# 1 Linked-read sequencing of gametes allows efficient genome-

# 2 wide analysis of meiotic recombination

- 3 Hequan Sun<sup>1,†</sup>, Beth A. Rowan<sup>2,†,\*</sup>, Pádraic J. Flood<sup>1</sup>, Ronny Brandt<sup>3</sup>, Janina Fuss<sup>3</sup>, Angela M.
- 4 Hancock<sup>1</sup>, Richard W. Michelmore<sup>2</sup>, Bruno Huettel<sup>3</sup>, Korbinian Schneeberger<sup>1,\*</sup>
- 5
- 6 <sup>1</sup>Department of Plant Developmental Biology, Max Planck Institute for Plant Breeding Research,
- 7 Carl-von-Linné-Weg 10, 50829 Cologne, Germany <sup>2</sup>The Genome Center and Department of Plant
- 8 Sciences, University of California, Davis, 451 East Health Sciences Drive, Davis, California 95616,
- 9 USA <sup>3</sup>Max Planck-Genome-center Cologne, Max Planck Institute for Plant Breeding Research, Carl-
- 10 von-Linné-Weg 10, 50829 Cologne, Germany
- 11
- <sup>†</sup>These authors contributed equally to this work.
- 13
- 14 \*Correspondence:
- 15 Korbinian Schneeberger (schneeberger@mpipz.mpg.de)
- 16 Beth A. Rowan (browan@ucdavis.edu)
- 17
- 18 Key words: meiotic recombination, crossover detection, single molecule, linked-read sequencing,
- 19 Arabidopsis thaliana
- 20
- 21

#### Sun and Rowan et al.

# 22 ABSTRACT

23 Meiotic crossovers (COs) ensure proper chromosome segregation and redistribute the genetic 24 variation that is transmitted to the next generation. Existing methods for CO identification are 25 challenged by large populations and the demand for genome-wide and fine-scale resolution. Taking 26 advantage of linked-read sequencing, we developed a highly efficient method for genome-wide 27 identification of COs at kilobase resolution in pooled recombinants. We first tested this method using 28 a pool of Arabidopsis F<sub>2</sub> recombinants, and obtained results that recapitulated those identified from 29 the same plants using individual whole-genome sequencing. By applying this method to a pool of 30 pollen DNA from a single  $F_1$  plant, we established a highly accurate CO landscape without 31 generating or sequencing a single recombinant plant. The simplicity of this approach now enables 32 the simultaneous generation and analysis of multiple CO landscapes and thereby allows for efficient 33 comparison of genotypic and environmental effects on recombination, accelerating the pace at 34 which the mechanisms for the regulation of recombination can be elucidated.

CO detection using linked-read sequencing

# 36 INTRODUCTION

37 During meiosis, existing genetic variation is reshuffled and passed down to offspring through 38 COs that exchange portions of homologous chromosomes. This results in new combinations of 39 alleles in the next generation that can generate novel phenotypic variation, which is the raw material 40 for both natural and artificial selection 13. In most organisms, the locations of COs along 41 chromosomes do not form a random distribution<sup>4-11</sup>. Thus, the local CO rate governs the types of allelic combinations that can arise through sexual reproduction. Both environmental and genetic 42 43 factors have been shown to affect CO rates and distributions; some examples are the local sequence divergence or rearrangements between the homologous chromosomes<sup>12-14</sup>, the chromatin 44 context<sup>15-18</sup>, variation in DNA repair mechanisms<sup>19-25</sup>, and environmental stress<sup>26,27</sup>. 45

Despite a large body of research on CO formation, our knowledge of what determines where 46 and how often COs occur is still incomplete, in part because the time, effort, and resources needed 47 48 to study this phenomenon have been limiting. Genotyping recombinant individuals, either by classical methods<sup>28</sup>, reduced representation sequencing<sup>29,30</sup>, or whole genome sequencing<sup>31</sup> and 49 performing cytological analysis of meiotic cells<sup>24</sup> represent the common methods for determining CO 50 51 locations and frequencies. However, none of these methods is easily suited for high-throughput analysis of thousands of COs in parallel. The use of recombination reporters<sup>32,33</sup> and pollen-typing<sup>34</sup> 52 53 enable rapid screening, but can only assess differences in CO frequency in a specific region of the 54 genome. The availability of an efficient method to assess the genome-wide distribution and 55 frequency of COs at a high resolution would greatly enhance our understanding of the processes 56 that govern meiotic recombination.

