

1 **Unlocking the diversity of genebanks: whole-genome** 2 **marker analysis of Swiss bread wheat and spelt**

3 Thomas Müller¹, Beate Schierscher-Viret², Dario Fossati², Cécile Brabant², Arnold
4 Schori², Beat Keller^{1*} and Simon G. Krattinger^{1,3*}

5 ¹Department of Plant and Microbial Biology, University of Zurich, Zurich, Switzerland

6 ²Department of Plant Production Sciences, Agroscope, Changins, Switzerland

7 ³Current address: King Abdullah University of Science and Technology (KAUST),
8 Biological and Environmental Sciences & Engineering Division (BESE), Thuwal,
9 23955-6900, Saudi Arabia

10

11 *Correspondence:

12 Beat Keller, E-mail: bkeller@botinst.uzh.ch, telephone: +41 44 6348230

13 Simon G. Krattinger, E-mail: simon.krattinger@kaust.edu.sa, telephone: +966 12
14 8082668

15 E-mail addresses:

16 Thomas Müller: thomas.mueller@botinst.uzh.ch

17 Beate Schierscher-Viret: beate.schierscher-viret@agroscope.admin.ch

18 Dario Fossati: dario.fossati@agroscope.admin.ch

19 Cécile Brabant: cecile.brabant@agroscope.admin.ch

20 Arnold Schori: arnold.schori@agroscope.admin.ch

21 **Key message**

22 High-throughput genotyping of Swiss bread wheat and spelt accessions revealed
23 differences in their gene pools and identified bread wheat landraces that were not used
24 in breeding.

25 **Keywords**

26 genebank; high-resolution genotyping; landraces; bread wheat; spelt; SNP

27 **Acknowledgements**

28 We thank Simon Fluckiger and Helen Zbinden for technical assistance with DNA
29 extraction and Prof. Eduard Akhunov for providing the genotypic data of the
30 international wheat accessions. This work was supported by the Swiss Federal Office
31 for Agriculture (BLW) in the framework of NAP-PGREL (national plan of action for
32 the conservation and sustainable utilization of plant genetic resources). S.G.K is
33 supported by an Ambizione fellowship of the Swiss National Science Foundation.

34

1 **(1) Abstract**

2 Genebanks play a pivotal role in preserving the genetic diversity present among old
3 landraces and wild progenitors of modern crops and they represent sources of
4 agriculturally important genes that were lost during domestication and in modern
5 breeding. However, undesirable genes that negatively affect crop performance are often
6 co-introduced when landraces and wild crop progenitors are crossed with elite cultivars,
7 which often limits the use of genebank material in modern breeding programs. A
8 detailed genetic characterization is an important prerequisite to solve this problem and
9 to make genebank material more accessible to breeding. Here, we genotyped 502 bread
10 wheat and 293 spelt accessions held in the Swiss National Genebank using a 15K wheat
11 SNP array. The material included both spring and winter wheats and consisted of old
12 landraces and modern cultivars. Genome- and sub-genome-wide analyses revealed that
13 spelt and bread wheat form two distinct gene pools. In addition, we identified bread
14 wheat landraces that were genetically distinct from modern cultivars. Such accessions
15 were possibly missed in the early Swiss wheat breeding program and are promising
16 targets for the identification of novel genes. The genetic information obtained in this
17 study is appropriate to perform genome-wide association studies, which will facilitate
18 the identification and transfer of agriculturally important genes from the genebank into
19 modern cultivars through marker-assisted selection.

20

21

22

1 **(2) Introduction**

2 Hexaploid wheat (*Triticum aestivum*) is one of the most important global crops (FAO
3 2015). The most widely cultivated subspecies is bread wheat (*T. aestivum* ssp.
4 *aestivum*). Spelt (*T. aestivum* ssp. *spelta*) is a close relative of bread wheat. It was the
5 main wheat subspecies grown in Central Europe since the Bronze Age before bread
6 wheat cultivation started to expand from the beginning of the 20th century (Akeret 2005;
7 Jacquemin 2011; Schilperoord 2013). Today, spelt is only grown as a niche product in
8 Central Europe. In contrast to free-threshing bread wheat, spelt is hulled and the kernels
9 are surrounded by tenacious glumes. Both wheat subspecies can be interbred and
10 crosses between winter wheat and spelt were systematically explored to improve the
11 agronomical value of spelt (Winzeler et al. 1991; Siedler et al. 1994). On the other hand,
12 spelt carries genes that were beneficial for bread wheat improvement (Kleijer et al.
13 2012). For example, the stripe rust resistance gene *Yr5*, first described in spelt
14 accessions, was transferred into bread wheat (Macer 1966; Sun et al. 2002). Similarly,
15 the leaf rust resistance gene *Lr65* was identified in a Swiss spelt from where it was
16 subsequently introduced into the bread wheat gene pool (Mohler et al. 2012).

17 Domestication of wild progenitors of modern wheat started in the Near East around
18 10,000 years ago (Salamini et al. 2002). Compared to wild wheat relatives and old
19 landraces, modern cultivars show a lower genetic diversity. This genetic bottleneck is
20 due to the limited number of wild wheat progenitors and landraces that gave rise to
21 modern wheat cultivars (Tanksley and McCouch 1997; Reif et al. 2005b; Feuillet et al.
22 2008; Fu and Somers 2009). Several strategies are used to counteract this problem and
23 to increase the genetic diversity in wheat breeding programs. For example, tetraploid

1 wheat can be crossed with the wild diploid D-genome progenitor *Aegilops tauschii*,
2 resulting in synthetic hexaploid wheat. This approach was widely explored by the
3 International Maize and Wheat Improvement Center (CIMMYT) and other breeding
4 programs (Mujeeb-Kazi et al. 1996; Dreisigacker et al. 2008). Another approach to
5 increase diversity in modern cultivars is through the use of the genetic diversity present
6 in landraces and wild wheat progenitors (Reynolds et al. 2006). Landraces are of
7 interest to breeders because they often carry genes with beneficial effects that were not
8 introduced into elite cultivars. For instance, the Swiss winter bread wheat landrace
9 Muenstertaler was identified as a source of resistance to snow molds (Gaudet and
10 Kozub 1991) and Swiss barley landraces from mountainous regions were identified as
11 sources of stem rust resistance (Steffenson et al. 2016).

12 The value of landraces and wild wheat progenitors has long been recognized, which
13 resulted in the systematic collection of plant genetic resources. Today, genebanks
14 worldwide maintain this agricultural diversity by storing and propagating seeds of
15 hundreds of thousands of wheat accessions (Börner et al. 2014). Hence, genebanks
16 provide an enormous resource that can be used in research and breeding to make wheat
17 more resilient to pests, diseases, or adverse climatic conditions. In Switzerland, this task
18 is managed by the Swiss National Genebank, which has been established in the early
19 1900s. Until 1950, the focus of the bread wheat collection was on Swiss landraces,
20 while spelt was also collected from Germany, Belgium, Austria, Liechtenstein, and
21 Spain. Today, the Swiss genebank contains the largest spelt collection worldwide with
22 more than 2,200 accessions (Kleijer et al. 2012). In addition, the genebank incorporates
23 more than 5,600 bread wheat accessions.

