

1

# Finding genetic variants in plants without complete

2

# genomes

3 Yoav Voichek, Detlef Weigel\*

4 Department of Molecular Biology, Max Planck Institute for Developmental Biology, 72076 Tübingen, Germany

5 \*for correspondence: [weigel@weigelworld.org](mailto:weigel@weigelworld.org)

6

## Abstract

7 Structural variants and presence/absence polymorphisms are common in plant genomes, yet they are  
8 routinely overlooked in genome-wide association studies (GWAS). Here, we expand the genetic variants  
9 detected in GWAS to include major deletions, insertions, and rearrangements. We first use raw  
10 sequencing data directly to derive short sequences, *k*-mers, that mark a broad range of polymorphisms  
11 independently of a reference genome. We then link *k*-mers associated with phenotypes to specific  
12 genomic regions. Using this approach, we re-analyzed 2,000 traits measured in *Arabidopsis thaliana*,  
13 tomato, and maize populations. Associations identified with *k*-mers recapitulate those found with  
14 single-nucleotide polymorphisms (SNPs), however, with stronger statistical support. Moreover, we  
15 identified new associations with structural variants and with regions missing from reference genomes.  
16 Our results demonstrate the power of performing GWAS before linking sequence reads to specific  
17 genomic regions, which allow detection of a wider range of genetic variants responsible for phenotypic  
18 variation.

19

## Introduction

20 Elucidating the link between genotype and phenotype is central to biological research, both in basic  
21 research as well as in translational medicine and agriculture. Correlating genotypic and phenotypic  
22 variability in genome-wide association studies (GWAS) has become the tool of choice for systematic  
23 identification of candidate loci in the genome that are causal for phenotypic differences. In plants, many  
24 species-centered projects are genotyping collections of individuals, for which different phenotypes can  
25 then be measured and analyzed. These include hundreds or thousands of strains from *Arabidopsis*  
26 *thaliana*, rice, maize, tomato, sunflower, and several other species (1001 Genomes Consortium, 2016;  
27 Bukowski et al., 2018; Hübner et al., 2018; Tieman et al., 2017; Wang et al., 2018).

28 A difficulty of working with plant genomes is that they are highly repetitive and feature excessive  
29 structural variation between members of the same species, mostly attributed to their active transposons  
30 (Bennetzen, 2000). For example, in the well-studied species *Arabidopsis thaliana*, natural accessions are  
31 missing 15% of the reference genome, indicating a similar fraction would be absent from the reference,  
32 but present in other accessions (1001 Genomes Consortium, 2016). Moreover, although *A. thaliana* has  
33 a small (140 Mb) and not very repetitive genome compared to many other plants, SNPs may be assigned  
34 to incorrect positions due to sequence similarity shared between unlinked loci (Long et al., 2013). The  
35 picture is even more complicated in other plant species, such as maize. The maize 2.3 Gb genome is  
36 highly repetitive, with transposons often inserted into other transposons, and 50%-60% of short read  
37 sequences can not be mapped uniquely to it, making the accurate identification of variants in the  
38 population a formidable challenge (Bukowski et al., 2018; Schnable et al., 2009). Furthermore, about 30%  
39 of low-copy genes present in the entire population are not found in the reference (Gore et al., 2009;  
40 Springer et al., 2018; Sun et al., 2018). Presence of large structural variants are ubiquitous all over the  
41 plant kingdom, and there are many examples for their effects on phenotypes (Saxena et al., 2014). The  
42 importance of structural variants in driving phenotypic variation has been appreciated from the early  
43 days of maize genetics (McClintock, 1950), though searching for them systematically is still an unsolved  
44 problem.

45 Correlating phenotypic and genotypic variation in GWAS is critically dependent on the ability to  
46 call individual genotypes. While short sequencing reads aligned to a reference genome can identify  
47 variants smaller than read length, such as SNPs and short indels, this approach is much less effective for  
48 larger structural variants. Moreover, variants such as SNPs can be in regions missing from the reference  
49 genome, which is frequently the case in plants. Organellar genomes are a special case, being left out of  
50 GWAS systematically although their genetic variation was shown to have strong phenotypic effects

51 (Davila et al., 2011; Joseph et al., 2013). Although not regularly used, short read sequencing can provide,  
52 in principle, information for many more variants in their source genomes than only SNPs and short  
53 indels (Iqbal et al., 2012).

54 While variants are typically discovered with short reads by mapping them to a target reference  
55 genome, one can also directly compare common subsequences among samples (Zielezinski et al., 2019).  
56 Such a direct approach is intuitively most powerful when the reference genome assembly is poor, or  
57 even non-existent. Because short reads result from random shearing of genomic DNA, and because they  
58 contain sequencing errors, comparing short reads between two samples directly is, however, not very  
59 effective. Instead, genetic variants in a population can be discovered by focusing on sequences of  
60 constant length  $k$  that are even shorter than typical short reads, termed  $k$ -mers. After  $k$ -mers have been  
61 extracted from all short reads, sets of  $k$ -mers present in different samples can be compared. Importantly,  
62  $k$ -mers present in some samples, but missing from others, can identify a broad range of genetic variants.  
63 For example, two genomes differing in a SNP (Fig. 1A) will have  $k$   $k$ -mers unique to each genome; this is  
64 true even if the SNP is found in a repeated region or a region not found in the reference genome.  
65 Structural variants, such as large deletions, inversions, translocations, transposable element (TE)  
66 insertion, etc. will also leave marks in the presence or absence of  $k$ -mers (Fig. 1A). Therefore, instead of  
67 defining genetic variants in a population relative to a reference genome, a  $k$ -mer presence/absence in raw  
68 sequencing data can be directly associated with phenotypes to enlarge the tagged genetic variants in  
69 GWAS (Lees et al., 2016).

70 Reference-free GWAS based on  $k$ -mers has been used for mapping genetic variants in bacteria,  
71 where each strain contains only a fraction of the genes present in the pan-genome (Lees et al., 2016,  
72 2017; Sheppard et al., 2013). This approach, not centered around one specific reference genome, can  
73 identify biochemical pathways associated with, for example, pathogenicity. This approach has also been  
74 applied in humans, where the number of unique  $k$ -mers is much higher than in bacterial strains, due to  
75 their larger genome (Rahman et al., 2018). However, this was restricted to case-control situations, and  
76 due to high computational load, population structure was corrected only for a subset of  $k$ -mers.

77 While  $k$ -mer based approaches are likely to be especially appropriate for plants, the large  
78 genome sizes, highly structured populations, and excessive genetic variation (Gordon et al., 2017; Minio  
79 et al., 2019; Sun et al., 2018) limit the application of previous  $k$ -mer methods to plants. A first attempt to  
80 nevertheless use  $k$ -mer based methods has recently been made in plants, but was limited to a small  
81 subset of the genome, and also accounting for population structure only for a small subset of all  $k$ -mers  
82 (Arora et al., 2019).

83 Here, we present an efficient method for  $k$ -mer-based GWAS and compare it directly to the  
84 conventional SNP-based approach on more than 2,000 phenotypes from three plant species with  
85 different genome and population characteristics - *A. thaliana*, maize and tomato. Most variants identified  
86 by SNPs can be detected with  $k$ -mers (and vice versa), but  $k$ -mers having stronger statistical support.  
87 For  $k$ -mer-only hits, we demonstrate how different strategies can be used to infer their genomic  
88 context, including large structural variants, sequences missing from the reference genome, and  
89 organeller variants. Lastly, we compute population structure directly from  $k$ -mers, enabling the analysis  
90 of species with poor quality or without a reference genome. In summary, we have inverted the  
91 conventional approach of building a genome, using it to find population variants, and only then  
92 associating variants with phenotypes. In contrast, we begin by associating sequencing reads with  
93 phenotypes, and only then infer the genomic context of these sequences. We posit that this change of  
94 order is especially effective in plant species, for which defining the full population-level genetic variation  
95 based on reference genomes remains highly challenging.

96

## Results

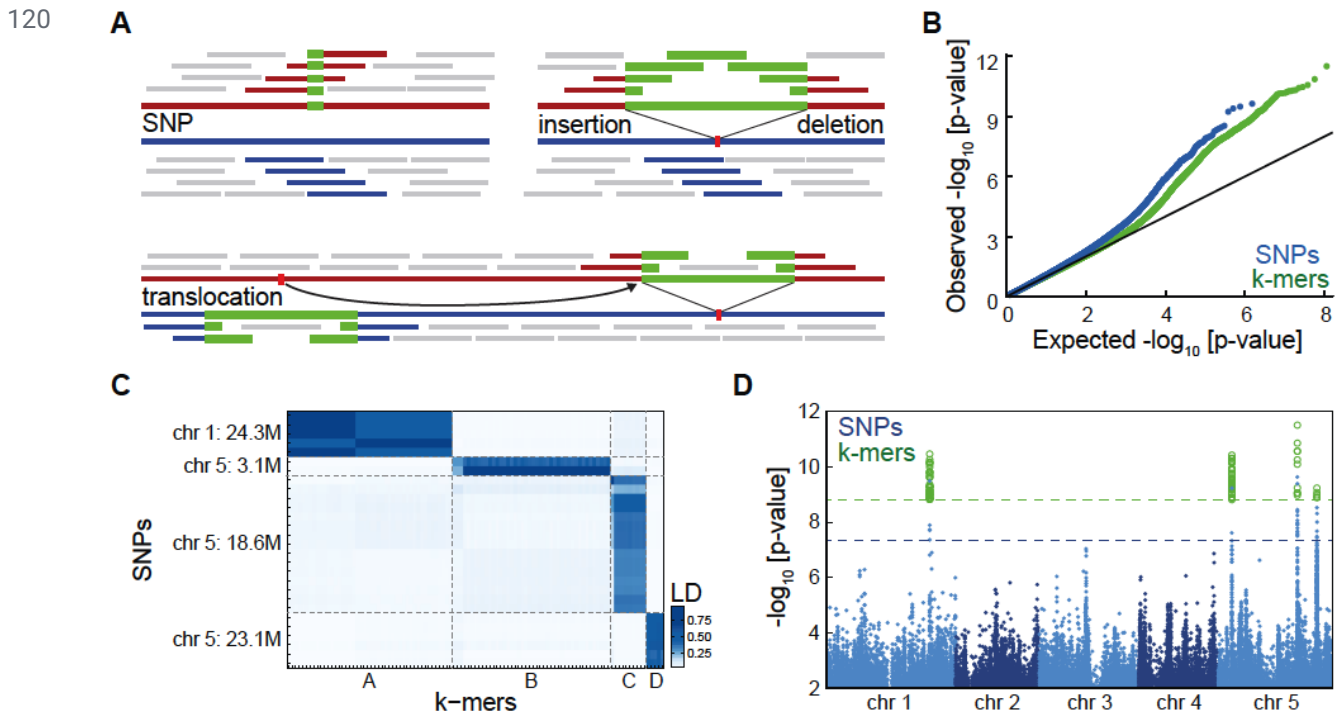
97

### Proof of concept: genetic variants for flowering of *A. thaliana*

98 As an initial proof of concept, we looked at the well-studied and well-understood trait in the model  
99 plant *A. thaliana*, flowering time. In *A. thaliana*, GWAS approaches have been used for almost 15 years  
100 (Aranzana et al., 2005), and 1,135 individuals, termed accessions, had their entire genomes resequenced  
101 several years ago (1001 Genomes Consortium, 2016). We used this genomic dataset to define the  
102 presence/absence patterns of 31 bp  $k$ -mers in these accessions (Fig. S1A). In order to minimize the  
103 effect of sequencing errors, for each DNA-Seq dataset we only considered  $k$ -mers appearing at least  
104 thrice. Out of a total 2.26 billion unique  $k$ -mers across the entire population, 439 million appeared in at  
105 least five accessions (Fig. S2A). These  $k$ -mers were not shared by all accessions, and we used the  
106 presence or absence of a  $k$ -mer as two alleles per variant to perform GWA with a linear mixed model  
107 (LMM) to account for population structure (Fig. S1B) (Zhou and Stephens, 2012). For comparison  
108 purposes, GWA was performed also with SNPs and short indels. In both cases statistically significant  
109 associations were detected (Fig. 1B).

110 To define a set of  $k$ -mers most likely to be associated with flowering time, we had to set a  
111 p-value threshold. A complication in defining such a threshold is that  $k$ -mers are often not independent,  
112 as a single genetic variant is typically tagged by several  $k$ -mers (Fig. 1A). For example, 180 million  $k$ -mers  
113 had a minor allele frequency above 5%, but these represented only 110 million unique presence/absence

114 patterns across accessions. Thus, a Bonferroni correction based on the number of all tests would be  
 115 inaccurate, as it would not accurately reflect the effective number of independent tests. To define a  
 116 threshold that accounts for the dependencies between  $k$ -mers we therefore used permutation of the  
 117 phenotype (Abney, 2015). This approach presents a computational challenge, as the full GWA analysis  
 118 has to be run multiple times. To this end, we implemented a LMM-based GWA specifically optimized for  
 119 the  $k$ -mer application (Fig. S1C) (Loh et al., 2015; Svishcheva et al., 2012).



121 **Figure 1. Flowering time associations in *A. thaliana***

122 **(A)** Presence and absence of  $k$ -mers marks a range of different genetic variants. Blue and red lines represent two  
 123 individuals genomes, and short bars above/below mark in color the  $k$ -mers unique to each genome due to genomic  
 124 differences or in grey ones shared between genomes.

125 **(B)** P-values quantile-quantile plot of SNPs and  $k$ -mers associations with flowering time measured in  $10^{\circ}\text{C}$ .  
 126 Deviation from the black line ( $y=x$ ) indicate stronger associations than expected by chance.

127 **(C)** LD (expressed as  $r^2$ ) between all SNPs and  $k$ -mers passing the p-value thresholds. Four highly linked families of  
 128 variants were identified with both methods. For SNP-to-SNP and  $k$ -mer-to- $k$ -mer LD, see Fig. S2B,C.

129 **(D)** Manhattan plot showing p-values of all SNPs (blue) and of the subset of  $k$ -mers passing the p-value threshold  
 130 (green) as a function of their genomic position. Dashed lines mark the p-value thresholds for SNPs (blue) and  
 131  $k$ -mers (green).

132 We calculated the p-value thresholds for SNPs and  $k$ -mers, set to a 5% chance of getting one  
 133 false-positive. The threshold for  $k$ -mers was more stringent than the one for SNPs (35-fold), but lower  
 134 than the increase in tests number (140-fold), as expected due to the higher dependency between

135 *k*-mers. Twenty-eight SNPs and 105 *k*-mers passed their corresponding thresholds. Using LD, we linked  
136 SNPs to *k*-mers directly without locating the *k*-mers genomic locations. Four distinct families of linked  
137 genetic variants were identified in both methods, with a clear one-to-one relationship between the four  
138 sets of SNPs and the four sets of *k*-mers (Fig. 1C, Fig. S2B,C). As expected, the *k*-mers aligned to the  
139 same genomic loci as the corresponding SNPs (Fig. 1D). For validation, we ran the analysis again with a  
140 *k*-mer length of 25 bp, obtaining a very similar result (Fig. S2D). Therefore, in this case, *k*-mer based  
141 GWAS identified the same genotype-phenotype associations as detected by SNPs.