57 Here we investigated the use of linked-read sequencing of bulk recombinants for high-58 throughput genome-wide determination of COs in *Arabidopsis thaliana*. We first developed and 59 assessed this approach using bulked  $F_2$  individuals within known recombination sites and then

#### Sun and Rowan et al.

applied this method to hybrid pollen to generate a genome-wide CO map with a single sequencing experiment and without even growing a single recombinant plant. These results show that the time and effort needed to generate genome-wide CO maps can be reduced to sequencing and analyzing a single DNA library, making it feasible to compare multiple CO maps and thereby determine the effect of genetic and environmental factors within a single study. We believe that this method will be widely applicable to different organisms and will have a huge impact on the design of experiments aiming to decipher how meiotic recombination is regulated.

# 67 **RESULTS**

### 68 **CO breakpoint detection from bulk recombinants**

To establish a set of COs for verifying our method, we first performed whole-genome sequencing of 50 individual  $F_2$  plants derived from a cross of two of the best-studied inbred lab strains of Arabidopsis, Col-0 and L*er* and used the haplotype reconstruction software TIGER<sup>31</sup> to determine a benchmark set of 400 COs across all 50 genomes (Fig. 1; Materials and Methods; Supplementary Tables 1 and 2).

74 We then bulked the identical 50  $F_2$  plants by pooling individual leaves of comparable size and extracting high molecular weight (HMW) DNA<sup>35</sup> (Fig. 1). After size selection and quality control 75 76 (Supplementary Fig. 1), we loaded 0.25 ng DNA into a 10X Genomics Chromium Controller. The 77 Chromium Controller encapsulates millions of gel beads as GEMs (Gel bead in EMulsion), each of 78 which is loaded with a small number of long DNA molecules. These long molecules are fragmented 79 and ligated with GEM-specific DNA barcodes to generate a 10X library suitable for Illumina 80 sequencing. This library, which we called P50L25, was whole-genome sequenced with 84 million 81 151 bp-read pairs (Supplementary Table 1).

#### CO detection using linked-read sequencing

After aligning the reads against the Col-0 reference sequence<sup>36</sup> using *longranger* (v2.2.2, 82 83 10X Genomics), we recovered 3.6 million molecules ( $\geq 1$  kb) including 116 million reads (Fig. 2a; Materials and Methods). On average, these molecules were ~45 kb in size and were covered by 84 ~21 read pairs, leading to a molecule base coverage of  $\sim 0.16x$  (Fig. 2b-d; Supplementary Table 1). 85 86 To avoid chimeras resulting from the accidental co-occurrence of two independent, but closely-87 spaced molecules with identical barcodes, we selected molecules which were smaller than 65 kb 88 that had fewer than 55 reads and no heterozygous genotypes. These thresholds were based on the distributions in molecule size and read number per molecule (Fig. 2b-d; Materials and Methods; 89 90 Supplementary Note 1). Overall, 2.7 million molecules passed all filtering.

91 Initially, we genotyped these molecules at the ~660,000 SNP markers predicted by 92 longranger. If an individual molecule was composed of two distinct clusters of different parental 93 alleles and fulfilled additional criteria regarding length and marker distribution (Fig. 2a; 94 Supplementary Table 3; Supplementary Note 1), the molecule was considered as a recombinant 95 molecule revealing a CO breakpoint. Overall, we predicted 1,786 recombinant molecules with a 96 median CO breakpoint resolution (distance between the two flanking markers) of 6.7 kb. However, a 97 comparison with the 400 benchmark COs revealed that only 674 of the molecules overlapped with 98 verified COs, while the remaining 1,112 were putative false positives (FP). Many of the FPs 99 appeared close or within structural rearrangements between the parental genomes, suggesting that 100 the molecule reconstruction using linked-read alignments is vulnerable to unrecognized structural 101 differences between the parental genomes and thereby leads to false predictions of recombinant 102 molecules (Supplementary Note 1).