1 For the most efficient use of genebank material it is important to have detailed genetic
2 information. For example, breeders may be interested in using landraces that are
3 genetically very different from current elite cultivars. Single nucleotide polymorphisms
4 (SNPs) distributed across the entire genome represent a valuable tool to assess genome-
5 wide diversity. It is nowadays possible to detect thousands of SNPs in a large number of
6 accessions with high-throughput technologies like SNP arrays or genotyping-by-
7 sequencing (GBS) in a reasonable time (Elshire et al. 2011; McCouch et al. 2012; Wang
8 et al. 2014). SNP arrays consist of a predefined set of polymorphisms, have low error
9 rates and low computational needs during data analysis. However, data generated by
10 SNP arrays may suffer from an ascertainment bias that is caused by the selection of
11 polymorphisms during array design (Albrechtsen et al. 2010). The 15K SNP array used
12 in this study was mainly designed with polymorphisms identified in bread wheat and
13 durum wheat accessions (Wang et al. 2014). The array consists of 13,006 gene-
14 associated SNPs and was already successfully applied to genotype durum and bread
15 wheat (Merchuk-Ovnat et al. 2016).

16 Here, we used the 15K wheat SNP array to genotypically characterize a core collection
17 of 502 bread wheat and 293 spelt accessions of the Swiss National Genebank. The
18 collection represents the history of Swiss wheat breeding and farming from the early
19 20th century to today. We show that bread wheat and spelt accessions represent two
20 separate gene pools based on a large number of genome-wide markers. Based on the
21 genotypic data, we identify two groups of bread wheat landraces that are genetically
22 different from modern bread wheat cultivars. The obtained genotypic data can be used
23 to narrow down the number of accessions to be used in breeding programs or to select
24 accessions for genome-wide association studies based on their genetic diversity.

1 (3) Material and methods

2 Plant material

3 Seven-hundred-ninety-five hexaploid wheat accessions (502 *T. aestivum* ssp. *aestivum*
4 and 293 *T. aestivum* ssp. *spelta*) from the Swiss National Genebank were used in this
5 study (Online resource 1). Among the bread wheat accessions were 367 landraces, 5
6 breeder's lines, and 120 cultivars from Switzerland. The remaining accessions came
7 from Italy, Japan, USA, Russia, Mexico, Turkey (one cultivar per country), and France
8 (one landrace, three cultivars). Accessions are defined as landraces if they were
9 collected prior to 1950. Cultivars are officially registered accessions and breeder's lines
10 originated from breeding or research projects. Cultivars and breeder's lines are both
11 representing modern material and were grouped together, i.e., each bread wheat
12 accession belongs to one of the following groups: winter bread wheat cultivar, winter
13 bread wheat landrace, spring bread wheat cultivar or spring bread wheat landrace (Supp.
14 Tab. S1).

15 The separation of spelt accessions into landraces and cultivars is difficult, because
16 registered spelt cultivars were often collected before 1950 and can also be considered as
17 landraces. In addition, modern spelt lines are in general crosses of spelt with bread
18 wheat accessions. Therefore, we grouped the spelt accessions into winter spelt and
19 spring spelt (242 accessions), representing accessions collected before 1950, and into
20 spelt/wheat crosses (51 accessions) representing accessions originating from breeding or
21 research programs (Supp. Tab. S1).

22

1 **DNA extraction and genotyping**

2 DNA from one plant per accession was extracted as described previously (Stein et al.
3 2001). Accessions were genotyped with an Illumina Infinium 15K wheat SNP array
4 (TraitGenetics GmbH, Gatersleben, Germany) consisting of 13,006 SNPs. Haplotype-
5 specific SNP markers were selected by allele frequency, functionality and sub-genome
6 specificity in hexaploid wheat from the wheat 90K iSelect assay (Wang et al. 2014).
7 Eleven spelt and eleven wheat accessions were genotyped twice. Ninety-eight SNPs
8 returned missing data in all accessions and 26 bread wheat and 11 spelt accessions
9 (including one sample of a replicate) had missing data for all markers. These 37
10 accessions were excluded from the analysis. For the replicates, we kept the sample with
11 fewer missing data and fewer heterozygous SNP calls for our analyses. Finally, 283
12 spelt and 476 bread wheat accessions were further analyzed. Genotyping data are
13 deposited on the website of the Swiss National Genebank (bread wheat:
14 <https://www.bdn.ch/lists/1701/export/>, spelt: <https://www.bdn.ch/lists/1699/export/>) and
15 in Online resource 2.

16 A randomly selected set of 29 wheat and 30 spelt accessions were in addition genotyped
17 using GBS (Supp. Tab. S2; Elshire et al. 2011; Poland et al. 2012). GBS was performed
18 by the Genomic Diversity Facility at Cornell University, USA, using *PstI* as restriction
19 enzyme. SNP calling was performed using the TASSEL 5 GBS v2 Pipeline of the
20 TASSEL package (Glaubitz et al. 2014) using the TGACv1 assembly of wheat as
21 reference (Clavijo et al. 2017).

22

1 **Data analysis**

2 Most analyses were performed in Python v3.4.3 using the libraries scipy v0.17.0,
3 numpy v1.11.0 (van der Walt et al. 2011), sklearn v0.17.1 (Pedregosa et al. 2011),
4 pandas v0.18.0 (McKinney 2010), matplotlib v1.5.1 (Hunter 2007) and ipython v4.2.0
5 (Perez and Granger 2007).

6 Genetic differentiation was determined using a simple Hamming distance (Hamming
7 1950; Wang et al. 2015), Rogers' distance (Rogers 1972; Reif et al. 2005a),
8 discriminant analysis of principal components (DAPC), and the fixation index F_{ST} .
9 Because the original marker data were notated in the IUPAC notation (Cornish-Bowden
10 1985), we converted each marker data point into a two character code, e.g., 'A' to 'AA',
11 'K' to 'GT' and concatenated them in the same order for each accession to calculate the
12 Hamming distance. The distance between two accessions was calculated as the number
13 of mismatches between those converted strings. Missing data were treated like matches.
14 For the calculation of Rogers' distance missing data were imputed with the mean of all
15 alleles. Rogers' distance was then calculated with the R-package *poppr* v2.5 (Kamvar et
16 al. 2014). R package *adegenet* v2.0.1 was used for DAPC (Jombart et al. 2010). DAPC
17 combines a principal component analysis (PCA) with linear discriminant analysis.
18 While PCA is based on the variation in the whole data set, DAPC results in clustering
19 the data in a way that maximizes the variation between and minimizes variation within
20 clusters. Pairwise F_{ST} values were calculated between the seven classes of accessions
21 using the Weir-Cockerham method implemented in *vcftools* v0.1.15 (Weir and
22 Cockerham 1984; Danecek et al. 2011). F_{ST} is a measure of population differentiation
23 that, compared to PCA, is found to be less affected by a potential ascertainment bias in

1 SNP array data (Albrechtsen et al. 2010). Principal coordinate analyses (PCoA) based
2 on pairwise Rogers' distances and mean pairwise F_{ST} values were performed with R-
3 package *ape* v3.5 (Paradis et al. 2004). Nei's G_{ST} was calculated using *vcfR* (Knaus and
4 Grünwald 2017).

5 Linkage disequilibrium (LD) of SNPs with a minor allele frequency greater than 0.05
6 was estimated by calculating the squared correlation coefficients (r^2) between genotypes
7 with *vcftools* v0.1.15 (Danecek et al. 2011). To determine LD decay the r^2 values were
8 plotted against genetic distances and an exponential curve was fitted in the data.