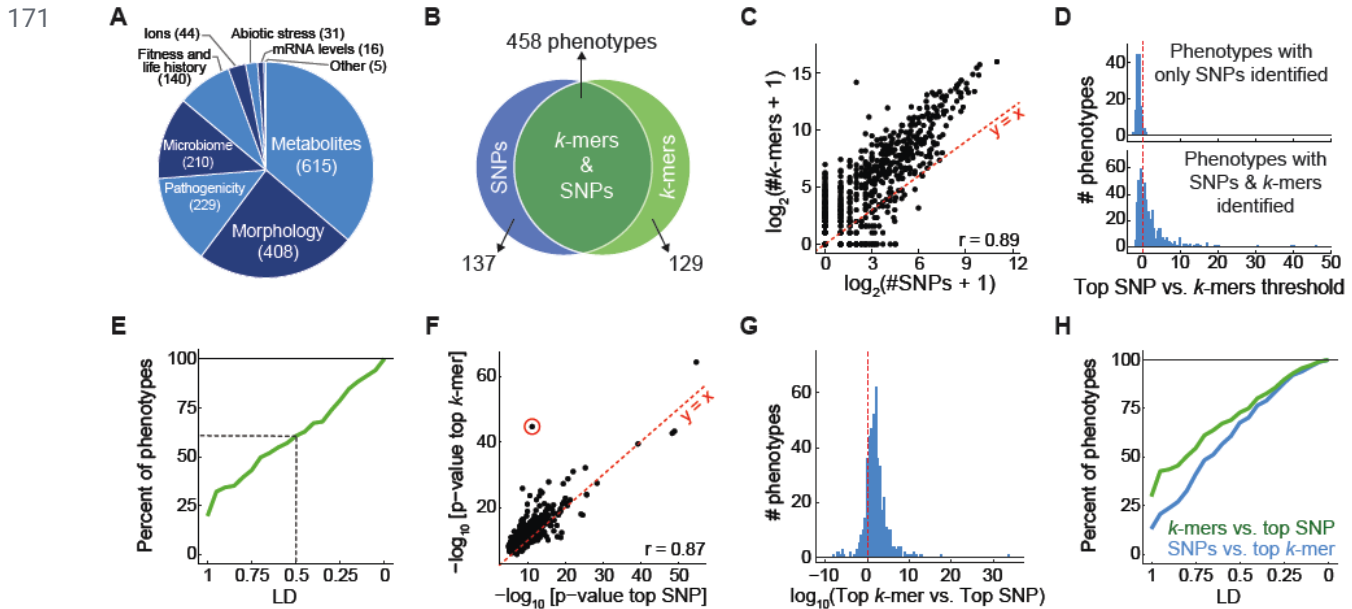
## 142 **Comparison of SNP- and *k*-mer-based GWAS on 1,697 *A. thaliana* phenotypes**

143 Flowering time is a very well studied trait, and it is unlikely that a new locus affecting it will be  
144 discovered by GWAS. To assess the potential of *k*-mer-based GWA to identify new associations, we set  
145 out to systematically compare it to the SNP-based method on a comprehensive set of traits. To this end,  
146 we collected 1,697 phenotypes from 104 *A. thaliana* studies (Table S1). This collection spans a  
147 representative sample of phenotypes regularly measured in plants (Fig. 2A). Eliminating phenotypes for  
148 which there are short read sequencing data from fewer than 40 accessions, we were left with 1,582  
149 traits to which both methods could be applied. All parameters affecting GWA analysis, such as minor  
150 allele frequency or relatedness between individuals, were the same, to obtain the most meaningful  
151 comparison. Moreover, as *A. thaliana* is a selfer, SNPs are homozygous, and their state is therefore  
152 comparable to the binary *k*-mer presence/absence.

153 We first wanted to learn whether the two methods identified similar associations. Indeed, there  
154 was substantial overlap between the traits for which associations were found (Fig. 2B). Also, the number  
155 of identified *k*-mers and SNPs per phenotype were correlated ( $r=0.89$ ), and as expected, more  
156 associated *k*-mers than SNPs were identified (Fig. 2C, Fig. S3A). For 137 phenotypes, only a significant  
157 SNP could be identified, due to the more stringent thresholds for *k*-mers, as the most significant SNPs in  
158 almost all of these phenotypes did not pass the *k*-mer threshold (Fig. 2D). Moreover, in most of these  
159 phenotypes, a *k*-mer passing the SNPs threshold was in high LD with the top SNP (Fig. 2E). Although the  
160 *k*-mer thresholds were more stringent than the SNPs thresholds (Fig. S3B), for 129 phenotypes only  
161 *k*-mers but no SNPs associations were identified. These cases were the best candidates for associations  
162 that cannot be captured with SNPs.

163 We next compared p-values of top SNPs to those of top *k*-mers; the two were correlated  
164 ( $r=0.87$ , Fig. 2F). Focusing on phenotypes for which both SNPs and *k*-mers were identified, the great  
165 majority, 86%, had stronger p-values for the top *k*-mer (Fig. 2G), a trend that had already been observed  
166 for flowering time (Fig. 1D). Lastly, we wanted to know how well top *k*-mers were tagged by significantly

167 associated SNPs and vice versa. We quantified this with the LD (as in Fig. 1C) between the top SNP and  
 168 the closest associated  $k$ -mer and the other way around. While SNPs tagged variants similar to top  
 169  $k$ -mers, associated  $k$ -mers were on average closer to top SNPs than associated SNPs to top  $k$ -mers (Fig.  
 170 2H). This was expected, as  $k$ -mers can represent SNPs but also capture other types of genetic variants.



172 **Figure 2. Comparison of SNP- and  $k$ -mer-based GWAS on 1,697 *A. thaliana* phenotypes**  
 173 (A) Assignment of 1,697 phenotypes to broad categories.  
 174 (B) Overlap between phenotypes with SNP and  $k$ -mer hits.  
 175 (C) Correlation of number of significantly associated  $k$ -mers vs. SNPs for all phenotypes.  
 176 (D) Ratios (in  $\log_{10}$ ) of top SNP p-value vs. the  $k$ -mers threshold for 137 phenotypes with only significant SNPs  
 177 (top), and for 458 phenotypes with both significant SNPs and  $k$ -mers (bottom).  
 178 (E) Fraction of phenotypes, from 137 phenotypes that had only significant SNP hits, for which a  $k$ -mer passing the  
 179 SNP threshold could be found within different LD cutoffs. For a minimum of LD=0.5 (dashed lines), 61% of  
 180 phenotypes had a linked  $k$ -mer that passed the SNP threshold.  
 181 (F) Correlation of p-values of top  $k$ -mers and SNPs for all phenotypes ( $r=0.87$ ). Red circle marks the strongest  
 182 outlier (see Fig. 3A, B for details on this phenotype).  
 183 (G) Ratio between top p-values (expressed as  $-\log_{10}$ ) in the two methods, for the 458 phenotypes with both  $k$ -mer  
 184 and SNP hits.  
 185 (H) Fraction of all phenotypes for which a significant SNP could be found within different LD cutoffs of top  $k$ -mer  
 186 (blue) and vice versa (green).

### 187 Specific case studies of $k$ -mer superiority

188 For some phenotypes,  $k$ -mers were more strongly associated with a phenotype than the top SNP,  
 189 although they represented the same variant (Fig. S4A). The goal of our study was, however, to identify  
 190 cases where  $k$ -mers provided a conceptual improvement. First, we looked into the phenotype



191 quantifying the fraction of dihydroxybenzoic acid (DHBA) xylosides among total DHBA glycosides (Li et  
192 al., 2014) (red circle in Fig. 2F). In this case, all significant *k*-mers mapped uniquely in the proximity of  
193 AT5G03490, encoding a UDP glycosyltransferase that was identified in the original study as causal (Fig.  
194 3A, Fig. S4B). The source of the stronger *k*-mers associations could be traced back to two  
195 non-synonymous SNPs, 4 bp apart, in the coding region of AT5G03490. Due to their proximity, one  
196 *k*-mer can hold the state of both SNPs, and their combined information is more predictive of the  
197 phenotype than each SNP on its own (Fig. 3B). This interaction between closely linked SNPs was not  
198 one of the types of genetic variants we had anticipated for *k*-mers.

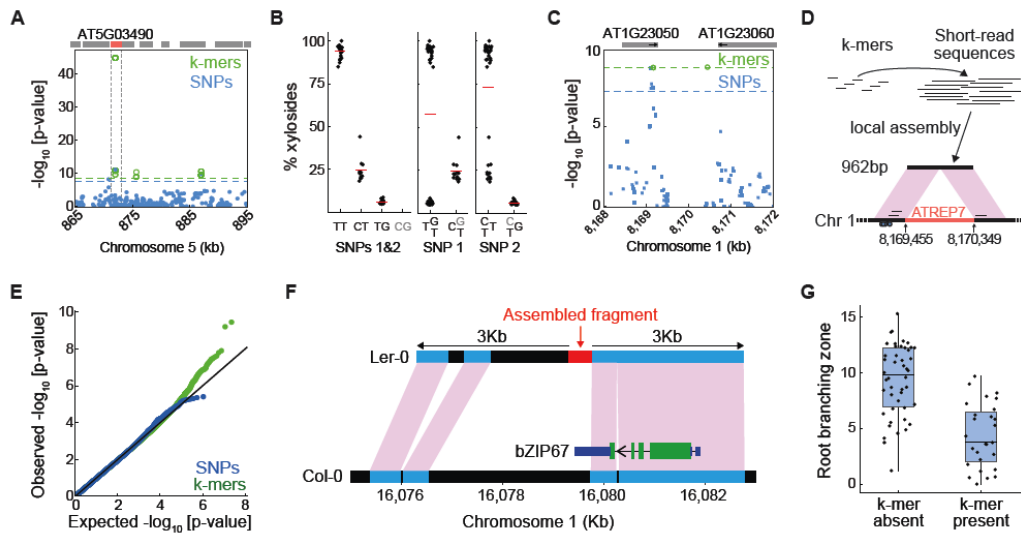
199 Our next case study involves inhibition of seedling growth in the presence of a specific *flg22*  
200 variant (Vetter et al., 2016), a phenotype for which we could map to the reference genome only three of  
201 the 10 significant *k*-mers; the three mappable *k*-mers were all located in the proximity of significant SNPs  
202 in AT1G23050 (Fig. 3C, Fig. S4C). To identify the genomic source of the remaining seven *k*-mers, we  
203 retrieved the short reads containing the *k*-mers from all relevant accessions and assembled them into a  
204 single 962 bp fragment. This fragment mapped to two genomic regions 892 bp apart, close to the three  
205 mapped *k*-mers (Fig. 3D). The junction sequence connecting the two regions could only be identified in  
206 accessions with the seven significant *k*-mers. We hypothesized that the 892 bp intervening fragment  
207 corresponds to a transposable element (TE), and a search of the Repbase database indeed identified  
208 similarity to helitron TE (Bao et al., 2015). Thus, the *k*-mers in this case marked an association with a  
209 structural variant, the presence or absence of a ~900 bp TE. While in this case the *k*-mer method did  
210 not identify a new locus, it more clearly revealed what is the likely genetic cause of variation in *flg22*  
211 sensitivity.

212 In the first two examples, hits with both *k*-mers and SNPs had been identified. Next, we looked  
213 for phenotypes for which we had only identified significant *k*-mers. One of these was germination in  
214 darkness and under low nutrient supply (Morrison and Linder, 2014). In this case, 11 *k*-mers but no  
215 significant SNPs had been found (Fig. 3E, Fig. S4D-E). However, neither the 11 *k*-mers nor the short  
216 reads they originated from could be mapped to the reference genome. The reads assembled into a 458  
217 bp fragment. A database search revealed a hit on chromosome 3 of Ler-0, a non-reference accession of  
218 *A. thaliana* with a high-quality genome assembly (Zapata et al., 2016). The flanking sequences were  
219 syntenic with region on chromosome 3 of the *A. thaliana* reference genome, with a 2 kb structural  
220 variant that included the 458 bp fragment we had assembled based on our *k*-mer hits (Fig. 3F). This  
221 variant affected the 3' untranslated region (UTR) of the bZIP67 transcription factor gene. bZIP67 acts  
222 downstream of LEC1 and upstream of DOG1, two master regulators of seed development (Bryant et al.,  
223 2019). Accumulation of bZIP67 protein but not *bZIP67* mRNA is affected by cold and thus likely



224 mediates environmental regulation of germination (Bryant et al., 2019). Structural variations in the 3'  
 225 UTR is consistent with translational regulation of bZIP67 being important. The bZIP67/germination case  
 226 study demonstrates directly the ability of our *k*-mer method to reveal associations with structural  
 227 variants that are not tagged by SNPs.

228



229

### Figure 3. Specific cases in which *k*-mers are superior to SNPs

230 (A) Associations with xyloside fraction in a region of chromosome 5. Grey boxes indicate genes with AT5G03490  
 231 marked in red.

232 (B) Xyloside fraction grouped by states at two SNPs (SNP1, 872,003 bp; SNP 2, 872,007 bp). One of the four  
 233 possible states (“CG”) does not exist, indicated in grey in left most plot, which shows grouping based on both  
 234 sites, as is possible with *k*-mers. Middle and right most plot show groupings based on only one of the two sites.

235 (C) Associations with seedling growth inhibition in the presence of flg22 near 8.17 Mb of chromosome 1. Absence  
 236 of SNPs in the central 1 kb region is likely due to the presence of a TE to which short reads cannot be  
 237 unambiguously mapped. Gene orientations indicated with short black arrows.

238 (D) Assembly of reads identified with the seven unmappable *k*-mers resulted in a 962bp fragment. This fragment  
 239 lacks the central 892 bp region in the reference genome encoding an ATREP7 helitron TE. Small circles on bottom  
 240 represent significant flanking SNPs, and short black bars above represent the three mappable significant *k*-mers.

241 (E) P-values quantile-quantile plot of associations with germination time in darkness and low nutrients. Only  
 242 *k*-mers show stronger-than-expected associations.

243 (F) Assembled reads (red bar) containing significant *k*-mers from GWA of germination time match a region on  
 244 chromosome 3 of Ler-0. Regions in addition to the red fragment that cannot be aligned to the Col-0 reference  
 245 genome are indicated in black. The 3' UTR of the gene encoding bZIP67 is indicated in dark blue. The extent of the  
 246 bZIP67 3' UTR in Ler-0 is not known. Green indicates coding sequences.

247 (G) Root branching zone length in millimeters in accessions that have the significant *k*-mer identified for this trait  
 248 and accessions that do not have this *k*-mer.

249 As a final case, we focused on the variation in the length of the root branching zone (Ristova et  
 250 al., 2018). While no significant SNPs could be identified, a single *k*-mer passed the significance threshold

251 (Fig. 3G, Fig. S4F). The  $k$ -mer and the reads containing it mapped to the chloroplast genome. When we  
252 lowered the threshold for the familywise error-rate from 5% to 10%, a second  $k$ -mer was identified,  
253 which also mapped to the chloroplast genome. Genetic variation in organelle genomes has been shown  
254 to affect phenotypic variation (Joseph et al., 2013), but they are often left out from GWA studies.

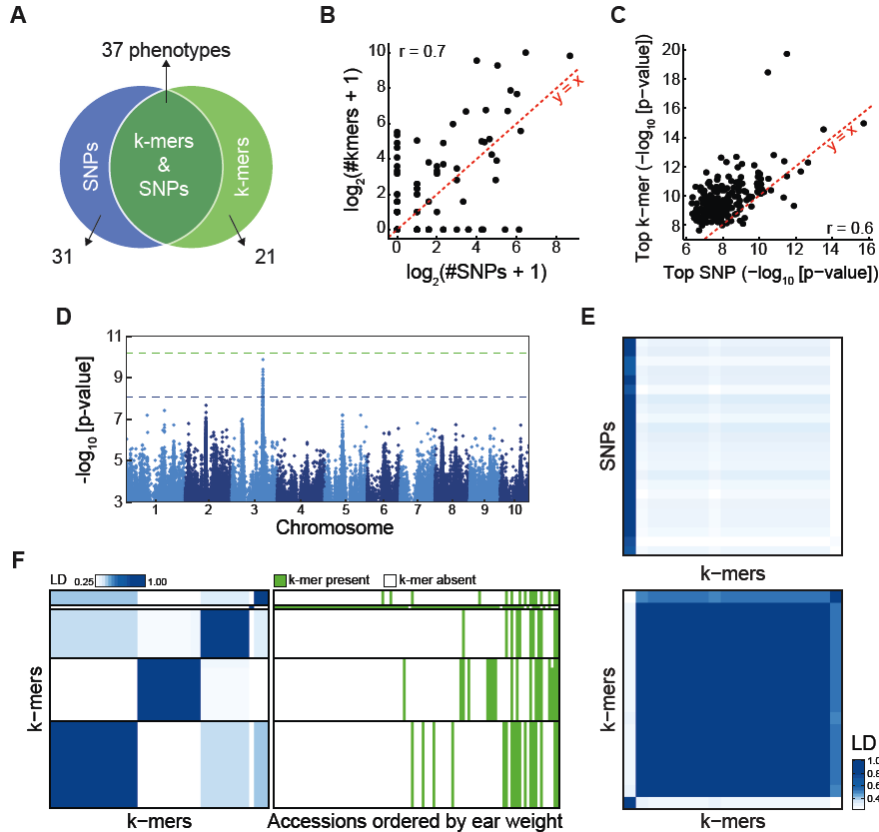
## 255 **Comparison of SNP- and $k$ -mer-based GWAS in maize**

256 While the results with *A. thaliana* were encouraging, its genome size and repeat content is not  
257 representative of many other flowering plants. We therefore wanted to evaluate our approach on larger,  
258 more complex genomes. This criterion is met by maize, with a reference genome of 2.3 Gb, ~85% of  
259 which consists of TEs and other repeats (Schnable et al., 2009). Moreover, individual maize genomes are  
260 highly divergent, with ~10% of genes being non-syntenic and many genes found in different accessions  
261 are missing from the reference genome (Gore et al., 2009; Springer et al., 2018; Sun et al., 2018).

