a leaf rust resistance locus. (a)** Manhattan plots showing GWAS results for leaf rust resistance based on SNPs (top row), SNPs scored as presence-absence markers (middle row) and *k*-mer markers. SNPs and *k*-mers were anchored to the reference sequence assembly of cv. Attraktion. Two regions at the distal ends of the long arms of chromosomes 4A and 5B were associated with leaf rust resistance. **(b)** Normalized read depth in 500 kb bins along chromosome 4A of Chinese Spring RefSeq V1.0 for representative elite varieties. **(c)** Phylogenetic tree constructed with NBS-LRR genes from the 5B locus together with cloned resistance genes of wheat. Two genes from the locus are highly homologous to *Lr21*. **(d)** Gene based-colinearity analysis of the Attraktion haplotype at the 5B locus with assemblies from the wheat pan-genome (Walkowiak et al. 2020), none of which carry the Attraktion haplotype.