9 Phylogenetic trees were constructed using R packages *adegenet* v2.0.1 and *poppr* v2.3.0
10 with 1,000 bootstrap replicates (Jombart and Ahmed 2011; Kamvar et al. 2014).

11 **Genome-wide association study**

12 A genome-wide association study (GWAS) was performed using EMMAX (Kang et al.
13 2010) and the binary trait 'type' (winter - spring) with the bread wheat accessions.
14 Heterozygous SNPs were set to missing, and SNPs with more than 20% missing data
15 were excluded. The applied genetic map consisted of 9,809 SNPs at 2,887 different
16 genetic positions (Wang et al. 2014). Only one SNP was kept, if SNPs at the same
17 genetic position had equal genotypes. Missing data was then imputed by MACH1 (Li et
18 al. 2010) and input files for EMMAX were generated with PLINK 1.9
19 (<http://pngu.mgh.harvard.edu/purcell/plink/>) (Purcell et al. 2007) filtering out SNPs
20 with minor allele frequencies below 5%. Balding-Nichols kinship matrix was used to
21 account for population structure in the GWAS. Manhattan plot was made with R

1 package *qqman* (Turner 2014). SNP reads, i.e., sequences around SNPs, were extracted
2 from Table S5 of Wang et al. (2014).

3 **(4) Results**

4 **Genetic distances reveal groups of highly similar accessions**

5 Genotyping with the 15K wheat SNP array was successful for 476 bread wheat and 283
6 spelt accessions and resulted in 12,895 polymorphic SNPs across all accession, 12,892
7 polymorphic SNPs across the bread wheat accessions alone, and 11,662 SNPs across
8 the spelt accessions alone. Missing data and heterozygous SNP calls were only slightly
9 higher in the spelt than in the bread wheat accessions (Wilcoxon rank sum test: $p <$
10 0.001 in both cases; Supp. Fig. S1). Mean heterozygosity rates in bread wheat and spelt
11 were 0.5% and 0.4% and the mean missing data rate was 0.6% in both subsets. Missing
12 data rates positively correlated with heterozygosity rates, which likely reflects problems
13 during probe annealing of these accessions rather than actual heterozygosity (Supp.
14 Fig. S1; Mengistu et al. 2016).

15 Hamming and Rogers' distances allow to calculate the dissimilarity between individual
16 accessions and consequently provide an estimate for the relatedness of accessions. The
17 mean Hamming and Rogers' distances between the bread wheat accessions (8,553.7 s.d.
18 495.7 and 0.335 s.d. 0.053, respectively) were higher than the mean distances between
19 the spelt accessions (4,747.7 s.d. 1027.2 and 0.187 s.d. 0.066, respectively). Winter and
20 spring spelt accessions were more similar to each other than to accessions resulting
21 from crosses of bread wheat and spelt or to bread wheat accessions (Supp. Fig. S2).
22 Replicates of the same accessions showed a high level of reproducibility (Supp.

1 Tab. S3). The mean Hamming distances of bread wheat and spelt replicates were 9.9
2 (s.d. 10.7) and 6 (s.d. 4), respectively. Based on these numbers, we selected a Hamming
3 distance threshold of 25 to identify highly similar accessions across the bread wheat and
4 spelt collections. This revealed 18 and 21 groups of highly similar spelt and bread
5 wheats consisting of 63 and 44 accessions, respectively (Online resource 3). Hence, the
6 fraction of highly similar accessions was estimated at 22% among the spelts and 9%
7 among the bread wheats.

8 Our results might indicate that the spelt wheat gene pool analysed in this study is
9 genetically narrower than the bread wheat gene pool, which is an observation that has
10 been made previously (Siedler et al. 1994; Bertin et al. 2001; Blatter et al. 2004).
11 Alternatively, it is possible that the lower diversity of the spelt wheat gene pool resulted
12 from an ascertainment bias that is associated with the selection of polymorphisms to
13 construct the 15K wheat SNP array. To check for a potential ascertainment bias, we
14 randomly selected and genotyped a subset of 59 bread wheat and spelt accessions using
15 GBS. In contrast to SNP arrays, GBS does not rely on a set of pre-selected SNPs and
16 consequently is less prone to an ascertainment bias comparable to SNP arrays.
17 However, GBS might introduce other biases related to the choice of the restriction
18 enzyme (Arnold et al. 2013). The GBS data confirmed that the spelt accessions showed
19 a lower genetic diversity than the bread wheat accessions (mean Hamming distances
20 wheat 4,012.82 s.d. 137.4; spelt 3,332.95 s.d. 383.42; mean Rogers' distances wheat
21 0.179 s.d. 0.016; spelt 0.143 s.d. 0.038).

22 A comparison of the minor allele frequency distribution between the 15k wheat SNP
23 array and the GBS data revealed differences for bread wheat while the minor allele

1 frequency distribution was more similar for spelt (Supp. Fig. S3). Nei's gene diversity
2 index G_{ST} was higher for the SNP array data compared to GBS data (Supp. Fig. S4),
3 indicating that the array may overestimate the diversity, probably due to the choice of
4 SNPs. On the other hand, it has been reported that GBS underestimates diversity
5 (Arnold et al. 2013). In summary, both genotyping methods revealed that the spelt gene
6 pool analyzed in this study was less diverse than the bread wheat gene pool and we
7 conclude that a potential ascertainment bias had no influence on these results.

8 **Genetic analyses revealed two distinct gene pools for bread wheat and spelt**

9 The bread wheat and spelt accessions were clearly separated in a PCA along the first
10 axis (Fig. [1a](#)). This separation of bread wheat and spelt remained even after including a
11 worldwide set of six spelt and 393 bread wheat accessions that were previously
12 genotyped with a 90K SNP array (Fig. [1b](#)) (Wang et al. 2014). These data indicate that
13 the separation was not due to the narrow geographic distribution of the bread wheat
14 accessions in the Swiss genebank. Spelt/wheat crosses, i.e., spelt accessions resulting
15 from breeding programs that carry introgressions of bread wheat, located between the
16 bread wheat and spelt clusters. The second axis of the PCA divided the bread wheat
17 accessions into spring and winter types (Fig. [1a](#)). In addition to PCA, we performed
18 DAPC, PCoA of pairwise mean F_{ST} values of the different groups of accessions and
19 PCoA of Rogers' distances with the SNP array data. The DAPC and both PCoA
20 analyses confirmed the PCA results and revealed a clear separation of bread wheat and
21 spelt (Supp. Fig. S5, S6, and S7). In addition, analysis of the 59 accessions genotyped
22 by GBS confirmed the separation of bread wheat and spelt, indicating that these results
23 are not caused by genotyping biases (Supp. Fig. S8).

1 Analyses of the bread wheat accessions alone revealed a separation of landraces from
2 cultivars (Fig. 2, Supp. Fig. S9, S10, and S11, Online resources 4 and 5). We identified
3 a cluster of winter bread wheat landraces and a cluster of spring bread wheat landraces
4 that were distinct from the cultivars. The accessions of the two clusters originated
5 mainly from mountainous regions and were most likely not used to generate the
6 cultivars analyzed in our set. The accessions of the two clusters are also grouped
7 together in a phylogenetic tree based on the genotypic data (Online resource 6). A PCA
8 of spelt accessions alone showed no separation between spring and winter spelt
9 accessions (Fig. 3), whereas DAPC and PCoA revealed a minor separation of winter
10 and spring types (Supp. Fig. S12, S13, and S14). Spring spelt accessions were grouped
11 together in the phylogenetic tree (Online resource 6). In summary, these results show
12 that bread wheat and spelt wheat form discrete gene pools using the 15K SNP array and
13 GBS. In addition, the identification of ‘unused’ bread wheat accessions confirms the
14 usefulness of high-throughput genotyping of genebank material for the selection of
15 accessions with potentially novel traits.