262 We set out to apply our  $k$ -mer-based GWAS approach to a set of 150 maize inbred lines with  
263 short read sequence coverage of at least 6x (Bukowski et al., 2018). There were 7.3 billion unique  
264  $k$ -mers in the population, of which 2.3 billion were present in at least five accessions, which were used  
265 for GWAS (Fig. S5A). As in *A. thaliana*, we sought to compare the  $k$ -mer- and SNP-based approaches. To  
266 this end, we applied both methods to 252 field measurements, mostly of morphological traits (Zhao et  
267 al., 2006). For 89 traits, significant associations were identified by at least one of the methods, and for 37  
268 by both (Fig. 4A). As in *A. thaliana*, the number of statistically significant variants as well as top  
269 associations between both methods were well correlated (Fig. 4B,C). Top  $k$ -mers had lower p-values  
270 than the top SNPs (Fig. S5D), and the  $k$ -mer method detected associations not found by SNPs.

271 To discern the added benefit of the  $k$ -mer-based approach, we compared SNPs and  $k$ -mers using  
272 LD, without attempting to locate  $k$ -mers in the genome. We used this comparison approach as SNPs  
273 were originally placed on the genomic map using external information in addition to short read mapping,  
274 due to the large proportion of short reads that do not map to unique places in the reference genome  
275 (Bukowski et al., 2018). We found several cases where a  $k$ -mer marked a common allele in the  
276 population with strong effect on a phenotype, but the allele could not be identified with the SNP dataset.  
277 For example, for days to tassel there was one clear SNP hit that was also tagged by  $k$ -mers (Fig. 4D,E),  
278 but a second genetic variant was only identified with  $k$ -mers. Another example is ear weight for which  
279 no SNPs passed the significance threshold (Fig. S5F), but several unlinked variants were identified with  
280  $k$ -mers (Fig. 4F). Thus, new alleles with high predictive power for maize traits can be revealed using  
281  $k$ -mers.

282



283

**Figure 4. Comparison of SNP- and *k*-mer-based GWAS in maize**

284

(A) Overlap between phenotypes with SNP and *k*-mer hits. See also Fig. S5B,C.

285

(B) Correlation of number of significantly associated *k*-mers vs. SNPs for all phenotypes. See also Fig. S5E.

286

(C) Correlation of p-values of top *k*-mers and SNPs for all phenotypes.

287

(D) Manhattan plot of SNP associations with days to tassel (environment 06FLI).

288

(E) LD between 23 significant SNPs and 18 *k*-mers (top) or *k*-mers to *k*-mers (bottom) for days to tassel. Order

289

of *k*-mers is the same in both heatmaps.

290

(F) LD between 45 *k*-mers associated with ear weight (environment 07A, left), and *k*-mer's presence/absence

291

patterns in different accessions ordered by their ear weight (right).

292

A major challenge in identifying causal variants in maize is the high fraction of short reads that

293

do not map uniquely to the genome. In the maize HapMap project, additional information had to be used

294

to find the genomic position of SNPs, including population LD and genetic map position (Bukowski et al.,

295

2018). The same difficulty of unique mappings also undermined the ability to identify the genomic source

296

of *k*-mers associated with specific traits. For example, we tried to locate the genomic position of the

297

*k*-mer corresponding to the SNP associated with days to tassel in chromosome 3 (Fig. 4D). The vast

298

majority of short reads from which the *k*-mer originated, 99%, could not be mapped uniquely to the

299

reference genome. However, when we assembled all these reads into a 924 bp contig, this fragment

300

could now be uniquely placed in the genome, to the same place as the identified SNPs. Thus, as we were

301 only interested in finding the genomic position after we already had an association in hand, we could use  
302 the richness of combining reads from many accessions to more precisely locate their origin without the  
303 use of additional genetic information, as had to be used for the SNPs.

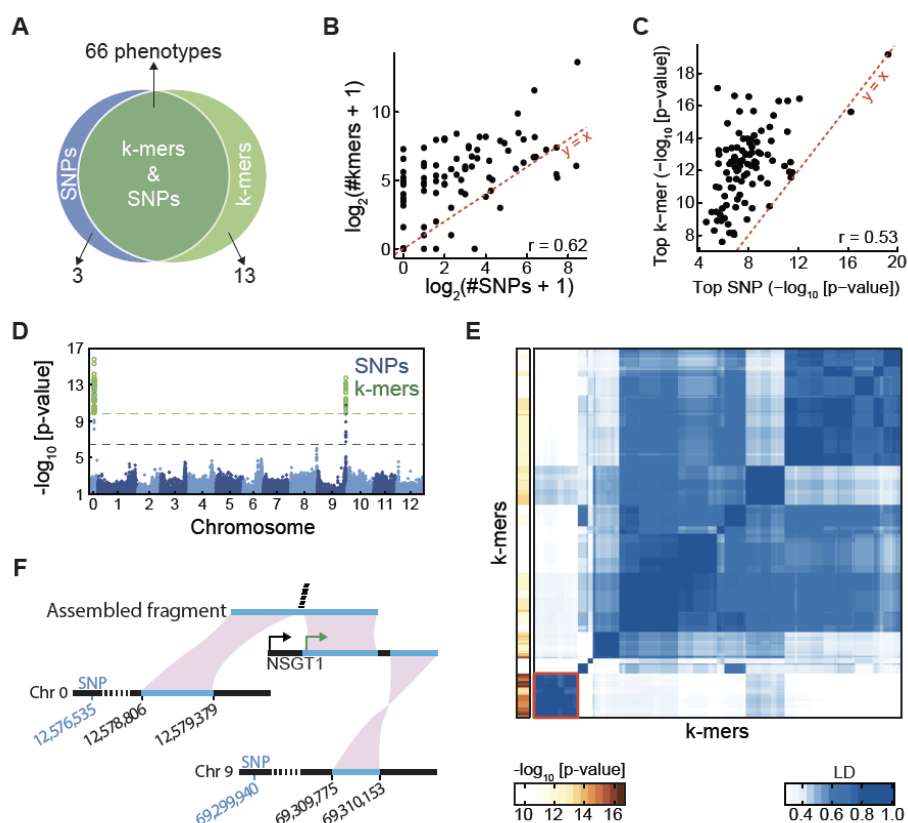
#### 304 **Comparison of SNP- and *k*-mer-based GWAS in tomato**

305 Tomato has a 900 Mb genome, which is intermediate between *A. thaliana* and maize, but it presents its  
306 own challenges, as modern tomatoes show a complex history of recent introgressions from wild  
307 relatives (Lin et al., 2014; Tomato Genome Consortium, 2012). Of 3.2 billion unique *k*-mers in all 246  
308 used accessions, 981 million were found in at least five accessions (Fig. S6A). We compared *k*-mer- and  
309 SNP-based GWAS on 96 metabolites measurements from two previous studies (Tieman et al., 2017;  
310 Zhu et al., 2018). For most metabolites, an association was identified by at least one method, with three  
311 metabolites having only SNP hits and 13 only *k*-mer hits (Fig. 5A). Similar to *A. thaliana* and maize, the  
312 number of identified variants as well as top p-values were correlated between methods (Fig. 5B,C). Top  
313 *k*-mers associations were also stronger than top SNPs (Fig. S5D), but even more so than in *A. thaliana*  
314 or maize, with an average difference of  $10^{4.4}$ , suggesting that in tomato the benefits of *k*-mer-based  
315 GWAS are also larger.

316 We next looked, as a case-study, at measurements of guaiacol, which results in a strong off-flavor  
317 and is therefore not desirable (Tieman et al., 2017). SNPs in two genomic loci were associated with it  
318 (Fig. 5D), one in chromosome 9 and the other in what is called “chromosome 0”, which corresponds to  
319 the concatenation of all sequence scaffolds that could not be assigned to one of the 12 nuclear  
320 chromosomes. From the 293 significant *k*-mers, 180 could be mapped uniquely to the genome, all close  
321 to significant SNPs. Among the remaining *k*-mers, of particular interest was a group of 35 *k*-mers in very  
322 high LD that had the lowest p-values, but could not be mapped to the reference genome (Fig. 5E).  
323 Assembly of the reads containing these *k*-mers resulted in a 1,172 bp fragment, of which the first 574 bp  
324 could be aligned near significant SNPs in chromosome 0 (Fig. 5F). The remainder of this fragment could  
325 not be placed in the reference genome, but there was a database match to the *NON-SMOKY*  
326 *GLYCOSYLTRANSFERASE 1 (NSGT1)* gene (Tikunov et al., 2013). The 35 significant *k*-mers covered  
327 the junction between these two mappable regions. Most of the *NSGT1* coding sequence is absent from  
328 the reference genome, but present in other accessions. *NSGT1* had been originally isolated as the causal  
329 gene for natural variation in guaiacol levels (Tikunov et al., 2013). Since *NSGT1* can be anchored to  
330 chromosome 9 near the identified SNPs (Fig. 5F), the significant SNPs identified in chromosomes 0 and 9  
331 apparently represent the same region, connected by the fragment we assembled from our set of 35

332 significant *k*-mers. Thus, we identified an association outside the reference genome, and linked the SNPs  
 333 in chromosome 0 to chromosome 9.

334



335

### Figure 5. Comparison of SNP- and *k*-mer-based GWAS in tomato

336 (A) Overlap between phenotypes with SNP and *k*-mer hits. See also Fig. S5B,C.

337 (B) Correlation of number of significantly associated *k*-mers vs. SNPs for all phenotypes. See also Fig. S5E.

338 (C) Correlation of p-values of top *k*-mers and SNPs for all phenotypes.

339 (D) Manhattan plot of SNPs and *k*-mers associations with guaiacol concentration.

340 (E) LD among 293 *k*-mers associated with guaiacol concentration (right), and the p-value of each *k*-mer (left). Red  
 341 square on bottom left indicates the 35 *k*-mers with strongest p-values and no mappings to the reference genome.

342 (F) The first part of a fragment assembled from the 35 unmapped *k*-mers (E) mapped to chromosome 0 and the  
 343 second part to the unanchored complete *NSGT1* gene. Only the 3' end of *NSGT1* maps to the reference genome,  
 344 to chromosome 9. The green and black arrows marks the start of the *NSGT1* ORF in the R104 “smoky” line and  
 345 “non-smoky” lines, respectively (Tikunov et al., 2013). Two SNPs are indicated, which are the significant SNPs  
 346 closest to the two regions of the reference genome.

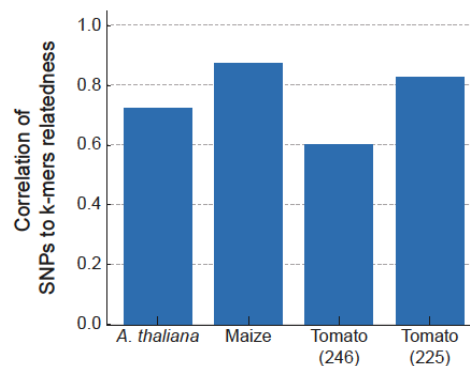
### 347 Calculation of relatedness between individuals based on *k*-mers

348 We have shown that we can assemble short fragments from *k*-mer-containing short reads and find hits  
 349 not only in the reference genome, but also in other published sequences. This opens the possibility to

350 apply our  $k$ -mer-based GWAS method to species without a high-quality reference genome. Draft  
351 genomes with contigs that include typically multiple genes can be relatively easily and cheaply generated  
352 using short read technology (Sohn and Nam, 2018). The major question with such an approach is then  
353 how one would correct for population structure in the GWAS step.

354 So far, we had relied on SNP kinship information. If one were to extend our method to species  
355 without high-quality reference genomes one would ideally be able to learn kinship directly from  $k$ -mers,  
356 thus obviating the need to map reads to a reference genome for SNP calling. With this goal in mind, we  
357 estimated relatedness using  $k$ -mers, applying the same method as with SNPs, with presence/absence as  
358 the two alleles. We calculated the relatedness matrices for *A. thaliana*, maize, and tomato and compared  
359 them to the SNP-based relatedness. In all three species there was agreement between the two methods,  
360 although initial results were clearly better for *A. thaliana* and maize than for tomato (Fig. 6). The inferior  
361 performance in tomato was due to 21 accessions (Fig. S7), which appeared to be more distantly related  
362 to the other accessions based on  $k$ -mer than what had been estimated with SNPs. This is likely due to  
363 these accessions containing diverged genomic regions that do poorly in SNP mapping, resulting in  
364 inaccurate relatedness estimates. Removing these 21 accessions increased the correlation between SNP-  
365 and  $k$ -mer-based relatedness estimates from 0.60 to 0.83. In conclusion,  $k$ -mers can be used to calculate  
366 relatedness between individuals, thus paving the way for GWAS in organisms without high-quality  
367 reference genomes.

368



369

#### Figure 6. Kinship matrix estimates with $k$ -mers

370 Relatedness between accessions was independently estimated based on SNPs and  $k$ -mers. The correlation between  
371 the two for tomato could be improved by removing 21 accessions that behaved differently between  $k$ -mers and  
372 SNPs (see Fig. S7).

373

374

## Discussion

375 The complexity of plant genomes makes identification of genotype-phenotype associations often  
376 challenging. To cope with this complexity, we followed a simple idea: most genetic variants leave a mark  
377 in the form of presence or absence of specific  $k$ -mers in whole genome sequencing data. Therefore,  
378 associating these  $k$ -mer marks with phenotypes will lead back to the genetic variants of interest. Our  
379 approach can identify associations found also by SNPs and short indels, but it excels when it comes to  
380 the detection of structural variants and variants not present in the reference genome. The expansion of  
381 variant types detected by our  $k$ -mer method complements SNP-based approaches, and greatly increases  
382 opportunities for finding and exploiting complex genetic variants driving phenotypic differences in plants,  
383 including improved genomic predictions.

384  $k$ -mers mark genetic polymorphisms in the population, but the types and genomic positions of  
385 these polymorphisms are initially not known. While one can also use  $k$ -mers for predictive models  
386 without knowing their genomic context, in many cases the genomic contexts of  $k$ -mers associated with  
387 certain phenotypes are of interest. The simplest solution is to align the  $k$ -mers or the short reads they  
388 originate from to a reference genome, an approach that was effective for some phenotypes we have  
389 studied, as it has been in bacteria (Pascoe et al., 2015). However, if  $k$ -mers can be mapped to the  
390 reference genome, the underlying variants are likely to be also tagged by SNPs, as we saw for *A. thaliana*  
391 flowering time. In case  $k$ -mers cannot be placed on the reference genome, one can first identify the  
392 originating short reads and assemble these into larger fragments. We found this to be a very effective  
393 path to uncovering the genomic context of  $k$ -mers. Particularly the combination of reads from multiple  
394 accessions can provide high local coverage around the  $k$ -mers of interest, increasing the chances that  
395 sizeable fragments can be assembled and located in the reference genome or in other sequence  
396 databases. For example, in the GWA of days to tassel in maize, reads containing the associated  $k$ -mers  
397 could not be assigned to a specific location in the genome, but the assembled fragment mapped to a  
398 unique genomic position. This approach, manually applied in this study, provides a framework to  
399 systematically elucidate  $k$ -mer's genomic context.

400 A main attraction of using  $k$ -mers as markers is that in principle they are able to tag many types  
401 of variants. A further improvement over our approach will be  $k$ -mers that tag heterozygous variants. In  
402 our current implementation, which relies on complete presence or absence of specific  $k$ -mers, only one  
403 of the homozygous states has to be clearly differentiated not only from the alternative homozygous  
404 state, but also from the heterozygous state. This did not affect comparisons between SNPs and  $k$ -mers  
405 in this study, as we only looked at inbred populations, where only homozygous, binary states are



406 expected. Another improvement will be to use  $k$ -mers to detect causal copy number variations. So far,  
407 we can only tag copy number variants, if the junctions produce unique  $k$ -mers, but it would be desirable  
408 to use also  $k$ -mers inside copy number variants. Therefore, a future improvement will be an  
409 implementation that uses normalized counts instead of presence/absence of  $k$ -mers, which will create a  
410 framework that can, at least in principle, detect almost any kind of genomic variation.