Making use of chromosome-level assemblies of both parental genomes<sup>36,37</sup>, we defined a new marker set composed of ~500,000 SNP markers, which only included SNPs in non-rearranged, co-linear regions between the parental genomes (Supplementary Table 4). We repeated the

#### Sun and Rowan et al.

106 analysis with this new marker set and filtered for molecules that appeared recombinant independent 107 of which parental genome was used as reference sequence. Overall, we recovered 558 108 recombinant molecules, of which 475 (85.1%) CO breakpoints (median resolution 5.3 kb) 109 overlapped with the benchmark COs. The remaining 83 (14.9%) COs were putative FPs most likely 110 due to the co-occurrence of different molecules with the same barcode (Supplementary Note 1). 111 However, we cannot rule out that some of them might have been true recombination events that were missed in the benchmark set (e.g. two closely spaced COs, gene conversion events<sup>31,38</sup> or 112 113 mitotic COs, which were only present in some cells of the sequenced material).

### 114 Increasing the number of molecules per library

115 To test whether the number of (recombinant) molecules per library would be increased by 116 increasing the DNA loading, we generated two additional 10X libraries from the same DNA pool by 117 loading 0.40 ng (P50L40) and 0.75 ng (P50L75) into the Chromium Controller (Fig. 1: Materials and 118 Methods). We sequenced these libraries with 104.8 and 212.3 million read pairs in anticipation of an 119 increased number of molecules within both libraries (Fig. 2b-d). Following the same analysis as for 120 the first library, the higher DNA loading drastically increased the number of recovered molecules. As 121 compared to the 2.7 million molecules for P50L25, we now found 4.7 and 10.5 million molecules for 122 P50L40 and P50L75 after filtering (Table 1). This also increased the number of recombinant 123 molecules to 1,012 in P50L40 and 2,519 in P50L75 as compared to 558 in P50L25. This, however, 124 had the cost of also increasing the FP rate by 5% in P50L40 and 10% in P50L75 (Table 1).

To investigate the effect of the FPs, we compared the distribution (i.e. recombination landscape) of all 2,519 COs in the library with the highest FP rate, P50L75, with two other landscapes: one calculated only from 1,874 true COs (those overlapping with the 400 benchmark COs) and one calculated only from the remaining 645 false recombinant molecules (Materials and Methods). The recombination landscape calculated for all COs (TP+FP) was nearly identical to the

#### CO detection using linked-read sequencing

130 landscape generated from true COs only (TP), while the distribution of the FPs was highly random 131 (Fig. 3a; Supplementary Fig. 2: *K-S* test, *p*-value=9.6e-01), suggesting that FPs hardly obscure true 132 recombination landscapes. In fact, correlating the sliding window-values of each of the three CO 133 landscapes revealed that the landscape calculated from all CO sites was almost perfectly correlated 134 to the one generated from real COs (Fig. 3b: Pearson's *r* 0.97, *p*-value<2.2e-16), while the 135 frequency of FP along the chromosomes was not correlated to the frequency of real COs and was 136 only marginally correlated to landscape of all COs.

137 To test the association of COs with genomic features, we checked all of the CO sites for their 138 annotation in the Col-0 reference sequence (Fig. 3c). In comparison to randomly placed CO sites. 139 the 1,874 true COs sites were significantly enriched in promoter regions (permutation test: p-value 140 5.0e-03) and intergenic regions (p-value 1.4e-03) and were significantly depleted in gene bodies (p-141 value 6.0e-04) and transposable elements (p-value 9.0e-04) recapitulating regional preferences which have been described before<sup>38,39</sup>. When additionally including the FP, the CO set showed the 142 143 same significant regional associations with one exception: COs were significantly enriched in gene 144 ends. Though this was not true when considering only the true COs, gene ends represented only a 145 minor fraction of all COs. To examine CO associations with genomic features at a finer scale, we 146 analyzed the regional preferences of the CO sites in individual transposable element super families because these were found to have strong differences in CO rates<sup>40</sup>. Though this only included 331 147 148 putative CO sites, we found that COs were significantly depleted within LTR-Gypsy and En-Spm 149 super families and enriched within Helitron and LINE elements, as has been previously shown<sup>40</sup>, 150 whether or not FP were included (Supplementary Figure 3). In addition, the presence of FP did not affect the previously reported relationships between COs and GC content or DNA methylation<sup>38</sup>. The 151 152 two datasets showed nearly identical correlations between CO frequency and GC content (Fig. 3d: 153 Pearson's r -0.46 and -0.45, both p-value<2.2e-16). Similarly, COs in both datasets were found in 154 regions with low levels of methylation (Fig. 3e).