16 We also performed PCA and PCoA for the three wheat sub-genomes individually based
17 on 4,186 A, 5,418 B, and 1,342 D genome-specific SNPs. The results for the sub-
18 genomes were similar to the results using the entire SNP set (Supp. Fig. S15, Supp.
19 Fig. S16). The separation of spelt and bread wheat was observed for each of the three
20 sub-genomes. The separation between spring and winter bread wheat on the other hand
21 was only observed for the A and B but not for the D sub-genome.

22

23

1 **The SNP data set is suitable for GWAS**

2 LD decay was calculated because the extent of LD determines the number of markers
3 required for GWAS (Flint-Garcia et al. 2003; Vos et al. 2017). The LD decays in the A
4 and B sub-genome were similar whereas LD decayed slowest in the D sub-genome
5 (Supp. Fig. S17), which is in agreement with previous studies (Chao et al. 2010; Wang
6 et al. 2014). The LD threshold of $r^2=0.1$ was reached after 1.94 cM in the wheat data
7 and after 6.34 cM in the spelt data. A conservative estimation of the genetic map size of
8 wheat is ~4,000 cM based on Wang et al. (2014). The size of the map, the LD decay and
9 the number of polymorphic SNPs lead to the conclusion that it is possible to use our
10 data set for GWAS.

11 An important difference between spring and winter accessions is the vernalization
12 requirement of winter accessions, i.e., the need of exposure to a cold period to induce
13 flowering. We used the simple binary trait 'type' (spring or winter) to test whether the
14 genotypic information and the SNP-density of the 15K SNP array are indeed sufficient
15 for GWAS. The GWAS was conducted on the bread wheat accessions using 9,737
16 SNPs (Online resource 1) and returned a peak on chromosome 5A (Fig. 4), in the region
17 where the main vernalization gene *VRN1A* is located (Yan et al. 2003). The ten SNPs
18 with lowest p-values were in the range of 89.56 cM to 91.3 cM (map positions based on
19 Wang et al. (2014); Supp. Tab. S4). The extended sequence of the top SNP
20 (wsnp_AJ612027A-Ta_2_1) produced a BLAST hit on a *T. monococcum* BAC clone
21 (GenBank: AF459639) that was used for positional cloning of *VRN1*, showing the
22 proximity of the GWAS peak to *VRN1* (Yan et al. 2003). These results demonstrate that

1 the SNP density in our data set is sufficient to perform GWAS, provided phenotypic
2 data are available.

3 **(5) Discussion**

4 Naturally occurring genetic variation offers an enormous potential for crop
5 improvement. In addition, landraces preserved in genebanks represent an important
6 resource to discover and introduce novel genes into modern crop cultivars (Gaudet and
7 Kozub 1991; Steffenson et al. 2016). However, a systematic phenotypic description of
8 the many thousand accessions stored in genebanks is not feasible. It is therefore
9 important that effective choices of genebank accessions can be made by breeders. The
10 genetic characterization of genebank material, which is relatively cheap and easy, is a
11 valuable strategy to identify groups of similar accessions or to select landraces for
12 phenotypic testing (Kilian and Graner 2012; Mason et al. 2015). For example, breeders
13 might be interested in testing landraces that are very diverse from the cultivars in a
14 specific breeding program to maximize chances to identify novel genes.

15 In our study, we found clusters of bread wheat landraces that were very diverse from the
16 modern wheat cultivars. A possible explanation is that these accessions, which originate
17 from the mountainous regions Wallis and Graubünden in Switzerland, were collected in
18 1943, after the Swiss wheat breeding program started around 1900 and these accessions
19 were then not introduced into the already advanced wheat breeding program (Martinet
20 1907; Schilperoord 2006). Such accessions might be sources of genes controlling frost
21 tolerance, snow molds resistance and early maturing. A major drawback that often
22 limits the use of landraces in modern breeding is the co-introduction of undesired traits
23 (Feuillet et al. 2008). In comparison to improved cultivars, landraces often show inferior

1 yield, are tall and thus susceptible to lodging. To make landraces accessible for modern
2 breeding it is essential that desired traits can effectively be separated from genes that
3 negatively affect crop performance. A GWAS allows to identify associations between
4 traits and molecular markers among hundreds of accessions. The identified markers can
5 then be used to introduce desired genes into modern cultivars through methods such as
6 marker-assisted backcrossing, thereby breaking the linkage drag (Collard et al. 2008). A
7 GWAS relies on genome-wide distributed polymorphisms and accurate phenotypes.
8 The SNPs of the 15K SNP array are located in genic regions. In addition, wheat is self-
9 pollinating and has a larger linkage disequilibrium than other cereals (Cavanagh et al.
10 2013; Wang et al. 2014; Sukumaran et al. 2015). Those factors make the 15K wheat
11 SNP array amenable for GWAS despite the wheat genome size of ~17 Gb and the
12 relatively small SNP density of the array (1 SNP per 1.3 Mb). We showed that the
13 amount and distribution of markers in our panel was sufficient to perform reliable
14 genotype-phenotype analyses.

15 Our genetic analyses conducted with two different genotyping methods revealed that
16 European spelt represents a gene pool that is distinct from the gene pool of bread wheat
17 landraces and cultivars. This observation is consistent with previous observations
18 (Siedler et al. 1994; Blatter et al. 2004; Dvorak et al. 2012). In contrast to these previous
19 studies that only used few markers or short gene sequences, our analysis assessed the
20 diversity of bread wheat and spelt on a genome-wide level with thousands of
21 polymorphisms. Analyses of minor allele frequencies and Nei's gene diversity index
22 G_{ST} showed that the SNP array data may suffer from an ascertainment bias. However,
23 our results were consistent using different analyses, inter alia a F_{ST} -based method which
24 is reported to be less affected by a potential ascertainment bias (Albrechtsen et al.

1 2010), and an additional genotyping method. Therefore, we conclude that our data are
2 not influenced by a strong ascertainment bias that would affect our conclusions.

3 **(6) Conclusion**

4 The main task of genebanks has traditionally been the conservation of agricultural
5 diversity and genebank material was not very frequently used in breeding in the past.
6 Linkage drag and the co-introduction of undesirable genes are probably the two most
7 important reasons for the limited use of old landraces and wild wheat progenitors in
8 modern breeding. We showed that it is now feasible, with the advances in wheat
9 genomics, to genotype large collections of spelt and bread wheat accessions from a
10 genebank and to search for diverse accessions. In combination with high-precision
11 phenotyping, this genotypic information can be used to identify novel genes through
12 GWAS. These genes can then be transferred into modern cultivars through marker-
13 assisted backcrossing, thereby avoiding linkage drag. These genomic advances will help
14 to transform genebanks from 'storage facilities' into active reservoirs for plant breeding.

15

16 **Conflict of interest**

17 The authors declare that they have no conflict of interest.

18 **Ethical standards**

19 The authors declare that this study complies with the current laws of the countries in
20 which the experiments were performed.