411 The comparison of the  $k$ -mer- and SNP- based GWAS provides an interesting view on tradeoffs  
412 in the characterization of genetic variability. The stronger top p-values obtained with  $k$ -mers in cases  
413 where a SNP is the actual underlying genetic-variant points to incomplete use of existing information in  
414 SNP calling. On the other hand by minimizing filtering of  $k$ -mers, we included in our analysis some  
415  $k$ -mers that represent only sequencing errors. Another potential source of noise comes from  $k$ -mers  
416 that are missed due to low coverage, which will be treated as absent. We reasoned that including these  
417 erroneous  $k$ -mers primarily has mostly computational costs, with some decrease in statistical power,  
418 since the chance of such  $k$ -mers generating an association signal is vanishingly small. Moreover, the high  
419 similarity of relatedness estimates using either SNPs (which are in essence largely filtered for sequencing  
420 errors) or all  $k$ -mers confirms that erroneous  $k$ -mers produce little signal. On the other hand, the  
421 higher effective number of  $k$ -mers compared to SNPs requires a more stringent threshold that takes the  
422 increased number of statistical tests into account and thereby decreases statistical power. This increase  
423 in test load is similar to the one that occurred when the genomics field moved from using microarray to  
424 next-generation sequencing in defining SNPs (100I Genomes Consortium, 2016; The 1000 Genomes  
425 Project Consortium, 2010; Weigel and Mott, 2009). Thus, the higher threshold is an inevitable result  
426 from increasing our search space to catch more genetic variants.

427  $k$ -mer associations inverts how GWAS is usually done. Instead of locating sequence variations in  
428 the genome and then associating them with a phenotype, we identify sequence-phenotype associations  
429 and only then find the genomic context of the sequence variations. Genome assemblies and genetic  
430 variant calling are procedures in which many logical decisions have to be made (Bradnam et al., 2013;  
431 Olson et al., 2015). These include high level decisions such as what information and software to use, as  
432 well as the many pragmatic thresholds chosen at each step of the way. Every community optimize these  
433 steps a bit differently, not least based on differences in the biology of the organisms they study, and  
434 surely these decisions affect downstream analyses (100I Genomes Consortium, 2016; Bukowski et al.,  
435 2018; Tieman et al., 2017). Here, we took a complementary path in which initially neither a genome  
436 reference nor variant calling is needed, trying to reduce arbitrary decisions to a bare minimum.  
437 Technological improvement in short- and long-read sequences as well as methods to integrate them into  
438 a population-level genetic variation data-structure will expand the covered genetic variants (Paten et al.,

439 2017; Schneeberger et al., 2009). While traditional GWAS methods will benefit from these technological  
440 improvements, so will *k*-mer based approaches, which will be able to use tags spanning larger genomic  
441 distances. Therefore, we posit that for GWAS purposes, *k*-mer based approaches are ideal because they  
442 minimize arbitrary choices when classifying alleles and because they capture more, almost optimal,  
443 information from raw sequencing data.

444

## 445 **Acknowledgment**

446 We thank the many colleagues who have shared *A. thaliana* phenotypic information with us. We thank in  
447 particular G. Zhu and S. Huang for help with tomato genotypic and phenotypic information and C.  
448 Romay, R. Bukowski, and E. Buckler for help with maize genotypes and phenotypes. We thank K. Swarts,  
449 F. Rabanal, I Soifer, and R. Schweiger for fruitful discussions and comments on the manuscript. This work  
450 was supported by ERC AdG IMMUNEMENSIS, DFG ERA-CAPS “1001 Genomes Plus” and the Max  
451 Planck Society.

452

## 453 **References**

- 453 1001 Genomes Consortium (2016). 1,135 Genomes Reveal the Global Pattern of Polymorphism in  
454 *Arabidopsis thaliana*. *Cell* *166*, 481–491.
- 455 Abney, M. (2015). Permutation testing in the presence of polygenic variation. *Genet. Epidemiol.* *39*,  
456 249–258.
- 457 Aranzana, M.J., Kim, S., Zhao, K., Bakker, E., Horton, M., Jakob, K., Lister, C., Molitor, J., Shindo, C., Tang,  
458 C., et al. (2005). Genome-wide association mapping in *Arabidopsis* identifies previously known flowering  
459 time and pathogen resistance genes. *PLoS Genet.* *1*, e60.
- 460 Arora, S., Steuernagel, B., Gaurav, K., Chandramohan, S., Long, Y., Matny, O., Johnson, R., Enk, J.,  
461 Periyannan, S., Singh, N., et al. (2019). Resistance gene cloning from a wild crop relative by sequence  
462 capture and association genetics. *Nat. Biotechnol.* *37*, 139–143.
- 463 Atwell, S., Huang, Y.S., Vilhjálmsson, B.J., Willems, G., Horton, M., Li, Y., Meng, D., Platt, A., Tarone, A.M.,  
464 Hu, T.T., et al. (2010). Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred  
465 lines. *Nature* *465*, 627–631.
- 466 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A.A., Dvorkin, M., Kulikov, A.S., Lesin, V.M., Nikolenko,  
467 S.I., Pham, S., Prjibelski, A.D., et al. (2012). SPAdes: a new genome assembly algorithm and its applications  
468 to single-cell sequencing. *J. Comput. Biol.* *19*, 455–477.
- 469 Bao, W., Kojima, K.K., and Kohany, O. (2015). Repbase Update, a database of repetitive elements in

- 470 eukaryotic genomes. *Mob. DNA* 6, 11.
- 471 Bennetzen, J.L. (2000). Transposable element contributions to plant gene and genome evolution. *Plant*  
472 *Mol. Biol.* 42, 251–269.
- 473 Bradnam, K.R., Fass, J.N., Alexandrov, A., Baranay, P., Bechner, M., Birol, I., Boisvert, S., Chapman, J.A.,  
474 Chapuis, G., Chikhi, R., et al. (2013). Assemblathon 2: evaluating de novo methods of genome assembly  
475 in three vertebrate species. *Gigascience* 2, 10.
- 476 Bryant, F.M., Hughes, D., Hassani-Pak, K., and Eastmond, P.J. (2019). Basic LEUCINE ZIPPER  
477 TRANSCRIPTION FACTOR67 Transactivates DELAY OF GERMINATION1 to Establish Primary Seed  
478 Dormancy in Arabidopsis. *Plant Cell* 31, 1276–1288.
- 479 Bukowski, R., Guo, X., Lu, Y., Zou, C., He, B., Rong, Z., Wang, B., Xu, D., Yang, B., Xie, C., et al. (2018).  
480 Construction of the third-generation *Zea mays* haplotype map. *Gigascience* 7, 1–12.
- 481 Chan, E.K.F., Rowe, H.C., Hansen, B.G., and Kliebenstein, D.J. (2010). The complex genetic architecture  
482 of the metabolome. *PLoS Genet.* 6, e1001198.
- 483 Cheng, C.-Y., Krishnakumar, V., Chan, A.P., Thibaud-Nissen, F., Schobel, S., and Town, C.D. (2017).  
484 Araport 11: a complete reannotation of the Arabidopsis thaliana reference genome. *Plant J.* 89, 789–804.
- 485 Danecsek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., Handsaker, R.E., Lunter, G.,  
486 Marth, G.T., Sherry, S.T., et al. (2011). The variant call format and VCFtools. *Bioinformatics* 27,  
487 2156–2158.
- 488 Davila, J.I., Arrieta-Montiel, M.P., Wamboldt, Y., Cao, J., Hagmann, J., Shedge, V., Xu, Y.-Z., Weigel, D., and  
489 Mackenzie, S.A. (2011). Double-strand break repair processes drive evolution of the mitochondrial  
490 genome in Arabidopsis. *BMC Biol.* 9, 64.
- 491 Devlin, B., and Risch, N. (1995). A comparison of linkage disequilibrium measures for fine-scale mapping.  
492 *Genomics* 29, 311–322.
- 493 Fordyce, R.F., Soltis, N.E., Caseys, C., Gwinner, R., Corwin, J.A., Atwell, S., Copeland, D., Feusier, J.,  
494 Subedy, A., Eshbaugh, R., et al. (2018). Digital Imaging Combined with Genome-Wide Association  
495 Mapping Links Loci to Plant-Pathogen Interaction Traits. *Plant Physiol.* 178, 1406–1422.
- 496 Gordon, S.P., Contreras-Moreira, B., Woods, D.P., Des Marais, D.L., Burgess, D., Shu, S., Stritt, C., Roulin,  
497 A.C., Schackwitz, W., Tyler, L., et al. (2017). Extensive gene content variation in the *Brachypodium*  
498 *distachyon* pan-genome correlates with population structure. *Nat. Commun.* 8, 2184.
- 499 Gore, M.A., Chia, J.-M., Elshire, R.J., Sun, Q., Ersoz, E.S., Hurwitz, B.L., Peiffer, J.A., McMullen, M.D.,  
500 Grills, G.S., Ross-Ibarra, J., et al. (2009). A first-generation haplotype map of maize. *Science* 326,  
501 1115–1117.
- 502 Hübner, S., Bercovich, N., Todesco, M., Mandel, J.R., Odenheimer, J., Ziegler, E., Lee, J.S., Baute, G.J.,  
503 Owens, G.L., Grassa, C.J., et al. (2018). Sunflower pan-genome analysis shows that hybridization altered  
504 gene content and disease resistance. *Nature Plants*.
- 505 Iqbal, Z., Caccamo, M., Turner, I., Flicek, P., and McVean, G. (2012). De novo assembly and genotyping of

- 506 variants using colored de Bruijn graphs. *Nat. Genet.* *44*, 226.
- 507 Joseph, B., Corwin, J.A., Li, B., Atwell, S., and Kliebenstein, D.J. (2013). Cytoplasmic genetic variation and  
508 extensive cytonuclear interactions influence natural variation in the metabolome. *Elife* *2*, e00776.
- 509 Kang, H.M., Zaitlen, N.A., Wade, C.M., Kirby, A., Heckerman, D., Daly, M.J., and Eskin, E. (2008). Efficient  
510 control of population structure in model organism association mapping. *Genetics* *178*, 1709–1723.
- 511 Kokot, M., Dlugosz, M., and Deorowicz, S. (2017). KMC 3: counting and manipulating k-mer statistics.  
512 *Bioinformatics* *33*, 2759–2761.
- 513 Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat. Methods* *9*,  
514 357–359.
- 515 Lees, J.A., Vehkala, M., Välimäki, N., Harris, S.R., Chewapreecha, C., Croucher, N.J., Marttinen, P., Davies,  
516 M.R., Steer, A.C., Tong, S.Y.C., et al. (2016). Sequence element enrichment analysis to determine the  
517 genetic basis of bacterial phenotypes. *Nat. Commun.* *7*, 12797.
- 518 Lees, J.A., Croucher, N.J., Goldblatt, D., Nosten, F., Parkhill, J., Turner, C., Turner, P., and Bentley, S.D.  
519 (2017). Genome-wide identification of lineage and locus specific variation associated with pneumococcal  
520 carriage duration. *eLife* *6*.
- 521 Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* *34*, 3094–3100.
- 522 Li, X., Svedin, E., Mo, H., Atwell, S., Dilkes, B.P., and Chapple, C. (2014). Exploiting natural variation of  
523 secondary metabolism identifies a gene controlling the glycosylation diversity of dihydroxybenzoic acids  
524 in *Arabidopsis thaliana*. *Genetics* *198*, 1267–1276.
- 525 Lin, T., Zhu, G., Zhang, J., Xu, X., Yu, Q., Zheng, Z., Zhang, Z., Lun, Y., Li, S., Wang, X., et al. (2014).  
526 Genomic analyses provide insights into the history of tomato breeding. *Nat. Genet.* *46*, 1220–1226.
- 527 Loh, P.-R., Tucker, G., Bulik-Sullivan, B.K., Vilhjálmsson, B.J., Finucane, H.K., Salem, R.M., Chasman, D.I.,  
528 Ridker, P.M., Neale, B.M., Berger, B., et al. (2015). Efficient Bayesian mixed-model analysis increases  
529 association power in large cohorts. *Nat. Genet.* *47*, 284–290.
- 530 Long, Q., Rabanal, F.A., Meng, D., Huber, C.D., Farlow, A., Platzer, A., Zhang, Q., Vilhjálmsson, B.J., Korte,  
531 A., Nizhynska, V., et al. (2013). Massive genomic variation and strong selection in *Arabidopsis thaliana*  
532 lines from Sweden. *Nat. Genet.* *45*, 884–890.
- 533 McClintock, B. (1950). The origin and behavior of mutable loci in maize. *Proc. Natl. Acad. Sci. U. S. A.*  
534 *36*, 344–355.
- 535 Minio, A., Massonnet, M., Figueroa-Balderas, R., Castro, A., and Cantu, D. (2019). Diploid Genome  
536 Assembly of the Wine Grape Carménère. *G3* *9*, 1331–1337.
- 537 Morrison, G.D., and Linder, C.R. (2014). Association mapping of germination traits in *Arabidopsis*  
538 *thaliana* under light and nutrient treatments: searching for G×E effects. *G3* *4*, 1465–1478.
- 539 Olson, N.D., Lund, S.P., Colman, R.E., Foster, J.T., Sahl, J.W., Schupp, J.M., Keim, P., Morrow, J.B., Salit,  
540 M.L., and Zook, J.M. (2015). Best practices for evaluating single nucleotide variant calling methods for

- 541 microbial genomics. *Front. Genet.* *6*, 235.
- 542 Pascoe, B., Méric, G., Murray, S., Yahara, K., Mageiros, L., Bowen, R., Jones, N.H., Jeeves, R.E.,  
543 Lappin-Scott, H.M., Asakura, H., et al. (2015). Enhanced biofilm formation and multi-host transmission  
544 evolve from divergent genetic backgrounds in *Campylobacter jejuni*. *Environ. Microbiol.* *17*, 4779–4789.
- 545 Paten, B., Novak, A.M., Eizenga, J.M., and Garrison, E. (2017). Genome graphs and the evolution of  
546 genome inference. *Genome Res.* *27*, 665–676.
- 547 Portwood, J.L., 2nd, Woodhouse, M.R., Cannon, E.K., Gardiner, J.M., Harper, L.C., Schaeffer, M.L., Walsh,  
548 J.R., Sen, T.Z., Cho, K.T., Schott, D.A., et al. (2019). MaizeGDB 2018: the maize multi-genome genetics  
549 and genomics database. *Nucleic Acids Res.* *47*, D1146–D1154.
- 550 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A.R., Bender, D., Maller, J., Sklar, P., de  
551 Bakker, P.I.W., Daly, M.J., et al. (2007). PLINK: a tool set for whole-genome association and  
552 population-based linkage analyses. *Am. J. Hum. Genet.* *81*, 559–575.
- 553 Rahman, A., Hallgrímsdóttir, I., Eisen, M., and Pachter, L. (2018). Association mapping from sequencing  
554 reads using k-mers. *Elife* *7*.
- 555 Ristova, D., Giovannetti, M., Metesch, K., and Busch, W. (2018). Natural Genetic Variation Shapes Root  
556 System Responses to Phytohormones in *Arabidopsis*. *Plant J.*
- 557 Robinson, J.T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E.S., Getz, G., and Mesirov, J.P.  
558 (2011). Integrative genomics viewer. *Nat. Biotechnol.* *29*, 24.
- 559 Saxena, R.K., Edwards, D., and Varshney, R.K. (2014). Structural variations in plant genomes. *Brief. Funct.*  
560 *Genomics* *13*, 296–307.
- 561 Schnable, P.S., Ware, D., Fulton, R.S., Stein, J.C., Wei, F., Pasternak, S., Liang, C., Zhang, J., Fulton, L.,  
562 Graves, T.A., et al. (2009). The B73 maize genome: complexity, diversity, and dynamics. *Science* *326*,  
563 1112–1115.
- 564 Schneeberger, K., Hagemann, J., Ossowski, S., Warthmann, N., Gesing, S., Kohlbacher, O., and Weigel, D.  
565 (2009). Simultaneous alignment of short reads against multiple genomes. *Genome Biol.* *10*, R98.
- 566 Seren, Ü., Grimm, D., Fitz, J., Weigel, D., Nordborg, M., Borgwardt, K., and Korte, A. (2017). AraPheno:  
567 a public database for *Arabidopsis thaliana* phenotypes. *Nucleic Acids Res.* *45*, D1054–D1059.
- 568 Sheppard, S.K., Didelot, X., Méric, G., Torralbo, A., Jolley, K.A., Kelly, D.J., Bentley, S.D., Maiden, M.C.J.,  
569 Parkhill, J., and Falush, D. (2013). Genome-wide association study identifies vitamin B5 biosynthesis as a  
570 host specificity factor in *Campylobacter*. *Proc. Natl. Acad. Sci. U. S. A.* *110*, 11923–11927.
- 571 Sohn, J.-I., and Nam, J.-W. (2018). The present and future of de novo whole-genome assembly. *Brief.*  
572 *Bioinform.* *19*, 23–40.
- 573 Springer, N.M., Anderson, S.N., Andorf, C.M., Ahern, K.R., Bai, F., Barad, O., Barbazuk, W.B., Bass, H.W.,  
574 Baruch, K., Ben-Zvi, G., et al. (2018). The maize W22 genome provides a foundation for functional  
575 genomics and transposon biology. *Nat. Genet.* *50*, 1282–1288.
- 576 Sun, S., Zhou, Y., Chen, J., Shi, J., Zhao, H., Zhao, H., Song, W., Zhang, M., Cui, Y., Dong, X., et al. (2018).