#### Sun and Rowan et al.

Together this suggests that our method is not only effective when facing an increasing amount of FP in libraries with a large number of molecules, but that it is also powerful enough to accurately identify chromosome-wide and local CO patterns.

### 158 Increasing the number of genomes per library

159 Another technical consideration is that many of the recombinant molecules were derived 160 from the same CO breakpoints, given the limited number of distinct COs within the 50-individual 161 pools. For instance, there were only 363 distinct COs recovered by the 1,874 recombinant 162 molecules in P50L75 (Table 1). Even though identifying a single CO breakpoint multiple times can 163 help to increase its resolution and reliability, it does reduce the overall number of distinct COs that 164 can be found with one library. Though the inclusion of more plants (i.e. increasing the number of 165 genomes in the pool) does not increase the number of recombinant molecules, it does increase the 166 number of *independent* recombinant molecules and thereby maximizes the number of distinct CO 167 events found with one library.

168 To test the effect of genome number on the number of distinct COs, we applied our method 169 to pools of 200 and 1,250 F<sub>2</sub> plants. We generated two individual 10X libraries each for a pool of 170 200 (P200R1, P200R2) and a pool of 1250 (P1250R1 and P1250R2) different  $F_2$  plants (Fig. 1; 171 Supplementary Figures 1 and 4; Materials and Methods). These libraries were sequenced with 172 168.2 - 194.4 million read pairs (Supplementary Table 1) and revealed 12.3 - 16.0 million 173 molecules. For P200R1 and P200R2, we found 2,538 and 2,334 recombinant molecules, which 174 identified 1,779 (69.5%) and 1,606 (68.8%) distinct CO sites. For the pool of 1,250 plants, there 175 were 2,880 and 3,430 recombinant molecules, which uncovered 2,386 (82.8%) and 2,788 (81.3%) 176 distinct CO sites. Comparison of these COs with 1,840 COs determined using individual whole 177 genome sequencing for an independent set of Col x Ler F<sub>2</sub>s<sup>20</sup> showed consistent genome-wide 178 patterns (Supplementary Figure 5). While the larger pools revealed more distinct CO sites, there

#### CO detection using linked-read sequencing

were still nearly 20% COs overlapping, but it does not necessarily imply that these are not unique CO events. As the genome size of *A. thaliana* is relatively small (~135 Mb<sup>41</sup>), independent COs have an unneglectable probability to overlap (Supplementary Note 2; Supplementary Table 5). According to simulations, given a pool of 1,250  $F_2s$ , ~15% of the independent CO events are expected to overlap (Supplementary Table 5 and 6) making up for large parts of the overlapping recombinant molecules in the 1,250 pools.

### 185 Estimating relative recombination frequency

186 Since the probability to identify a CO breakpoint within a molecule depends on the average 187 recombination frequency among the pooled genomes, it might be possible to estimate the average 188 recombination frequency from the fraction of observed recombinant molecules. However, the 189 probability to identify a CO breakpoint also depends on the length and sequencing coverage of the 190 molecules within a library. It is therefore only meaningful to calculate an average recombination 191 frequency that is relative to the underlying molecule characteristics (relative recombination 192 frequency). To test for significant differences in relative recombination frequencies between libraries, 193 it is essential that they have close-to-identical molecule characteristics. For libraries with differences 194 in molecule length and sequencing coverage distributions, subsampling can be used to generate 195 identical distributions. Once the distributions are similar, the relative recombination frequency can be 196 determined as the number of recombinant molecules per million molecules ( $C^{M}$ ). Repeated 197 subsampling even allows for the calculation of confidence intervals for  $C^{M}$  values and thereby allows 198 for testing significant differences between the average recombination frequencies of different pools.