1 **Author's contributions**

2 SGK, BSV, and BK designed the experiment. SGK, DF, CB, and AS provided seeds of
3 accessions and information about these accessions. TM analyzed the data. Manuscript
4 was written by TM, SGK, and BK. All authors read and approved the final manuscript.

5 **(7) References**

- 6 Akeret Ö (2005) Plant remains from a Bell Beaker site in Switzerland, and the
7 beginnings of *Triticum spelta* (spelt) cultivation in Europe. *Veg Hist Archaeobot*
8 14:279–286. doi: 10.1007/s00334-005-0071-1
- 9 Albrechtsen A, Nielsen FC, Nielsen R (2010) Ascertainment biases in SNP chips affect
10 measures of population divergence. *Mol Biol Evol* 27:2534–47. doi:
11 10.1093/molbev/msq148
- 12 Arnold B, Corbett-Detig RB, Hartl D, Bomblies K (2013) RADseq underestimates
13 diversity and introduces genealogical biases due to nonrandom haplotype
14 sampling. *Mol Ecol* 22:3179–3190. doi: 10.1111/mec.12276
- 15 Bertin P, Grégoire D, Massart S, de Froidmont D (2001) Genetic diversity among
16 European cultivated spelt revealed by microsatellites. *Theor Appl Genet* 102:148–
17 156. doi: 10.1007/s001220051630
- 18 Blatter RHE, Jacomet S, Schlumbaum A (2004) About the origin of European spelt
19 (*Triticum spelta* L.): Allelic differentiation of the HMW Glutenin B1-1 and A1-2
20 subunit genes. *Theor Appl Genet* 108:360–367. doi: 10.1007/s00122-003-1441-7
- 21 Börner A, Landjeva S, Nagel M, et al (2014) Plant genetic resources for food and
22 agriculture (PGRFA) maintenance and research. *Genet Plant Physiol* 4:13–21.

- 1 Cavanagh CR, Chao S, Wang S, et al (2013) Genome-wide comparative diversity
2 uncovers multiple targets of selection for improvement in hexaploid wheat
3 landraces and cultivars. *Proc Natl Acad Sci* 110:8057–8062. doi:
4 10.1073/pnas.1217133110
- 5 Chao S, Dubcovsky J, Dvorak J, et al (2010) Population- and genome-specific patterns
6 of linkage disequilibrium and SNP variation in spring and winter wheat (*Triticum*
7 *aestivum* L.). *BMC Genomics* 11:727. doi: 10.1186/1471-2164-11-727
- 8 Clavijo BJ, Venturini L, Schudoma C, et al (2017) An improved assembly and
9 annotation of the allohexaploid wheat genome identifies complete families of
10 agronomic genes and provides genomic evidence for chromosomal translocations.
11 *Genome Res* 27:885–896. doi: 10.1101/gr.217117.116
- 12 Collard BCY, Mackill DJ, B PTRS (2008) Marker-assisted selection: An approach for
13 precision plant breeding in the twenty-first century. *Phil Trans R Soc B* 363:557–
14 572. doi: 10.1098/rstb.2007.2170
- 15 Cornish-Bowden A (1985) Nomenclature for incompletely specified bases in nucleic
16 acid sequences: Recommendations 1984. *Nucleic Acids Res* 13:3021–30. doi:
17 doi.org/10.1093/nar/13.9.3021
- 18 Danecek P, Auton A, Abecasis G, et al (2011) The variant call format and VCFtools.
19 *Bioinformatics* 27:2156–2158. doi: 10.1093/bioinformatics/btr330
- 20 Dreisigacker S, Kishii M, Lage J, et al (2008) Use of synthetic hexaploid wheat to
21 increase diversity for CIMMYT bread wheat improvement. *Aust J Agric Res*
22 59:413. doi: 10.1071/AR07225
- 23 Dvorak J, Deal KR, Luo MC, et al (2012) The origin of spelt and free-threshing

1 hexaploid wheat. *J Hered* 103:426–441. doi: 10.1093/jhered/esr152

2 Elshire RJ, Glaubitz JC, Sun Q, et al (2011) A robust, simple genotyping-by-sequencing
3 (GBS) approach for high diversity species. *PLoS One* 6:e19379. doi:
4 10.1371/journal.pone.0019379

5 FAO (2015) FAO statistical pocketbook. Food and Agriculture Organization of the
6 United Nations, Rome

7 Feuillet C, Langridge P, Waugh R (2008) Cereal breeding takes a walk on the wild side.
8 *Trends Genet* 24:24–32. doi: 10.1016/j.tig.2007.11.001

9 Flint-Garcia SA, Thornsberry JM, Buckler ES (2003) Structure of Linkage
10 Disequilibrium in Plants. *Annu Rev Plant Biol* 54:357–374. doi:
11 10.1146/annurev.arplant.54.031902.134907

12 Fu YB, Somers DJ (2009) Genome-wide reduction of genetic diversity in wheat
13 breeding. *Crop Sci* 49:161–168. doi: 10.2135/cropsci2008.03.0125

14 Gaudet DA, Kozub GC (1991) Screening winter wheat for resistance to cottony snow
15 mold under controlled conditions. *Can J Plant Sci* 71:957–966. doi:
16 10.4141/cjps91-138

17 Glaubitz JC, Casstevens TM, Lu F, et al (2014) TASSEL-GBS: A high capacity
18 genotyping by sequencing analysis pipeline. *PLoS One*. doi:
19 10.1371/journal.pone.0090346

20 Hamming RW (1950) Error detecting and error correcting codes. *Bell Syst Tech J*
21 29:147–160. doi: 10.1002/j.1538-7305.1950.tb00463.x

22 Hunter JD (2007) Matplotlib: A 2D graphics environment. *Comput Sci Eng* 9:90–95.
23 doi: 10.1109/MCSE.2007.55

1 Jacquemin JM (2011) Wheat breeding in Belgium. In: Bonjean AP, Angus WJ, van
2 Ginkel M (eds) The world wheat book: A history of wheat breeding. Volume 2.
3 Lavoisier, Paris,

4 Jombart T, Ahmed I (2011) adegenet 1.3-1 : New tools for the analysis of genome-
5 wide SNP data. *Bioinformatics* 1–2.

6 Jombart T, Devillard S, Balloux F (2010) Discriminant analysis of principal
7 components: A new method for the analysis of genetically structured populations.
8 *BMC Genet* 11:94. doi: 10.1186/1471-2156-11-94

9 Kamvar ZN, Tabima JF, Grünwald NJ (2014) Poppr: an R package for genetic analysis
10 of populations with clonal, partially clonal, and/or sexual reproduction. *PeerJ*
11 2:e281. doi: 10.7717/peerj.281

12 Kang HM, Sul JH, Service SK, et al (2010) Variance component model to account for
13 sample structure in genome-wide association studies. *Nat Genet* 42:348–54. doi:
14 10.1038/ng.548

15 Kilian B, Graner A (2012) NGS technologies for analyzing germplasm diversity in
16 genebanks. *Brief Funct Genomics* 11:38–50.

17 Kleijer G, Schori A, Schierscher B (2012) Die nationale Genbank von Agroscope ACW
18 gestern, heute und morgen. *Agrar Schweiz* 3:408–413.