- 577 Extensive intraspecific gene order and gene structural variations between Mo17 and other maize  
578 genomes. *Nat. Genet.* *50*, 1289–1295.
- 579 Svishcheva, G.R., Axenovich, T.I., Belonogova, N.M., van Duijn, C.M., and Aulchenko, Y.S. (2012). Rapid  
580 variance components–based method for whole-genome association analysis. *Nat. Genet.* *44*, 1166.
- 581 The 1000 Genomes Project Consortium (2010). A map of human genome variation from  
582 population-scale sequencing. *Nature* *467*, 1061–1073.
- 583 Tieman, D., Zhu, G., Resende, M.F.R., Jr, Lin, T., Nguyen, C., Bies, D., Rambla, J.L., Beltran, K.S.O., Taylor,  
584 M., Zhang, B., et al. (2017). A chemical genetic roadmap to improved tomato flavor. *Science* *355*,  
585 391–394.
- 586 Tikunov, Y.M., Molthoff, J., de Vos, R.C.H., Beekwilder, J., van Houwelingen, A., van der Hooft, J.J.J.,  
587 Nijenhuis-de Vries, M., Labrie, C.W., Verkerke, W., van de Geest, H., et al. (2013). Non-smoky  
588 glycosyltransferase I prevents the release of smoky aroma from tomato fruit. *Plant Cell* *25*, 3067–3078.
- 589 Tomato Genome Consortium (2012). The tomato genome sequence provides insights into fleshy fruit  
590 evolution. *Nature* *485*, 635–641.
- 591 Vetter, M., Karasov, T.L., and Bergelson, J. (2016). Differentiation between MAMP Triggered Defenses in  
592 *Arabidopsis thaliana*. *PLoS Genet.* *12*, e1006068.
- 593 Wang, W., Mauleon, R., Hu, Z., Chebotarov, D., Tai, S., Wu, Z., Li, M., Zheng, T., Fuentes, R.R., Zhang, F.,  
594 et al. (2018). Genomic variation in 3,010 diverse accessions of Asian cultivated rice. *Nature* *557*, 43–49.
- 595 Weigel, D., and Mott, R. (2009). The 1001 genomes project for *Arabidopsis thaliana*. *Genome Biol.* *10*,  
596 107.
- 597 Zapata, L., Ding, J., Willing, E.-M., Hartwig, B., Bezdan, D., Jiao, W.-B., Patel, V., Velikkakam James, G.,  
598 Koornneef, M., Ossowski, S., et al. (2016). Chromosome-level assembly of *Arabidopsis thaliana* Ler  
599 reveals the extent of translocation and inversion polymorphisms. *Proc. Natl. Acad. Sci. U. S. A.* *113*,  
600 E4052–E4060.
- 601 Zhao, W., Canaran, P., Jurkuta, R., Fulton, T., Glaubitz, J., Buckler, E., Doebley, J., Gaut, B., Goodman, M.,  
602 Holland, J., et al. (2006). Panzea: a database and resource for molecular and functional diversity in the  
603 maize genome. *Nucleic Acids Res.* *34*, D752–D757.
- 604 Zhou, X., and Stephens, M. (2012). Genome-wide efficient mixed-model analysis for association studies.  
605 *Nat. Genet.* *44*, 821–824.
- 606 Zhu, G., Wang, S., Huang, Z., Zhang, S., Liao, Q., Zhang, C., Lin, T., Qin, M., Peng, M., Yang, C., et al.  
607 (2018). Rewiring of the Fruit Metabolome in Tomato Breeding. *Cell* *172*, 249–261.e12.
- 608 Zielesinski, A., Girgis, H.Z., Bernard, G., Leimeister, C.-A., Tang, K., Dencker, T., Lau, A.K., Röhling, S.,  
609 Choi, J.J., Waterman, M.S., et al. (2019). Benchmarking of alignment-free sequence comparison methods.  
610 *Genome Biol.* *20*, 144.



611

## Methods

612

### Curation of an *A. thaliana* phenotype compendium

613 Studies containing phenotypic data on *A. thaliana* accessions were located by searching NCBI PubMed using  
614 a set of general terms. For most studies, relevant data was obtained from the supplementary information or  
615 an online repository. Requests were sent to the corresponding authors of studies for which data could not  
616 be found in the public domain. Data already uploaded to the AraPheno dataset (Seren et al., 2017)  
617 downloaded from there. Phenotypic data in PDF format was extracted using Tabula software. Different sets  
618 of naming for accessions were converted to accession indices. In case an index for an accession could not  
619 be located, we omitted the corresponding data point. In case an accession could potentially be assigned to  
620 different indices, we first checked if it was part of the 1001 Genomes project; if so, we used the 1001  
621 Genomes index. In case the accession was not part of it, one of the possible indices was assigned at  
622 random. Phenotypes of metabolite measurements from two studies, (Fordyce et al., 2018) and (Chan et al.,  
623 2010), were filtered to a reduced set by the following procedure: take the first phenotype, sequentially  
624 retain phenotypes if correlation with all previously taken phenotypes is lower than 0.7. Data from the  
625 second study (Chan et al., 2010), were further filtered for phenotypes with a title. Assignment of categories  
626 for each phenotype was done manually (Table S1). All processed phenotypic data can be found in Dataset  
627 S1.

628

### Whole genome sequencing data and variant calls of *A. thaliana*

629 Whole genome short reads for 1,135 *A. thaliana* accessions were downloaded from NCBI SRA (accession  
630 SRP056687). Accessions with fewer than  $10^8$  unique *k*-mers, a proxy for low effective coverage, were  
631 removed, resulting in a set of 1,008 accessions. The 1001 Genomes project VCF file with SNPs and short  
632 indels was downloaded from <http://1001genomes.org/data/GMI-MPI/releases/v3.1> and was condensed into  
633 these 1,008 accessions, using vcftools v0.1.15 (Danecek et al., 2011). We required a minimum minor allele  
634 count (MAC) of 5 individuals, resulting in 5,649,128 genetic variants. The VCF file was then converted to a  
635 PLINK binary file using PLINK v1.9 (Purcell et al., 2007). In case more than two alleles were possible in a  
636 specific location, PLINK keeps the reference allele and the most common alternative allele. The TAIR10  
637 reference genome was used for short read and *k*-mer alignments. Coordinates for genes in figures were  
638 taken from Araport1.1 (Cheng et al., 2017).

639

### Whole genome sequencing data and variant calls of maize

640 Whole genome short reads of maize accessions corresponded to the “282 set” part of the maize  
641 HapMap3.2.1 project (Bukowski et al., 2018). Sequencing libraries “x2” and “x4” were downloaded from  
642 NCBI SRA (accession PRJNA389800) and combined. Coverage per accession was calculated as number of  
643 reads multiplied by read length and divided by the genome size, only data for 150 accessions with coverage  
644 >6x was used. Phenotypic data for 252 traits measured for these accessions were downloaded from Panzea  
645 (<https://www.panzea.org>) (Zhao et al., 2006).

646 Two of these phenotypes were constant over more than 90% of the 150 accessions, these two were  
647 removed from further analysis (“NumberofTilleringPlants\_env\_07A”,



648 “TilleringIndex-BorderPlant\_env\_07A”). The HapMap3.2.1 VCF files (c\*\_282\_corrected\_onHmp321.vcf.gz)  
649 of SNPs and indels were downloaded from Cyverse. Variant files were filtered using vcftools v0.1.15 to the  
650 relevant 150 accessions. Variants were further filtered for MAC of  $\geq 5$ , resulting in a final set of 35,522,659  
651 variants. The B73 reference genome, version AGPv3 (Portwood et al., 2019), that was used to create the  
652 VCF file was downloaded from MaizeGDB and used for short read and  $k$ -mer alignments (Portwood et al.,  
653 2019).

654

## 655 **Whole genome sequencing data and variant calls of tomato**

656 Whole genome short reads were downloaded for 246 accessions with coverage  $>6x$ , from NCBI SRA and  
657 EBI ENA (accession numbers SRP045767, PRJEB5235 and PRJNA353161). A table with coverage per  
658 accession was shared by the authors (Tieman et al., 2017). Metabolite measurements were taken from  
659 (Tieman et al., 2017) (only adjusted values) and a subset of metabolites from (Zhu et al., 2018). These were  
660 filtered to a reduced set by the following procedure: take the first phenotype, sequentially retain  
661 phenotypes if correlation with all previously taken phenotypes is lower than 0.7. Metabolites were ordered  
662 as reported originally (Zhu et al., 2018). Only one repeat, the one with more data points and requiring at  
663 least 40 data points was retained. The VCF file with SNPs and short indels (Tieman et al., 2017) was  
664 obtained from the authors and filtered for the relevant 246 accessions. Variants were further filtered for  
665 MAC of  $\geq 5$ , resulting in a final set of 2,076,690 variants. Reference genome SL2.5 (Tomato Genome  
666 Consortium, 2012) ([https://www.ncbi.nlm.nih.gov/assembly/GCF\\_000188115.3/](https://www.ncbi.nlm.nih.gov/assembly/GCF_000188115.3/)) used to create the VCF file  
667 was used for short read and  $k$ -mer alignments.

## 668 **$k$ -mer counting and initial processing**

669 For each accession from each of the three species all sequencing data from different runs were combined.  
670 The number of times each  $k$ -mer ( $k=25\text{bp}/31\text{bp}$ ) appeared in the raw sequencing reads were counted using  
671 KMC v3 (Kokot et al., 2017).  $k$ -mers were counted twice, first counting canonical  $k$ -mers representation,  
672 which is the lower lexicographically for a  $k$ -mer and its reverse-complement. This list contains only  $k$ -mer  
673 appearing at least twice (maize and tomato) or thrice (*A. thaliana*) in the sequence reads. The second count  
674 includes all  $k$ -mers and without canonization. The KMC binary outputs of  $k$ -mers counts in the two lists  
675 were read using KMC C++ API, to keep all calculations in binary representation. For each  $k$ -mer in the first  
676 list, the information of which form (canonized, not-canonized, or both) it appeared in was extracted from  
677 the second list. This form information was coded in two bits, where the first/second bit indicates if the  $k$ -mer  
678 was observed in its canonized/non-canonized form, respectively. These two bits were inserted in the 2  
679 most-significant-bits of the  $k$ -mer bit representation, as  $k$ -mers are of maximal length of 31bp, all  
680 information could be coded in a 64-bit word. The 64-bit  $k$ -mers representation were sorted according to  
681 the  $k$ -mer lexicographic order and saved to a file in binary representation.

682 For each species, the latter  $k$ -mers lists from all accessions were combined into one list according to the  
683 following criteria: only  $k$ -mers appearing in at least 5 accessions, and for a  $k$ -mer appearing in  $N$  accession  
684 it had to be observed in both canonized and non-canonized form in at least  $0.2*N$  of the accessions. There  
685 were  $2.26*10^9$ ,  $2.21*10^9$ ,  $3.23*10^9$ , and  $7.28*10^9$  unique  $k$ -mers in all accessions in the first type of counting,  
686 i.e. before filtering, and  $439*10^6$ ,  $393*10^6$ ,  $981*10^6$ , and  $2.33*10^9$  passed the second criteria for *A. thaliana*

687 (31-mers), *A. thaliana* (25-mers), tomato (31-mers), and maize (31-mers), respectively. The final filtered  
688 *k*-mers were outputted in binary format to a file, the histogram of number of *k*-mers appearances was  
689 calculated and saved during this process as well (e.g. Fig. S2A).

## 690 **Combining *k*-mers from different accessions to a *k*-mers presence/absence table**

691 Tables containing the presence/absence per *k*-mer per accession in binary format were created, for each  
692 specie and *k*-mer size. The tables were organized as follows: *k*-mers information was written in serialized  
693 blocks of N+1 64-bit words. In each block, the first word codes the *k*-mer ( $k < 32$ bp), the next N 64-bit  
694 blocks codes for the presence/absence of the *k*-mer in the different accessions: 1 in position *i* denoting the  
695 *k*-mer was found in accession *i* and 0 otherwise. N is the number of accessions divided by 64, rounded up.  
696 The last remaining padding bits not used were set to 0. Calculation of tables was done as follows: *k*-mers  
697 lists for all accession were opened together, in each step all the *k*-mers up to a threshold were read. *k*-mers  
698 were then combined in a sub-table to create the presence/absence patterns and then outputted in the  
699 described format with lexicographically ordered *k*-mers. This process was designed to minimize the  
700 memory load, and could be achieved due to the sorted *k*-mers in all separate lists.

## 701 **Counting and filtering unique presence absence patterns of *k*-mers**

702 To check if a specific presence/absence pattern was already observed, the following method was used. This  
703 was done in order to count or filter the patterns. Each pattern, represented by a vector of N 64-bit words  
704 was inputted in a hash function which outputs a single 64-bit word. The hashed value was then stored in a  
705 set structure built on a hash-table. The size of the set was an indication of the number of unique patterns.  
706 Moreover, it was used continuously to filter patterns, by checking if a pattern (its hashed value) was already  
707 observed. The probability that two different patterns had the same hash value is very low: if we have *n*  
708 patterns, the space is of size  $S = 2^{64}$ , the probability that at least one collision occurs randomly is:

$$709 \quad p = 1 - (2^{64}/2^{64})((2^{64} - 1)/2^{64}) \dots ((2^{64} - n + 1)/2^{64}) \approx 1 - e^0 e^{-1/2^{64}} e^{-2/2^{64}} \dots e^{-(n-1)/2^{64}} = 1 - e^{-(n-1)n/2^{65}}$$

710 If  $n = 2^{30} > 1,000,000,000$  then  $p \approx 1 - e^{-2^{60}/2^{65}} < 0.031$ , so there is ~97% chance that not even one  
711 collision occurred for 1 billion distinct *k*-mers.

## 712 **Calculate and comparison of kinship matrices**

713 Kinship matrix of relatedness between accessions was calculated as in EMMA (Kang et al., 2008), with  
714 default parameters. The algorithm was re-coded in C++ to read directly PLINK binary files for improved  
715 efficiency. For *k*-mers based relatedness the same algorithm was used, coding presence/absence as two  
716 alleles. For comparison of *k*-mers- to SNPs-based relatedness we correlated (pearson) the values for all  $\binom{n}{2}$   
717 pairs, for *n* accessions. For tomato, 3492 pairs had a relatedness more than 0.15 lower for *k*-mer than for  
718 SNPs. 3,298 (94.4%) of these pairs were between a set of 21 accessions and all other 225 accessions. We  
719 calculated the correlation twice: for all pairs, and only between pairs of these 225 accessions.

## 720 **GWA on SNPs and short indels or on full $k$ -mers table**

721 Genome-wide association on the full set of SNPs and short indels was conducted using linear mixed models  
722 with the kinship matrix, using GEMMA version 0.96 (Zhou and Stephens, 2012). Minor allele frequency  
723 (MAF) was set to 5% and MAC was set to 5, with a maximum of 50% missing values (-miss 0.5). Kinship  
724 matrix was used to account for population structure. To run GWA on the full set of  $k$ -mers (e.g. in Fig. 1B),  
725  $k$ -mers were first filtered for  $k$ -mers having only unique patterns on the relevant set of accessions, MAF of  
726 at least 5%, and MAC of at least 5. Presence/absence patterns were then condensed to only the relevant  
727 accessions and output as a PLINK binary file directly. GEMMA was then run using the same parameters as  
728 for the SNPs GWA described above.

## 729 **Phenotype covariance matrix estimation and phenotypes permutation**

730 EMMA (emma.REMLE function) was used to calculate the variance components which were used to  
731 calculate the phenotypic covariance matrix (Kang et al., 2008). We then calculated 100 permutations of the  
732 phenotype using the mvnpermute R package (Abney, 2015). The  $n\%$  (e.g.  $n=5$  gives 5%) family-wise error  
733 rate threshold was defined by taking the  $n$ -th top p-value from the 100 top p-value of running GWA on  
734 each permutation. In all cases, unless indicated otherwise, where a threshold is referred to, it is the 5%  
735 threshold.