First, we compared the CO frequencies of P200R1 and P200R2, which were two independent libraries generated from the same DNA (Fig. 1; Fig. 2b-d; Materials and Methods). After subsampling (Supplementary Figure 6), there was no significant difference in  $C^{M}$  values (26.7±0.4 and 26.8±0.4), as expected for these libraries (Table 2). We repeated this comparison for

#### Sun and Rowan et al.

203 P1250R1 and P1250R2, which differed greatly in their original molecule characteristics (Fig. 2b-d). 204 After subsampling to the same molecule distributions as for the smaller pools (Supplementary Figure 6), the  $C^{M}$  values of 27.2±0.6 and 27.1±0.7 were also not significantly different from each 205 206 other (Table 2). Moreover, when comparing the 200 and 1,250 recombinant pools, we also found no 207 significant difference between any of the pools (all confidence intervals overlapped), which is 208 expected, given that all libraries were generated from individuals of the same  $F_2$  population. Together, these results show that  $C^{M}$  values are stable against differences in the molecule 209 210 characteristics and pool sizes and allow for determining and comparing average recombination 211 frequencies between samples.

### 212 Estimating recombination frequency and landscapes from gametes

213 Crosses between divergent strains provide the simplest opportunity for determining COs, but 214 typically require the generation of a recombinant population after a single round of meiosis in a 215 hybrid context. Observing COs directly in the products of meiosis (gametes) would greatly expedite 216 the study of recombination, especially in inbred species with long generation times. To test how our 217 method performs on recombinant gametes, we extracted HMW DNA from Col x Ler F<sub>1</sub> hybrid pollen, 218 created a 10X library (P8000) and sequenced it with 319.9 million read pairs (Fig.1; Supplementary 219 Figure 1f; Supplementary Table 1). Following the same analysis as for the  $F_2$  pools, we identified 20 220 million molecules with an average molecule base coverage of 0.26x (Fig. 2b-d). Among those, there 221 were 3.246 recombinant molecules with a median CO breakpoint resolution of 8.0 kb.

We compared relative recombination frequencies in pollen with the P1250R1 and P1250R2 pools after re-sampling molecules to obtain comparable characteristics (Supplementary Figure 6; Materials and Methods). The CO frequency estimate ( $C^{M}$ ) for pollen was significantly higher than for F<sub>2</sub>s (Table 2), consistent with the higher male recombination rate in this species<sup>42</sup>.

#### CO detection using linked-read sequencing

The distributions of the CO breakpoints in pollen and  $F_2s$  were highly correlated in most regions of the genome (Pearson's *r* 0.77~0.83, both *p*-value<2.2e-16), but exhibited a few local differences (Fig. 4). The increased COs in pollen therefore followed a similar genome-wide distribution as in  $F_2s$ , as the regions with locally enriched COs in pollen accounted for only a minor fraction of the all COs.

231 We further compared the recombination landscape of our pollen sample to a previously 232 published recombination landscape from a Col-0 x Ler backcross population, where all COs were 233 derived from male meiosis<sup>42</sup>. This data set, consisting of 7,418 COs at an average resolution of ~320 kb. was generated from 1,505 individuals genotyped with 380 SNP markers that were evenly 234 235 spaced throughout the genomes. After binning our CO data into the same windows, we found that 236 the recombination landscapes were broadly similar (Supplementary Figure 7: Pearson's r 0.41, p-237 value<2.2e-16). Unexpected differences, however, were observed mostly at the end of 238 chromosomes. COs formed by male meiosis and identified in recombinant BCF<sub>1</sub> plants showed significant increases at the end of some chromosomes<sup>42</sup>, while most of these increases were not 239 240 observed in pollen data.

In general, this shows that a single sequencing library is sufficient to derive genome-wide CO patterns from pooled gametes without the need for producing any recombinant individuals. Such CO maps have the advantage to measure sex-specific meiotic recombination without doing any backcrosses and are less influenced by post-meiotic biases.