19 Knaus BJ, Grünwald NJ (2017) vcfr : a package to manipulate and visualize variant call
20 format data in R. *Mol Ecol Resour* 17:44–53. doi: 10.1111/1755-0998.12549

21 Li Y, Willer CJ, Ding J, et al (2010) MaCH: Using sequence and genotype data to
22 estimate haplotypes and unobserved genotypes. *Genet Epidemiol* 34:816–34. doi:
23 10.1002/gepi.20533

- 1 Macer RCF (1966) The formal and monosomic genetic analysis of stripe rust (*Puccinia*
2 *striiformis*) resistance in wheat. In: Mackey I (ed) Proceedings of the Second
3 International Wheat Genetics Symposium, Lund, Sweden 1963. *Hereditas Suppl* 2.
4 pp 127–142
- 5 Martinet G (1907) Expériences sur la sélection des céréales.
- 6 Mason AS, Zhang J, Tollenaere R, et al (2015) High-throughput genotyping for species
7 identification and diversity assessment in germplasm collections. *Mol Ecol Resour*
8 15:1091–1101. doi: 10.1111/1755-0998.12379
- 9 McCouch SR, McNally KL, Wang W, Hamilton RS (2012) Genomics of gene banks: A
10 case study in rice. *Am J Bot* 99:407–23. doi: 10.3732/ajb.1100385
- 11 McKinney W (2010) Data structures for statistical computing in Python. In: Millman S
12 van der W and J (ed) Proceedings of the 9th Python in Science Conference. pp 51–
13 56
- 14 Mengistu DK, Kidane YG, Catellani M, et al (2016) High-density molecular
15 characterization and association mapping in Ethiopian durum wheat landraces
16 reveals high diversity and potential for wheat breeding. *Plant Biotechnol J*
17 14:1800–1812. doi: 10.1111/pbi.12538
- 18 Merchuk-Ovnat L, Barak V, Fahima T, et al (2016) Ancestral QTL alleles from wild
19 emmer wheat improve drought resistance and productivity in modern wheat
20 cultivars. *Front Plant Sci* 7:452. doi: 10.3389/fpls.2016.00452
- 21 Mohler V, Singh D, Singrün C, Park RF (2012) Characterization and mapping of Lr65
22 in spelt wheat “Altgold Rotkorn.” *Plant Breed* 131:252–257. doi: 10.1111/j.1439-
23 0523.2011.01934.x

- 1 Mujeeb-Kazi A, Rosas V, Roldan S (1996) Conservation of the genetic variation of
2 *Triticum tauschii* (Coss.) Schmalh. (*Aegilops squarrosa* auct. non L.) in synthetic
3 hexaploid wheats (*T. turgidum* L. s.lat. x *T. tauschii*; 2n=6x=42, AABBDD) and its
4 potential utilization for wheat improvement. *Genet Resour Crop Evol* 43:129–134.
5 doi: 10.1007/BF00126756
- 6 Paradis E, Claude J, Strimmer K (2004) APE: Analyses of phylogenetics and evolution
7 in R language. *Bioinformatics* 20:289–90. doi:
8 10.1093/BIOINFORMATICS/BTG412
- 9 Pedregosa F, Varoquaux G, Gramfort A, et al (2011) Scikit-learn: Machine learning in
10 Python. *J Mach Learn Res* 12:2825–2830.
- 11 Perez F, Granger BE (2007) IPython: A system for interactive scientific computing.
12 *Comput Sci Eng* 9:21–29. doi: 10.1109/MCSE.2007.53
- 13 Poland JA, Brown PJ, Sorrells ME, Jannink J-L (2012) Development of high-density
14 genetic maps for barley and wheat using a novel two-enzyme genotyping-by-
15 sequencing approach. *PLoS One* 7:e32253. doi: 10.1371/journal.pone.0032253
- 16 Purcell S, Neale B, Todd-Brown K, et al (2007) PLINK: A tool set for whole-genome
17 association and population-based linkage analyses. *Am J Hum Genet* 81:559–75.
18 doi: 10.1086/519795
- 19 Reif JC, Melchinger AE, Frisch M (2005a) Genetical and Mathematical Properties of
20 Similarity and Dissimilarity Coefficients Applied in Plant Breeding and Seed Bank
21 Management. *Crop Sci* 45:1. doi: 10.2135/cropsci2005.0001
- 22 Reif JC, Zhang P, Dreisigacker S, et al (2005b) Wheat genetic diversity trends during
23 domestication and breeding. *Theor Appl Genet* 110:859–864. doi: 10.1007/s00122-

1 004-1881-8

2 Reynolds M, Dreccer F, Trethowan R (2006) Drought-adaptive traits derived from
3 wheat wild relatives and landraces. *J Exp Bot* 58:177–186. doi: 10.1093/jxb/erl250

4 Rogers JS (1972) Measures of genetic similarity and genetic distance. *Stud Genet*
5 7:145–153.

6 Salamini F, Ozkan H, Brandolini A, et al (2002) Genetics and geography of wild cereal
7 domestication in the Near East. *Nat Rev Genet* 3:429–441. doi: 10.1038/nrg817

8 Schilperoord P (2013) *Kulturpflanzen in der Schweiz - Dinkel*. Verein für alpine
9 Kulturpflanzen, Alvaneu

10 Schilperoord P (2006) Die Bedeutung des Getreidearchivs der Forschungsanstalt für
11 Agrarökologie und Landwirtschaft Zürich-Reckenholz für die nationale
12 Samenbank. Arbeitsbericht III NAP 02-231.

13 Siedler H, Messmer MM, Schachermayr GM, et al (1994) Genetic diversity in European
14 wheat and spelt breeding material based on RFLP data. *Theor Appl Genet* 88:994–
15 1003. doi: 10.1007/BF00220807

16 Steffenson BJ, Solanki S, Brueggeman RS (2016) Landraces from mountainous regions
17 of Switzerland are sources of important genes for stem rust resistance in barley.
18 *Alp Bot* 126:23–33. doi: 10.1007/s00035-015-0161-3

19 Stein N, Herren G, Keller B (2001) A new DNA extraction method for high-throughput
20 marker analysis in a large-genome species such as *Triticum aestivum*. *Plant Breed*
21 120:354–356. doi: 10.1046/j.1439-0523.2001.00615.x

22 Sukumaran S, Dreisigacker S, Lopes M, et al (2015) Genome-wide association study for
23 grain yield and related traits in an elite spring wheat population grown in temperate

1 irrigated environments. *Theor Appl Genet* 128:353–363. doi: 10.1007/s00122-014-
2 2435-3

3 Sun Q, Wei Y, Ni Z, et al (2002) Microsatellite marker for yellow rust resistance gene
4 Yr5 in wheat introgressed from spelt wheat. *Plant Breed* 121:539–541. doi:
5 10.1046/J.1439-0523.2002.00754.X

6 Tanksley S, McCouch S (1997) Seed banks and molecular maps: Unlocking genetic
7 potential from the wild. *Science* (80-) 277:1063–1066. doi:
8 10.1126/science.277.5329.1063

9 Turner SD (2014) qqman: An R package for visualizing GWAS results using Q-Q and
10 manhattan plots. *bioRxiv*. doi: dx.doi.org/10.1101/005165

11 van der Walt S, Colbert SC, Varoquaux G (2011) The NumPy array: A structure for
12 efficient numerical computation. *Comput Sci Eng* 13:22–30. doi:
13 10.1109/MCSE.2011.37