## 736 **Scoring p-values from GWA for similarity to uniform distribution and filtering phenotypes**

737 Each SNP-based GWA run was scored for a general bias in p-value distribution, similar to Atwell et al.  
738 (Atwell et al., 2010). All SNPs p-values were collected, the 99% higher p-values were tested against the  
739 uniform distribution using a kolmogorov-smirnov test, and the test statistic was used to filter phenotypes  
740 for which distribution deviated significantly from the uniform distribution. A threshold of 0.05 was used,  
741 filtering 89, 0, and 295 phenotypes for *A. thaliana*, maize and tomato, respectively.

## 742 **$K$ -mers genome-wide associations**

743 Association of  $k$ -mers was done in two steps, with the aim of getting the most significant  $k$ -mers p-values.  
744 The first step was based on the approach used in Bolt-Imm-inf and GRAMMAR-Gamma (Loh et al., 2015;  
745 Svishcheva et al., 2012). For phenotypes  $y$ , genotypes  $g$ , and a covariance matrix  $\Omega$ , the  $k$ -mer score is:

$$746 \quad T_{score}^2 = \frac{1}{\gamma} \frac{(\tilde{g}^T \Omega^{-1} \tilde{y})^2}{\tilde{g}^T \tilde{g}}$$

747 Where  $\tilde{g} = g - E(g)$  and  $\tilde{y} = y - E(y)$ . The first step was used only to filter a fixed number of top  
748  $k$ -mers, thus we could use any score monotonous with  $T_{score}^2$ , and specifically  $\frac{(\tilde{g}^T \Omega^{-1} \tilde{y})^2}{\tilde{g}^T \tilde{g}}$  which is  
749 independent of  $\gamma$  (see supplementary note on calculation optimization). To keep used memory low, only  
750 best  $k$ -mers were stored in a priority queue data structure of constant size. The  $k$ -mers-table was uploaded  
751 to the memory in small chunks and associations were done with the phenotype and its permuted  
752 phenotypes for all  $k$ -mers in each chunk. The association step was implemented with the use of threads.  
753 After all  $k$ -mers were scored for associations with the phenotype and all its permutations, the  $k$ -mers-table

754 was loaded again in chunks. The top  $k$ -mers with their genotype patterns were outputted in binary PLINK  
 755 format, for the phenotype and each permutation separately. In the second step, the best  $k$ -mers were run  
 756 using GEMMA to calculate the likelihood ratio test p-values (Zhou and Stephens, 2012).

757 The number of  $k$ -mers filter in the first step was set to 10,000 for *A. thaliana* and 100,000 for maize and  
 758 tomato. Both steps associate  $k$ -mers while accounting for population structure, while the first step uses an  
 759 approximation, the second use an exact model. Therefore, real top  $k$ -mers might be lost as they would not  
 760 pass the first filtering step. To control for this, we first defined the 5% family-wise error-rate threshold  
 761 based on the phenotype permutations, and then identified all the  $k$ -mers which passed the threshold. Next,  
 762 we used the following criteria to minimize the chance of losing  $k$ -mers: we checked if all identified  $k$ -mers  
 763 were in the top  $N/2$   $k$ -mers from the ordering of the first step ( $N=10,000$  or  $100,000$  dependent on  
 764 species). For example, in maize all  $k$ -mers passing the threshold in the second step should be in the top  
 765 50,000  $k$ -mers from the first step. The probability that this will happen randomly is  $2^{-m}$ , where  $m$  is  
 766 number of identified  $k$ -mers, in most phenotypes this is very unlikely. In 8.5% of phenotypes from *A. thaliana*  
 767 the criteria was not fulfilled, for these phenotypes we re-run the two-steps with 100x more  $k$ -mers filtered  
 768 in the first step, that is 1,000,000  $k$ -mers. For 6 phenotypes the criteria still did not hold, these phenotypes  
 769 were not used in further analysis. In tomato, 33% of phenotypes did not fulfill these criteria, in these cases  
 770 we re-run the first step with 100x more  $k$ -mers filtered (10,000,000), 17 phenotypes still did not pass the  
 771 threshold and were omitted from further analysis. The permutations were not re-run, and the threshold  
 772 defined using 100,000  $k$ -mers was used, as the top  $k$ -mer used to define the threshold tended to be high in  
 773 the list. For maize all phenotypes passed the criteria and no re-running was needed.

#### 774 **Optimizing of initial $k$ -mers scoring**

775 For:  $N$  – number of individuals,  $\Omega$  – covariance matrix,  $y$  – phenotype,  $g$  – genotype (for  $k$ -mers taking  
 776 the values 0 for absence and 1 for presence), and  $\gamma$  - GRAMMAR-Gamma factor which depends on the  
 777 phenotype and relatedness between individuals, but not on specific  $g$  (Svishcheva et al., 2012).

778  $\tilde{y} = y - E(y)$  and  $\tilde{g} = g - E(g)$

779  $r = \Omega^{-1}\tilde{y}$  the transformed phenotype

780 The GRAMMAR-Gamma score of association  $T_{score}^2$  is distributed according to  $\chi^2$  with 1 d.f. and satisfies:

781 
$$T_{score}^2 = \frac{1}{\gamma} \left( \frac{\tilde{g}^T \Omega^{-1} \tilde{y}}{\tilde{g}^T \tilde{g}} \right)^2 = \frac{1}{\gamma} \frac{(\tilde{g}^T r)^2}{\tilde{g}^T \tilde{g}} = \frac{1}{\gamma} \frac{(\sum (g_i - \frac{\sum g_i}{N}) r_i)^2}{\sum (g_i - \frac{\sum g_i}{N})^2} = \frac{1}{\gamma} \frac{(\sum g_i r_i - \frac{\sum g_i}{N} \sum r_i)^2}{\sum (g_i^2 - 2g_i \frac{\sum g_i}{N} + (\frac{\sum g_i}{N})^2)}$$

782 
$$\frac{1}{\gamma} \frac{(N \sum g_i r_i - (\sum g_i)(\sum r_i))^2}{N^2 \sum g_i^2 - 2N \sum g_i \sum r_i + N (\sum g_i)^2} = \frac{1}{\gamma} \frac{(N \sum g_i r_i - (\sum g_i)(\sum r_i))^2}{N^2 \sum g_i^2 - N (\sum g_i)^2}$$

783 A  $k$ -mer can only be present or absent but not missing or heterozygous, thus  $g_i = g_i^2$  and we get:

784

$$T_{score}^2 = \frac{1}{\gamma N} \frac{\left( N \sum g_i r_i - \left( \sum g_i \right) \left( \sum r_i \right) \right)^2}{N \sum g_i^2 - \left( \sum g_i \right)^2}$$

785 As we used the GRAMMAR-Gamma score only to filter the top  $k$ -mers, we did not need to calculate the  
786 p-value of  $T_{score}^2$  and could calculate a score that is monotonous with  $T_{score}^2$ , that is:

787

$$K_{score} = \frac{\left( N \sum g_i r_i - \left( \sum g_i \right) \left( \sum r_i \right) \right)^2}{N \sum g_i^2 - \left( \sum g_i \right)^2}$$

788 The summation  $\sum r_i$  can be calculated once per phenotype. Moreover, as we use permutation of  
789 phenotypes we can further optimize the scoring by calculating  $\sum g_i$  only once per  $k$ -mer.

790 For calculating the score of a specific  $k$ -mer, once  $\sum r_i$ ,  $\sum g_i$ , and  $\sum g_i r_i$  were calculated, we were left with  
791 8 basic mathematical operations to obtain  $K_{score}$ . Therefore, most of the computational load will be spent  
792 in the calculation of  $\sum g_i r_i$ , which requires  $2N$  basic operations.

793 To computationally optimize the calculation of  $\sum g_i r_i$ , we used the Streaming SIMD Extensions 4 (SSE4)  
794 CPU instruction set. This implementation can be further optimized on a CPU that has AVX2, likely getting  
795 another 2-fold increase in efficiency with only small modifications to the code, however, we have not tested  
796 this option.

797

798 To optimize the GRAMMAR-Gamma filtering of SNPs we cannot benefit from the same optimizations as for  
799  $k$ -mers. This is due to missing and heterozygous values a SNP can take. Therefore, in this case  $g_i \neq g_i^2$ . For  
800 SNPs our score will take the same form as  $T_{score}^2$ :

801

$$S_{score} = \frac{1}{\gamma N} \frac{\left( N \sum g_i r_i v_i - \left( \sum g_i \right) \left( \sum v_i r_i \right) \right)^2}{N \sum g_i^2 v_i^2 - \left( \sum g_i v_i \right)^2}$$

802 In this case  $N$  is different for different SNPs, and so as  $\sum r_i$ . This later summation can be written as  $\sum v_i r_i$ ,  
803 by defining  $v_i = 0$  for  $g_i = \text{missing}$  and  $v_i = 1$  for  $g_i \neq \text{missing}$ .

804

805 Thus,  $\sum v_i r_i$  is specific for each SNP's score and as  $g_i$  can also get the value 0.5, we separated  $\sum g_i r_i$  to  
806 two separate dot-products in our implementation, as genotypes are coded by bit vectors.

### 807 **SNPs-based GWAS on phenotype permutations**

808 To calculate thresholds for SNPs-based GWAS we used the two step approach used for  $k$ -mers. The  
809 permuted phenotypes were run in two steps as we were only interested in the top p-value to define  
810 thresholds. We filtered 10,000 variants in the first step which were then run using GEMMA to get exact  
811 scores (Zhou and Stephens, 2012). The non-permuted phenotype were run using GEMMA on all the  
812 variants.

### 813 **Calculation of linkage-disequilibrium (LD)**

814 For two variants,  $x$  and  $y$ , each can be a  $k$ -mer or a SNP, LD measure was calculated using the  $r^2$  measure  
815 (Devlin and Risch, 1995). For a  $k$ -mer, variants were coded as 0/1, if absent or present, respectively. For  
816 SNPs one variant was coded as 0 and the other as 1. If one of the variants had a missing or heterozygous  
817 value in a position, this position was not used in the analysis. The LD value was calculated using the formula:

$$818 \quad r^2 = \frac{(p(x=1 \& y=1) - p(x=1)*p(y=1))^2}{p(x=1)*p(y=1)*p(x=0)*p(y=0)}$$

### 819 **LD cumulative graph (Fig 2E,H)**

820 For a set of phenotypes and for every  $l = 0, 0.05, \dots, 1$  we calculated the percentage of phenotypes for  
821 which exists a  $k$ -mer or a SNP in the pre-defined group which is in  $LD \geq l$  with top SNP or top  $k$ -mer,  
822 respectively. The pre-defined groups are: (1) all the  $k$ -mers which passed the SNPs defined threshold in  
823 Figure 2E or (2) all the SNPs or  $k$ -mers which passed their own defined thresholds in Figure 2H. The  
824 percentage is then plotted as a function of  $l$ .

### 825 **Retrieving source reads of a specific $k$ -mer and assembling them**

826 For a  $k$ -mer identified as being associated with a phenotype we first looked in the  $k$ -mers-table and  
827 identified all accessions taking part in the association analysis and having this  $k$ -mer present. For each of  
828 these accessions we went over all sequencing reads and filtered out all paired-end reads which contained  
829 the  $k$ -mer or its reverse-complement. To assemble paired-reads, SPAdes v3.11.1 was used with  
830 "--careful" parameter (Bankevich et al., 2012).

### 831 **Alignment of reads or $k$ -mers to the genome**

832 Paired-end reads were aligned to the genome using bowtie2 v2.2.3, with the "--very-sensitive-local"  
833 parameter.  $k$ -mers were aligned to the genome using bowtie v1.2.2 with "--best --all --strata" parameters  
834 (Langmead and Salzberg, 2012).

### 835 **Analysis of flowering time in 10C (Figure 1, Figure S2)**

836 To find the location in the genome of the 105 identified *k*-mers, *k*-mers were first mapped to the *A. thaliana*  
837 genome. 84 of the *k*-mers had a unique mapping, one *k*-mer was mapped to multiple locations and 20 could  
838 not be mapped. For the 21 *k*-mers with no unique mapping we located the sequencing reads they originated  
839 from, and mapped the reads to the *A. thaliana* genome. For each of the *k*-mers we looked only on the reads  
840 with the top mapping scores. For the one *k*-mer which had multiple possible alignment also the originating  
841 reads did not have a consensus mapping location in the genome. For every *k*-mer from the 20 non-mapped  
842 *k*-mers, all top reads per *k*-mer, in some cases except one, mapped to a specific region spanning a few  
843 hundred base pairs. The middle of this region was defined as the *k*-mer position for the Manhattan plot in  
844 Figure 1D.

845 To find the location of the 93 associated *k*-mers of length 25bp, presented in supplementary Figure S2D, we  
846 followed the same procedure. 87 of the *k*-mers had a unique mapping, one was mapped multiple times and  
847 5 could not be mapped. For the 5 *k*-mers with no mapping and the *k*-mer with non unique mapping, we  
848 located the originated short reads and aligned them to the genome. For each of the 5 *k*-mers with no  
849 mapping, all reads with top mapping score mapped to a specific region of a few hundred base pairs, we took  
850 the middle of the region as the *k*-mer location in the Manhattan plot. For the *k*-mer with multiple mappings,  
851 15 out of the 17 reads mapped to the same region and we used this location. All *k*-mers mapped to the 4  
852 location in the genome for which SNPs were identified except one - AAGCTACTTGGTTGATAATACTAAT.  
853 The reads from which this *k*-mer originated mapped to the same region in chromosome 5 position  
854 3191745-3192193 and we used the middle of this region as the *k*-mer location.

### 855 **Analysis of xylosides percentage (Figure 3A,B)**

856 All *k*-mers passing the threshold, were mapped uniquely to chromosome 5 in the region 871,976 –  
857 886,983. Of the 123 identified *k*-mers, 27 had the same minimal *p*-value ( $-\log_{10}(p\text{-value}) = 44.7$ ).  
858 These *k*-mers mapped to chromosome 5 in positions 871,976 to 872,002, all covering the region  
859 872,002-872,007. For the 60 accessions used in this analysis, all reads from the 1001G were mapped to  
860 the reference genome. The mapping in region 872,002-872,007 of chromosome 5 were examined  
861 manually by IGV in all accessions (Robinson et al., 2011), and the 2 SNPs 872,003 and 872,007 were  
862 called manually without knowledge of the phenotype value.

### 863 **Analysis of growth inhibition in presence of flg22 (Figure 3C,D)**

864 The phenotype in the original study was labeled “flgPsHRp” (Vetter et al., 2016). For each of the 7 *k*-mers  
865 which could not be mapped uniquely to the genome, the originated reads from all accessions were  
866 retrieved and assembled. All the seven cases resulted in the same assembled fragment (SEQ1, table S2).  
867 Using NCBI BLAST we mapped this fragment to chromosome 1: position 40-265 were mapped to  
868 8169229-8169455 and position 262-604 were mapped to 8170348-8170687. For every accession from the  
869 106 that were used in the GWAS analysis we tried to locally assemble this region, to see if the junction  
870 between chromosome 1 8169455 to 8170348 could be identified. We used all the 31bp *k*-mers from the  
871 above assembled fragment as bait, and located all the reads for each accession separately. For 11 out of the  
872 13 accessions that had all 10 identified *k*-mers we got a fragment from the assembly process. In all 11 cases



873 the exact same junction was identified. For 1 of the 4 accessions that had only part of the 10 identified  
874 *k*-mer we got a fragment from the assembler, which had the same junction. For 43 of the 89 accessions that  
875 had none of the identified *k*-mers the assembly process resulted in a fragment, in none of these cases the  
876 above junction could be identified.

#### 877 **Analysis of germination in darkness and low nutrients (Figure 3E, F)**

878 The phenotype in the original study was labeled “k\_light\_0\_nutrient\_0” (Morrison and Linder, 2014).  
879 The 11 identified *k*-mers had two possible presence/absence patterns, separating them into two groups  
880 of 4 and 7 *k*-mers. The short-read sequences containing the 4 or 7 *k*-mers were collected separately and  
881 assembled, resulting in the same 458bp fragment (SEQ2, table S2). This fragment was used as a query in  
882 NCBI BLAST search, resulting in alignment to Ler-0 chromosome 3 (LR215054.1) positions 15969670 to  
883 15970128. The LR215054.1 sequence was downloaded and the region between (15969670-3000) to  
884 (15970128+3000) was retrieved and used as query to a NCBI BLAST search. The BLAST search resulted  
885 in a mapping to Col-0 reference genome chromosome 3 (CP002686.1). Region 1-604 mapped to  
886 16075369-16075968, region 930-1445 mapped to 16076025-16076532, region 3446-3946 mapped to  
887 16079744-16080244, and region 3958-6459 mapped to 16080301-16082781.

#### 888 **Analysis of root branching zone (Figure 3G)**

889 The phenotype in the original study was labeled “Mean(R)\_C”, that is Branching zone in no treatment  
890 (Ristova et al., 2018). No SNPs and 1 *k*-mer (AGCTACTTTGCCACCCACTGCTACTAACTCG) passed  
891 their corresponding 5% thresholds. The *k*-mer mapped the chloroplast genome in position 40297, with 1  
892 mismatch. No SNPs and another *k*-mer (CCGGCGATTACTAGAGATTCCGGCTTCATGC) passed the  
893 10% family-wise error-rate threshold. This *k*-mer mapped non-uniquely to two place in the chloroplast  
894 genome: 102285 and 136332.

#### 895 **Analysis of Lesion by *Botrytis cinerea* UKRazz (Figure S3A)**

896 The Lesion by *Botrytis cinerea* UKRazz phenotype was labeled as “Lesion\_redgrn\_m\_theta\_UKRazz”. In  
897 the GWAS analysis 19 *k*-mers and no SNPs were identified. All *k*-mers had the same presence/absence  
898 pattern. The short-read sequences from which the *k*-mers originated were mapped to chromosome 3  
899 around position 72,000bp, and contained a 1-bp deletion of a T nucleotide in position 72,017. Whole  
900 genome sequencing reads were mapped to the genome for the 61 accessions with phenotypes used in these  
901 analyses. We manually observed the alignment around position 72,017 of chromosome 3, without the prior  
902 knowledge if the accession had the identified *k*-mers. For 20 accessions, we observed the 1-bp deletion in  
903 position 72,017, all 19 accessions containing the *k*-mers were part of these 20.

#### 904 **Analysis of days to tassel and ear weight in maize (Figure 4)**

905 Ear weight phenotype was labeled “EarWeight\_env\_07A” in original dataset (Zhao et al., 2006). Days to  
906 tassel was measured in growing degree days (GDD) and was labeled as “GDDDaystoTassel\_env\_06FL1” in  
907 original dataset. In comparison of LD between *k*-mers and SNPs in days to tassel (Fig. 4E, upper panel), two  
908 SNPs were filtered out as having more than 10% heterozygosity and one as having, exactly, 50% missing  
909 values. In days to tassel the *k*-mer which was similar to identified SNPs was

910 AGAAGATATCTTATGAACTCCTCACCAGTAA. The 171 paired-end reads from which this *k*-mer  
911 originated mapped to the genome as follows - 2 (1.17%) aligned concordantly 0 times, 2 (1.17%) aligned  
912 concordantly exactly 1 time, and 167 (97.66%) aligned concordantly >1 times. The assembly of these reads  
913 produce two fragments, the first of length 273bp with coverage of 1.23 and the second of length 924bp and  
914 with coverage of 27.41 (SEQ3, table S2). We aligned this fragment to the genome using Minimap2, with the  
915 default parameters (Li, 2018). Minimap2 reported only 1 hit to chromosome 3 (NC\_024461.1) in positions  
916 159141222-159142137.

### 917 **Analysis of guaiacol concentration in tomato (Figure 5)**

918 Guaiacol concentration was labeled “log3\_guaiacol” in the original study. From the 293 *k*-mers passing the  
919 threshold, 184 could be mapped uniquely to the genome: 135 to chromosome 0 between position  
920 12573795-12576534, and 45 to chromosome 9 between position 69301436-69305717, 3 to chromosome 6  
921 between position 8476136-8476138, and 1 to chromosome 4 at position 53222324. The 4 *k*-mers mapped  
922 to chromosome 4 and 6 were checked manually by locating the reads containing them and aligning the reads  
923 to the genome, in all cases no reads were able to be aligned to the genome (>99.5% of reads). For the 35  
924 *k*-mers not mapping to genome and in high LD, visualized in Figure 5E, all reads containing at least one of  
925 the *k*-mers were retrieved and assembled (SEQ4, table S2). NCBI Blast search of this fragment resulted in:  
926 positions 1-574 mapped to positions: 12578806-12579379 in chromosome 0 of the tomato genome  
927 (CP023756.1) and positions 580-1169 mapped to positions 289-878 in NSGT1 (KC696865.1). The R104  
928 “smoky” accession NSGT1 ORF starts at position 307, as reported previously (Tikunov et al., 2013). NCBI  
929 BLAST of NSGT1 (KC696865.1), identified mapping to chromosome 9 of the tomato genome  
930 (CP023765.1), from positions 975-1353 to positions 69310153-69309775.

931

### 932 **Code availability**

933 Code is available in <https://github.com/voichek/kmersGWAS>.

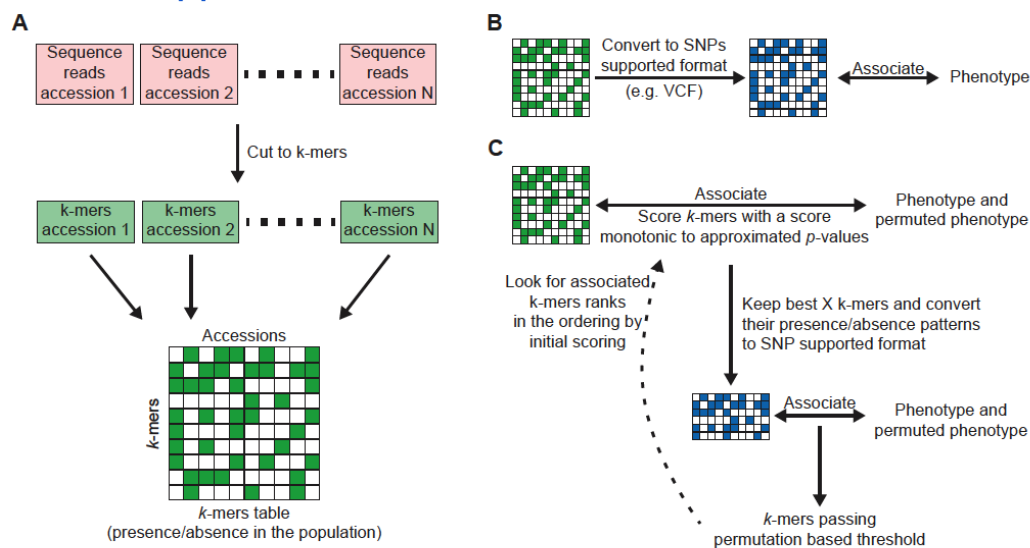
934

## Supplementary materials

935

Figure S1: Scheme of pipeline for *k*-mer-based GWAS

936



937 **(A)** Creating the *k*-mer presence/absence table: Each accession's genomic DNA sequencing reads are cut into  
 938 *k*-mers of constant length using KMC (Kokot et al., 2017). Only *k*-mers appearing at least twice/thrice in a  
 939 sequencing library are used. *k*-mers are further filtered to retain only those present in at least 5 accessions, and  
 940 ones that are also found in their reverse-complement form in at least 20% of accessions they appear in. *k*-mer lists  
 941 from all accessions are then combined into a *k*-mer presence/absence table. This table is encoded in a binary  
 942 format, with each cell represented as a single bit.

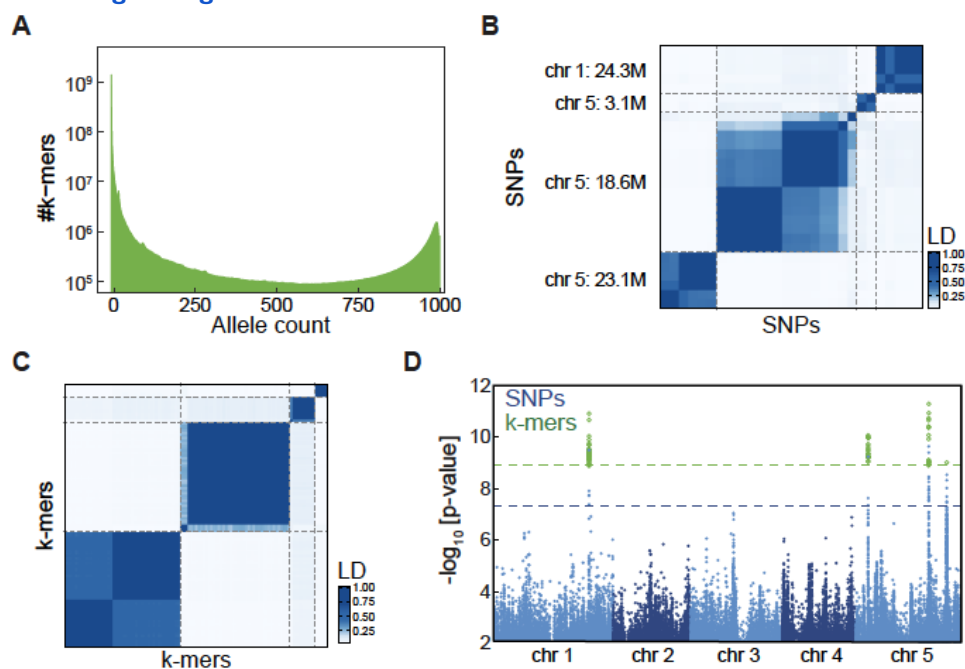
943 **(B)** Genome-wide associations on the full *k*-mer table using SNP-based software: *k*-mer table can be converted  
 944 into PLINK binary format, which can be used directly as input for association mapping in various software for  
 945 SNP-based GWA (Purcell et al., 2007; Zhou and Stephens, 2012).

946 **(C)** GWA optimized for the *k*-mers presence/absence table: *k*-mers presence/absence patterns are first associated  
 947 with the phenotype and its permutations using a linear-mixed model to account for population structure (Loh et  
 948 al., 2015; Svishcheva et al., 2012). This first step is done by calculating a score monotonic to an approximation of  
 949 the exact model. This scoring system is ultra-fast and is built for the high computational load coming from the large  
 950 number of *k*-mers and many permutations of phenotypes. Best *k*-mers from this first step (e.g. 100,000 *k*-mers)  
 951 are used in the second step. In the second step an exact *p*-value is calculated (Zhou and Stephens, 2012) for all  
 952 *k*-mers for both the phenotype and its permutations. A permutation-based threshold is calculated and all *k*-mers  
 953 passing this threshold are checked for their rank in the scoring from the first step. If not all *k*-mers hits are in the  
 954 top 50% of the initial scoring, then the entire process is rerun from the beginning, passing more *k*-mers from the  
 955 first to the second step. This last test is built to confirm that the approximation of the first step will not remove  
 956 true associated *k*-mers.

957

**Figure S2: Flowering time genetic associations in *A. thaliana* identified with *k*-mers**

958



959 **(A)** Histogram of *k*-mer allele counts: For every  $N=1..1008$ , plotted how many *k*-mers appeared in exactly  $N$   
960 accessions.

961 **(B)** LD between SNPs associated with flowering time. Dashed lines represent the four variant types, as in Figure  
962 IC.

963 **(C)** LD between *k*-mers associated with flowering time, Dashed lines represent the four variant types, as in Figure  
964 IC.

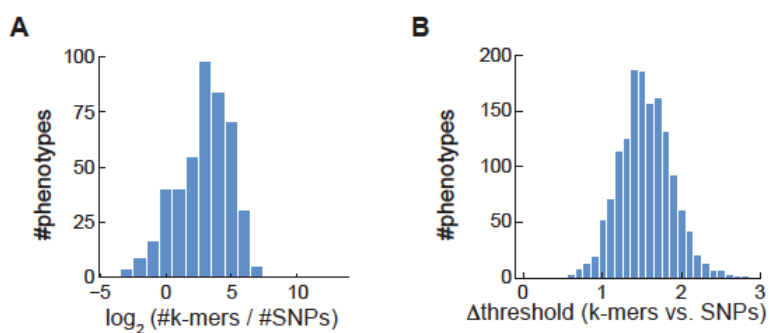
965 **(D)** Manhattan plot of SNPs and *k*-mer associations with flowering time in  $10^{\circ}\text{C}$  as in Figure ID for *k*-mers of  
966 length 25bp.

967

**Figure S3: Comparison of SNP- and  $k$ -mer-GWAS on phenotypes from 104 studies on *A. thaliana* accessions**

968

969



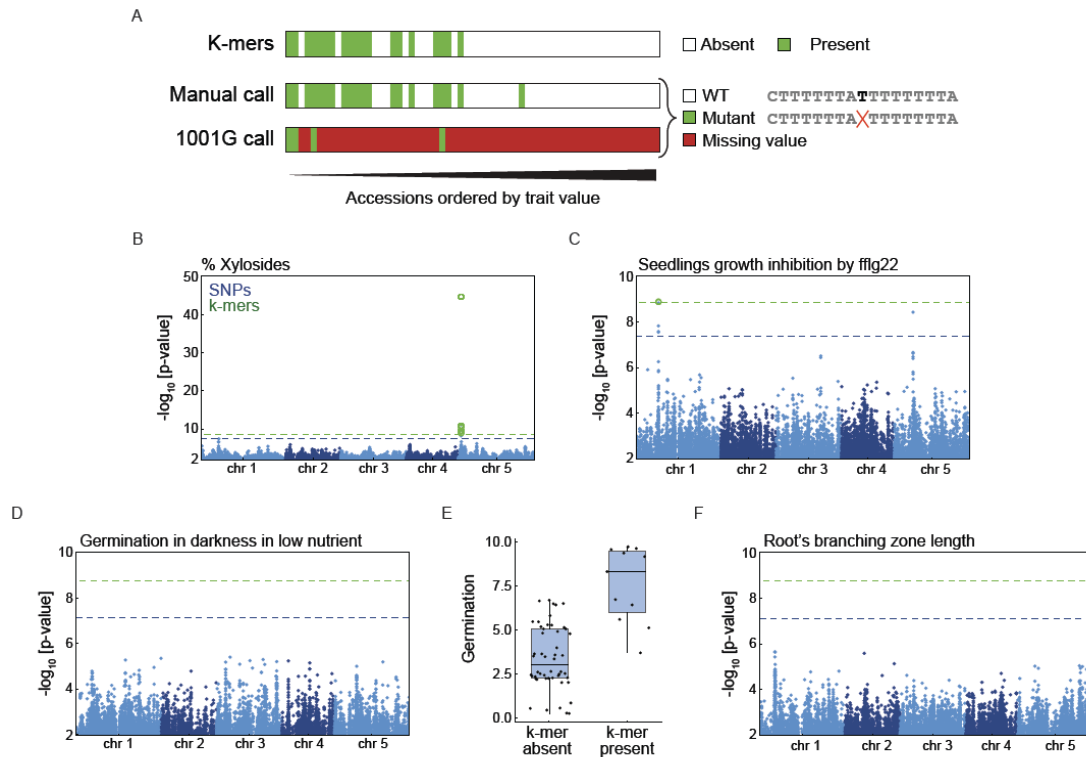
970 **(A)** Histogram of the number of identified  $k$ -mers vs. identified SNPs (in  $\log_2$ ) for *A. thaliana* phenotypes. Only  
971 the 458 phenotypes with both variant types identified were used.

972 **(B)** Histogram of thresholds difference of  $k$ -mers vs. SNPs of all *A. thaliana* phenotypes. Thresholds were  
973  $-\log_{10}$  transformed.

974

### Figure S4: Specific case studies in which *k*-mers are superior to SNPs

975



976 **(A)** Results from GWAS on measurements of lesion by *Botrytis cinerea* UKRazz strain (Fordyce et al.,  
 977 2018), an example of *k*-mers having better hold on genetic-variants present in the SNPs/indels table. We  
 978 identified 19 *k*-mers and no SNPs as being associated with this phenotype. All the *k*-mers had the same  
 979 presence/absence pattern (top row). The short sequence reads containing the *k*-mers mapped to  
 980 chromosome 3 in proximity to position 72,000. The reads contained a single T nucleotide deletion in  
 981 position 72,017, relative to the reference genome. The T nucleotide was part of an 8 T's strach, the  
 982 reference and mutated sequence around the deletion are indicated to the right of the manual calling for all  
 983 accessions (middle row) and to the calls from the 1001G project (bottom). In the 1001G only 4  
 984 accessions were called out of the 61 accessions part of the analysis, for the other accessions, the tabled  
 985 contained missing values.

986 **(B)** Manhattan plot, for xyloside percentage. A focused view on region with identified associations is  
 987 presented in Figure 3A.

988 **(C)** Manhattan plot, for seedling growth inhibition by flg22. A focused view on region with identified  
 989 associations is presented in Figure 3C.

990 **(D)** Manhattan plot, for germination in darkness in low nutrient conditions. All identified *k*-mers could  
 991 not be mapped to the genome.

992 **(E)** The germination phenotype is plotted for accessions which have the top associated *k*-mer and those  
 993 that do not.