14 Vos PG, Paulo MJ, Voorrips RE, et al (2017) Evaluation of LD decay and various LD-
15 decay estimators in simulated and SNP-array data of tetraploid potato. *Theor Appl*
16 *Genet* 130:123–135. doi: 10.1007/s00122-016-2798-8

17 Wang C, Kao W-H, Hsiao CK, et al (2015) Using Hamming distance as information for
18 SNP-sets clustering and testing in disease association studies. *PLoS One* 10:e
19 0135918. doi: 10.1371/journal.pone.0135918

20 Wang S, Wong D, Forrest K, et al (2014) Characterization of polyploid wheat genomic
21 diversity using a high-density 90 000 single nucleotide polymorphism array. *Plant*
22 *Biotechnol J* 12:787–796. doi: 10.1111/pbi.12183

23 Weir BS, Cockerham CC (1984) Estimating F-statistics for the analysis of population

1 structure. *Evolution* (N Y) 38:1358. doi: 10.2307/2408641

2 Winzeler H, Schmid J, Winzeler M, Rügger A (1991) Neue Aspekte der
3 Dinkelzüchtung (*Triticum spelta* L.) in der Schweiz. In: 2. Hohenheimer
4 Dinkelkolloquium, Universität Hohenheim. pp 11–25

5 Yan L, Loukoianov A, Tranquilli G, et al (2003) Positional cloning of the wheat
6 vernalization gene *VRN1*. *Proc Natl Acad Sci U S A* 100:6263–8. doi:
7 10.1073/pnas.0937399100

8

9

10

11

12

13

14

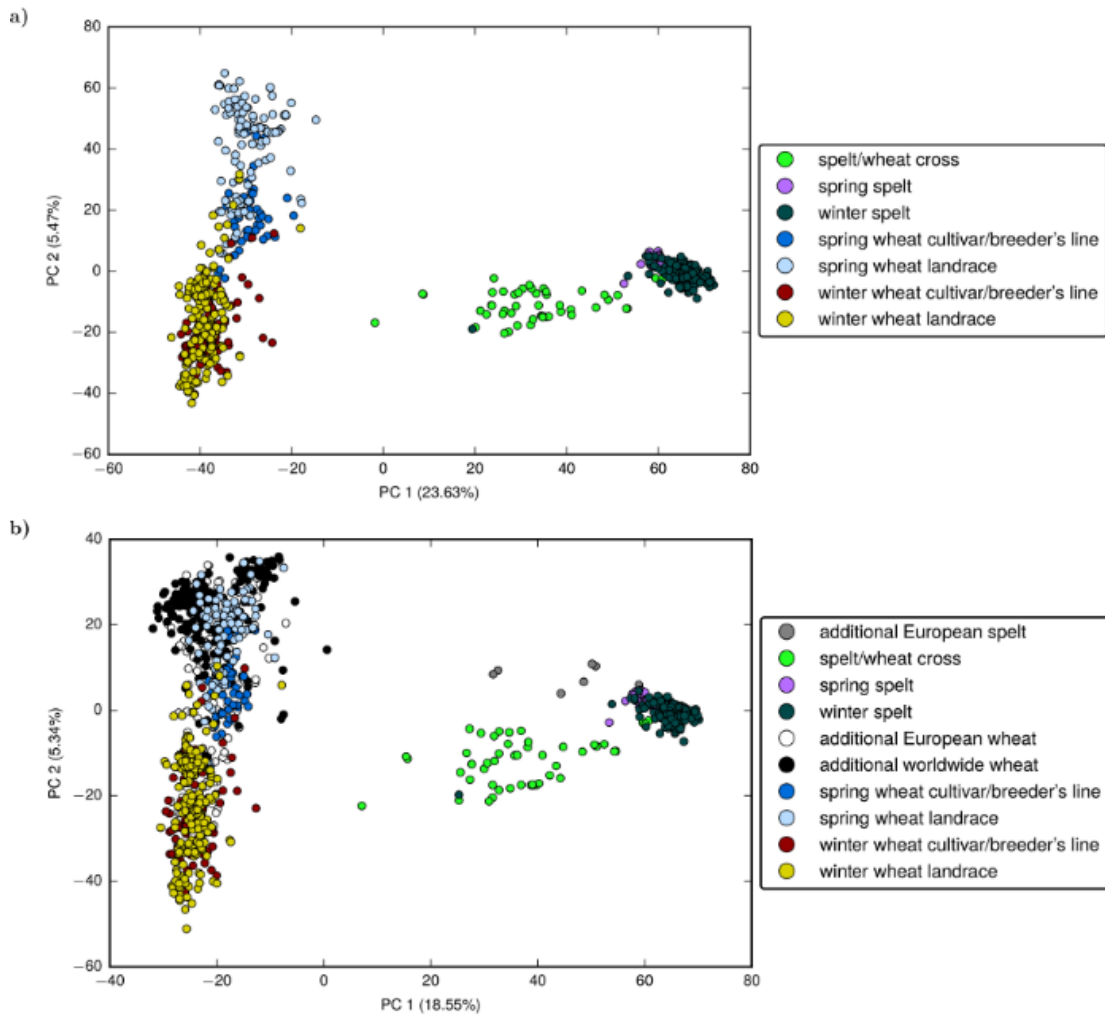
15

16

17

18

1 (8) Figure legends



2

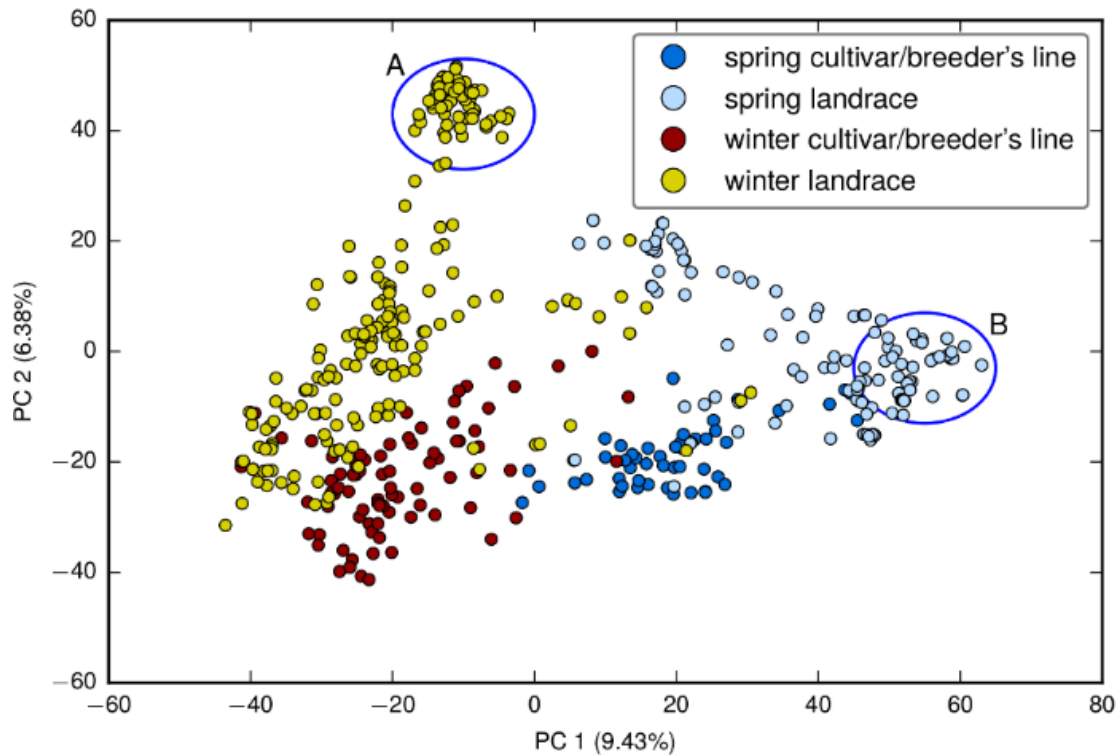
3 **Figure 1:** Principal component analysis of bread wheat and spelt accessions. a) Bread
4 wheat (left) and spelt (right) accessions are separated in Swiss material. b) PCA
5 including a set of worldwide bread wheat and spelt accessions. 9,991 SNPs that were
6 present in both the 90K SNP and the 15K SNP array were considered. Each dot
7 represents an accession according to the color coding giving in the legend.