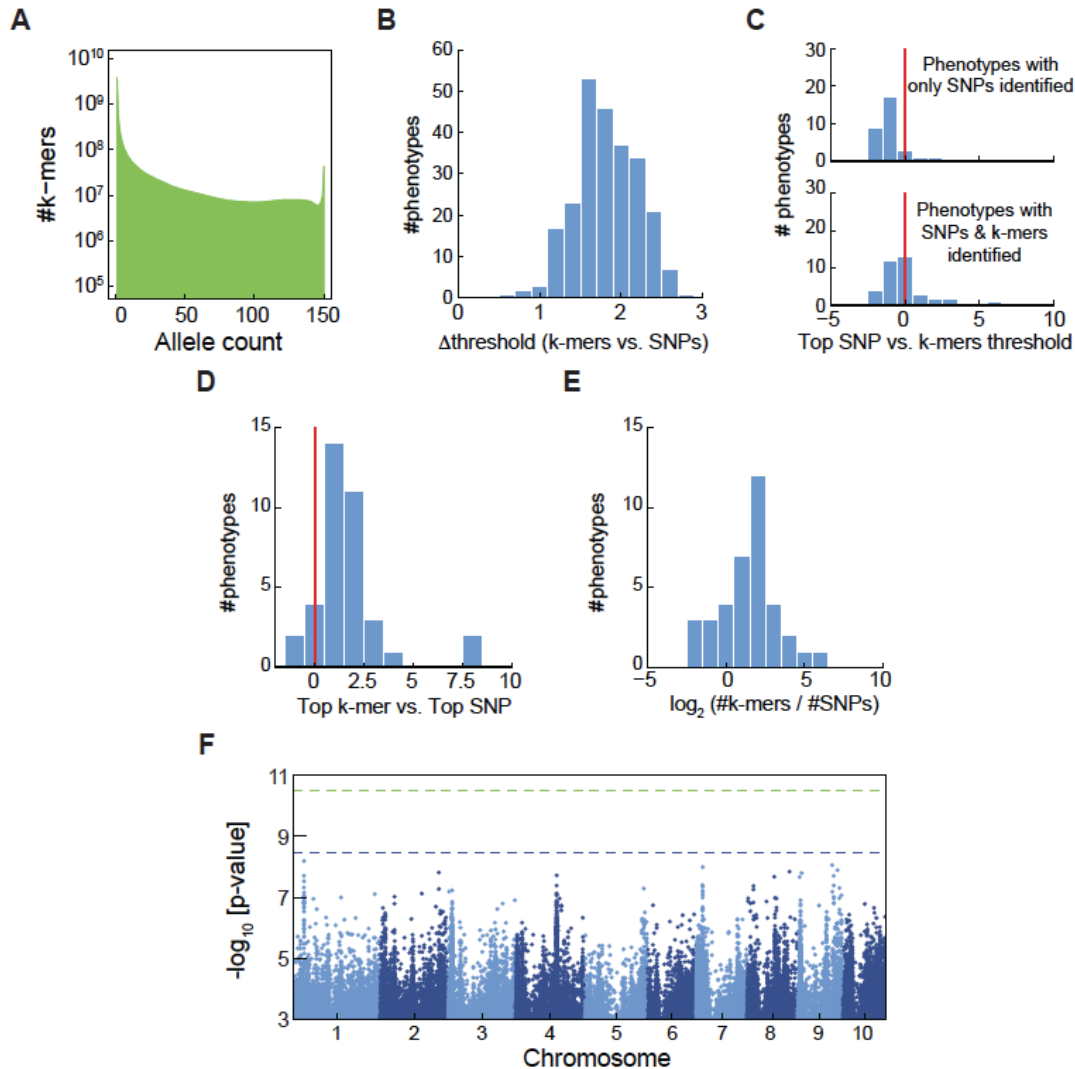
994 **(F)** Manhattan plot, for root's branching zone length. Identified *k*-mer mapped the chloroplast genome,  
 995 and thus not present in the graph.

996

997

**Figure S5: Comparison of SNP- and *k*-mer based GWAS in maize**

998



999

**(A)** Histogram of *k*-mer allele counts for maize accessions.

1000

**(B)** Histogram of difference between threshold values of SNPs and *k*-mers for maize phenotypes.

1001

**(C)** Histogram of the top SNP p-value divided by the *k*-mers defined threshold, in  $(-\log_{10})$ , for maize phenotypes. Plotted for phenotypes with only identified SNPs (upper panel) or for phenotypes with both SNPs and *k*-mers identified (lower panel).

1002

**(D)** Histogram of the difference between top  $(-\log_{10})$  p-values in the two methods for maize phenotypes identified by both methods. Plotted as in Figure 2G.

1003

**(E)** Histogram of the number of identified *k*-mers vs. identified SNPs for maize phenotypes.

1004

**(F)** Manhattan plot of association with ear weight (environment 07A). Associated *k*-mers genomic location were not located, and are thus not presented.

1005

1006

1007

1008

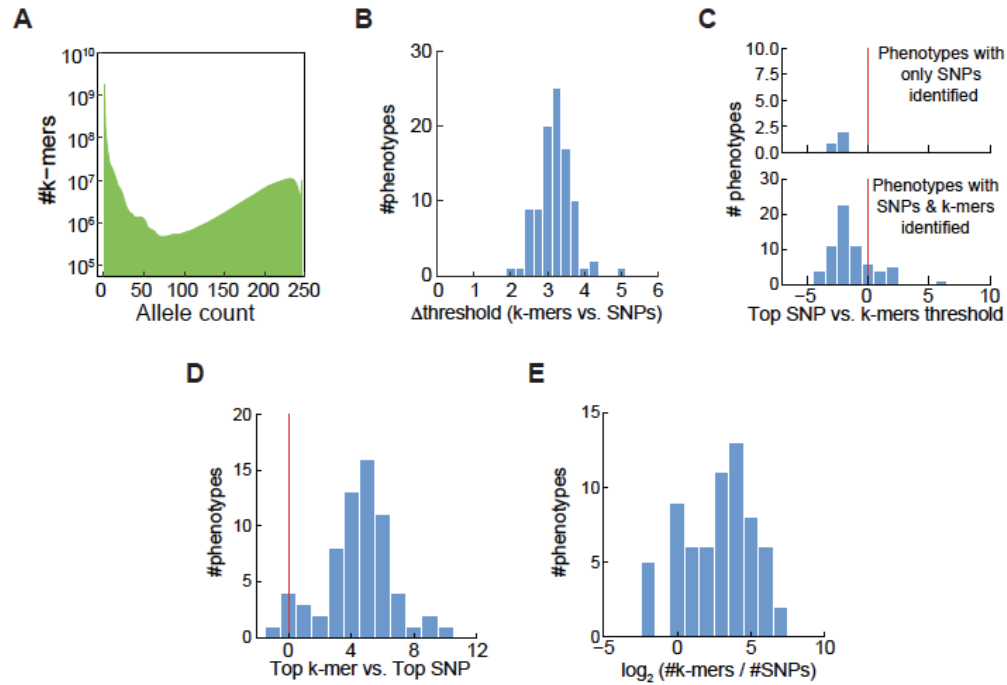
1009



1010

**Figure S6: Comparison of SNP- and *k*-mer based GWAS in tomato**

1011



1012 (A) Histogram of *k*-mers allele counts for tomato accessions.

1013 (B) Histogram of difference between threshold values of SNPs and *k*-mers for tomato phenotypes.

1014 (C) Histogram of the top SNP p-value divided by the *k*-mers defined threshold, in (-log<sub>10</sub>), for tomato phenotypes. Plotted for phenotypes with only identified SNPs (upper panel) or for phenotypes with both SNPs and *k*-mers identified (lower panel).

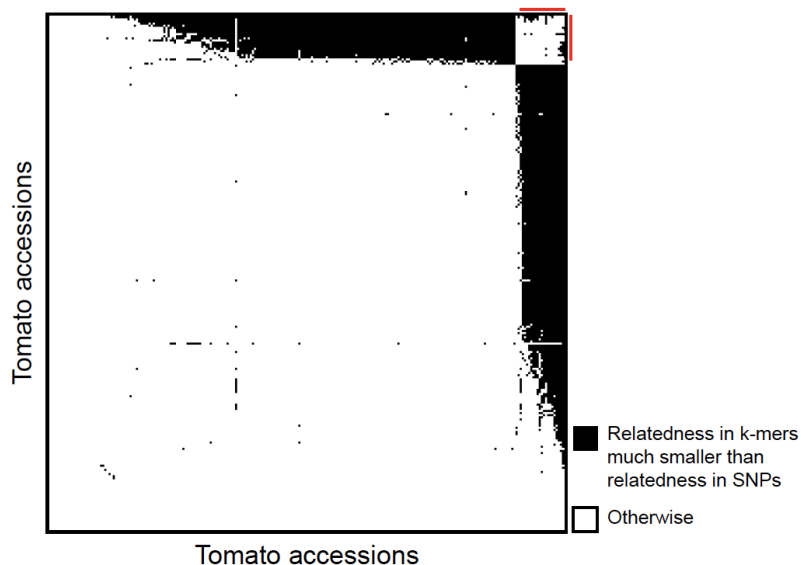
1017 (D) Histogram of the difference between top (-log<sub>10</sub>) p-values in the two methods for tomato phenotypes.

1018 (E) Histogram of the number of identified *k*-mers vs. identified SNPs for tomato phenotypes.

1019

**Figure S7: Kinship matrix calculation based on *k*-mers**

1020



1021

1022 Identification of pairs of tomato accessions for which relatedness as measured with *k*-mers is much lower  
1023 than relatedness as measured with SNPs. For every pair among the 246 accessions, a black square is  
1024 plotted if the difference in relatedness between SNPs and *k*-mers is larger than 0.15. Accessions are  
1025 ordered by the number of black square in their row/column. Red lines mark the 21 accessions with most  
1026 black squares, that is, those for which the *k*-mer/SNP difference in relatedness is larger than 0.15 for the  
most pairs.

1027

**Table S2: Assembled fragments from retrieved reads**

Sequence identifier	sequence
SEQ1	ACTGTAGCAGAAAAAATTGTTGATTGAATTAGGAGAGGCTAAGAACATTATTCG AAGTATTTCTTGTATTATTTAGATATTTACTCATTATATTGATACGGTAAGACAC AAATATGCAATTTAAAAGTTACATCACATAATTATTTGTCGCGATCCATGAATTA GGATAAGCACGAGCAACATCAATAACGTCACCTTTTCGTGGGTGAGTTCATAGA TAGTGGACAACAGTATGGAGGTTACGAATGGACAAAAGGATTAATAATAATTA TAATAGACTCTTTTATCATGTGGAACTCATGCAAGCAGAAAATGAAAGTATAT GGAGGCCGCCTCGAATCAAATTAGTTGAAAATCAGAATTAATAAATTAACGTTGT ATGGAAAAACAGAGGGTTTTTATTTTTGGGTTTTGCACAAAAAATCTTAGTCTT GAGTATTTTTGTTCAAATAAGTGTCTTTCAAGTTTCTAATATAAATTTTCAA AATTCAAACCAGCTTTATAATTTACCCCTTACCAAAGCTAATAAACTTGTTTT TTTTTTTTTATAGTATATTTATACAGTTAATTTTTTTTTTAATATTTGAAATGTGT AATA
SEQ2	CTTCTTGATTTTCATATAGAGTTCGTATACAATAATAGTTACCAAAAAAGTACTG ATACATAGTCTTACGAAGTATTGTATGGACGAGCATGTCAGACGCCCTTGATT GGACATCGGTGGACGAACAAATGCTATTTGGTTCAGAAATTGTGGACGAAACA ATAAAAAGATGAAATTCCTTTAAAGTTAAGTTAAAAGAGGTCTAAGACCGACA AAAACGTTATGCATATAGACATCGGAAGAAGCTAAAATTAAGTGGGAGATTT AGTGTACCTAAAGGCGGTGACTTACAAGGAGAGCAGACGTTTTTCCAAGAGGA AAAAGCTAATTACAAATACATGGTGCCATACAACTGCACGAACGAATTGGAG CCGTGGCTTACAAGCTTGATTTACCCTCAAAGTTGGACGCGTTTCATAAAGTTT TTCATGTATCGCAATTTAGGAAATGCCT
SEQ3	CATAAGAACATAATGATGACTGACAGGCCACTCGAACTACTCCGCATGGACCT ATTCGGCCCAATCGCTTATATAAGCATCGGCGGGAGTAAGTACTGTCTTATTAT TGTGGATGATTATTCTCGCTTCACTTAGGTATTCTTCTTGCAGGAAAAATCTCA AACCAAGAACTTTAAAGAGATTCTTGAGACGAGCTCAAATGAGTTCAGATT GAGAATCAAAAAGATTAGAAGCGATAATGGGACGGAGTTCAAGAATTCACAAA TTAAAGGATTTCTTGAGGAGGAGGGCATCAAGCATGAGTTCTCTTCTCCCTAC ACACCTCAACAAAATGGTGTAGTGGAGAGGAAGAATGGAATCTATTGGACAT GGCAAGAACCATGCTTGATGAGTACAAGACACCAGACCTGTTTTGGGCGGAG GCGATTAACACCGCCTGCTACTCCATCAACCGTTATATCTTACCAGAACTCT CAAGAAGATATCTTATGAACTCCTCACCAGTAAAAGCCCAATGTTTTATATTT AGAGTCTTTGGTAGCAAATGCTTTATTCTTGTTAAAAGAGGTAGAAGTTCTAAA TTTGCTCCTAAGGCTGTTGAAGGCTTTTTACTTGGTTATGACTCAAACACAAGG GCATATAGAGTCTTCAACAGGTCCACTGGACTAGTTGAAGTTTCTTGTGACATT GTGTTTGTGAGACTAGTGGCTCCCAAGTGGAGCAAGTTGATCTTGATGAATTA GATGATGAAGAGGCTCCGTGCATCGCGCTAAGGAACATGTCCATTGGGGATGT GTATCCTAAGGAATCCGAAGAGCCCAATAATGCACAAGATCAACCATCATCTT CCATGCAAGCATCTCCACCAACCAAGATGAGGATCAAGCTCAAG

SEQ4	GCGTAATTCTTCTCTCTACAACCGATTTTTAAGAGCGTGAGTTAGATTCAAAT ATTGATTTAACATGATATTAGATCTTTTTAATGATAGTTAACTATTTAATAGTAT GAAAATAGGGAAAAGGGTTGAAATATTACCTAACTTTGACCGAAATTGCTGTA ACAATCTCAAATTCTGATCATGACTTATTATCCGTCTGCACTATTTAATAGTGTA TTTTAAAGGAATATATATGCTCACATGGACACTTTACTATTTATAATGATGTAAT ATCTATGATGTCCACGTGTTACATATATACCTTTAAAATACACTATTAATAAT ACATGAAGTAACAAATTCTTTCAAAGTTCAGATTTGTTATAACAATTTCAATTA AATTAAGTTTTGAATATATTTCAAAAAAAGTTGCAAAAAATATAATAGGGATC TATGTCAAACCCTATGTCACCACAAGGTGGATCAAAAAAATAGTAAGAATAAA GTAATTATTGATAATGTCATTAAATTTGAAAGAGAAAGAAAAAGGTTTATAATTT TGGAGGTAGTTGTTAAAGATGGTACCTAACCTTATTCAAGCCTTTCAAATGGC TTCTTCAAATTTCCAAGCATAATTGAAACCCTAAACCTAACTTGATTATATAT GATGGGTCCAACCATGGGTAGCAACTATGGCTTCATCATAACAGTATTCATGCT ATTATGTTTTATGTTTCTTCAACTTCTGGTCTTGCCTACATTTACCACCAATTC TTCATGGGAGTTCAAGCCTTACATCTTTCCATTTTCTTCCATATACCTTCATGA CCATGAGATCAAGAAATTAGGCATACAACCAATAAAACCACGCGATGAGAAAG CTTTTGCATACATAATCCTTGAGTCTTTTGAACAATCTCACAACATTGTTTTGTT GAACACTTGTAGGGAGACTGAGGGGAAGTATATAGATTATGTTTCTACAATAG GAAAGAAAGAGTTGATACCAATTGGACCATTAATTCGCGAGGCGATGATAGAT GAGGAGGAGGATTGGGGGACAATTCAATCTTGGCTAGACAAGAAGGATCAATT ATCATGTGTTTATGTATCATTTGGAAGTGAAAGCTTCTTGTCAAAGCAAGAAAT TGAAGAGATAGCAAAGGGCTTGAGCTCAG
------	---

1028