8

9

10

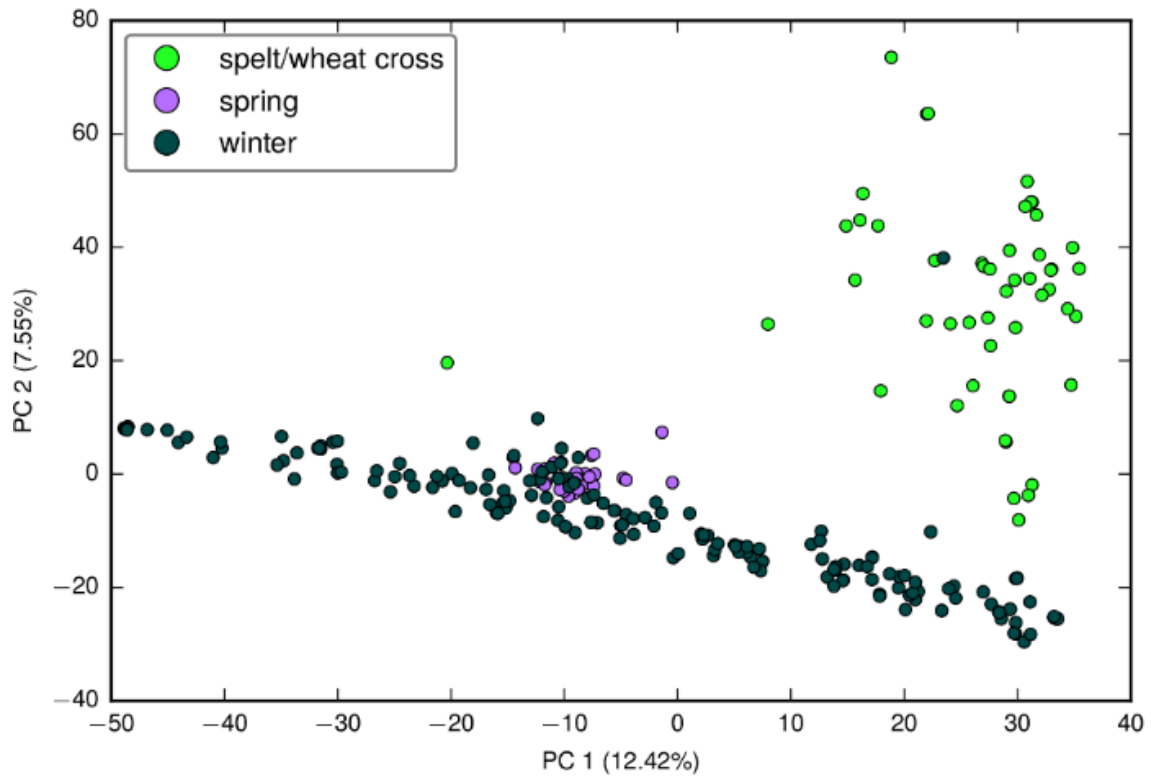
11



1

2 **Figure 2:** Principal component analysis of bread wheat accessions alone. Two clusters
 3 of landraces originating mainly from the Wallis region (winter accessions, circle A,
 4 Online resource 4) and from the Wallis and Graubünden regions (spring accessions,
 5 circle B, Online resource 5) show no (A) or only little (B) overlap to cultivars.

6



1

2 **Figure 3:** Principal component analysis of spelt accessions alone.

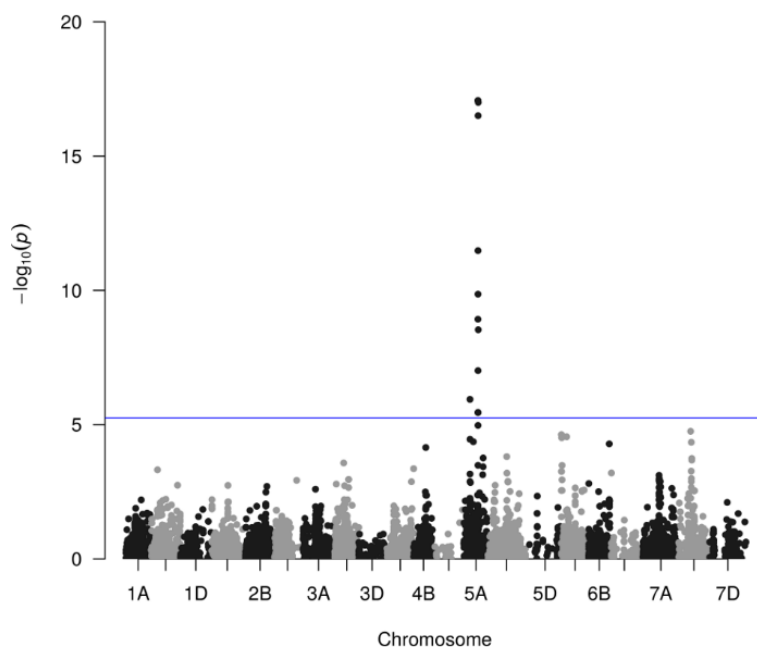
3

4

5

6

7



1

2 **Figure 4:** Manhattan plot of GWAS using EMMAX with trait 'type' (winter or spring
 3 bread wheat). The blue line corresponds to a p-value of 5.14×10^{-06} (Bonferroni
 4 correction at a significance level of 0.05).

5

6 **(9) Supporting information**

7 **Online resource 1 — Accession and passport data**

8 Additional information about the accessions used in this study.

9 **Online resource 2 — Genotypes**

10 Genotypes of accessions.

11 **Online resource 3 — Similar accessions**

12 Additional information about highly similar accessions.

1 **Online resource 4 — Winter bread wheat landraces with no overlap with cultivars**

2 Passport data of winter bread wheat landraces of circle A in Figure 2.

3 **Online resource 5 — Spring bread wheat landraces with only little overlap with**
4 **cultivars**

5 Passport data of spring bread wheat landraces of circle B in Figure 2.

6 **Online resource 6 — Phylogenetic tree**

7 Phylogenetic tree of all accessions. Tree was constructed using R packages *adegenet*
8 and *poppr*. Color codes: purple: spring spelt, dark green: winter spelt, light green:
9 wheat/spelt cross, yellow: wheat winter landrace, red: wheat winter cultivar, light blue:
10 wheat spring landrace, dark blue: wheat spring cultivar.

11 **Online resource 7 — Supporting figures and tables**

12 Supporting figures S1-S17 and supporting tables S1-S4.

13