# 245 **DISCUSSION**

For the past century, the study of recombination has relied on inferring or determining the genotypes at (limited) marker positions along the chromosomes of recombinant individuals. Here, we developed an efficient and accurate method for detecting recombination breakpoints in pooled DNA using linked-read sequencing. This method can be applied directly to gametes, avoiding the

#### Sun and Rowan et al.

need for generating or genotyping recombinant populations. In turn, this simplification allows for the production of multiple different genome-wide CO maps with relatively little effort. Thus, it is now feasible to directly compare CO maps generated from many different genotypes or environments. Long-read sequencing or single-cell DNA-seq of individual pollen would be two alternative methods, but the current costs of library preparation and sequencing are prohibitive for comparing multiple samples.

256 The small genome size of the species we used in this study, A. thaliana, was potentially an 257 obstacle to accurately scoring CO events using linked-read sequencing, as this technology assigns 258 the same barcode to a small number of individual DNA molecules that are tens of kb in length. Since 259 this process is random, two molecules that originate from different genomes in the pool that happen 260 to align near each other can receive the same barcode. The probability of these "collision" events is 261 much higher in small genomes and can mimic COs if the two molecules have different parental 262 genotypes, giving rise to false positives (FP). However, the FP had little effect on the overall 263 accuracy of the CO predictions (Fig. 3). Even though the FP rate did not pose a substantial problem, 264 making several libraries using a reduced input would result in high molecule numbers at the lower 265 FP rate. This would incur higher costs of library production while maintaining the same sequencing 266 cost. Moreover, when applying this method to larger genomes, the chance of "collision" of two 267 molecules with the same barcode would be lower, leading to a lower likelihood of FPs.

The spatial distributions of COs identified with our method recapitulated known recombination landscapes and precisely co-localized with genomic and epigenetic features that have been reported to be associated with meiotic recombination<sup>34,37,39</sup>. The resolution of the CO breakpoints was very precise, with the median CO interval less than 10 kb in each of the samples. This accuracy, however, could only be achieved by filtering the marker list for co-linear regions between the parental genomes as genomic rearrangements between the parents were hotspots for FPs.

#### CO detection using linked-read sequencing

The careful validation studies of our method using  $F_2$ s provided the basis for determining the COs directly in gametes. Linked-read sequencing of bulk  $F_1$  pollen led to the discovery of ~3,500 CO events that showed a similar distribution as was recently published for male meiosis (Supplementary Figure 7). Dreissig et al. previously assessed the recombination rate in barley pollen using whole-genome sequencing of individual pollen nuclei<sup>46</sup>. While this allowed for the analysis of CO interference, their method employed a more technically challenging library preparation that does not scale to high numbers and achieved a much lower CO resolution.

282 Interrogating COs in gametes (pollen, in the case of plants or sperm in animals) has many 283 advantages over current approaches. It avoids the laborious and time-consuming step of 284 growing/rearing recombinant populations. It reduces the number of generations required, greatly 285 facilitating recombination studies in inbred organisms with long generation times. Because entire 286 recombination landscapes can be generated from single libraries, a recombination landscape can 287 now be studied as a single trait. Multiple CO landscapes that can be replicated either in the same or 288 different genetic backgrounds and environments and compared in a single study. In consequence, 289 this now allows for a more complex and sophisticated experimental design to test hypotheses 290 regarding the regulation of recombination.

### 291 Materials and Methods

### 292 **F**<sub>2</sub> **DNA extraction and library preparation**

F<sub>2</sub> seeds from Col-0 x L*er*-0 were stratified for 7 days at 4°C, sown on soil in 24-pot trays, and grown under 16 h light, 8 h dark cycles at 20°C for three weeks. DNA pools of up to 1,250  $F_2$ individuals were constructed (Fig. 1). HMW DNA was extracted from pools with 50 distinct  $F_2$  plants using a published protocol<sup>35</sup>. One 50- $F_2$  pool was selected for validating the method, for which DNA was also extracted from each individual (a<sub>1-50</sub>) and WGS libraries were prepared using the Illumina

#### Sun and Rowan et al.

298 DNA TruSeq protocol. The DNA of four 50-F<sub>2</sub> pools with similar fragment size distribution 299 (Supplementary Figure 4a) were merged based on equal concentration, from which replicates 300 P200R1 and P200R2 were obtained. In addition, 25 50-F<sub>2</sub> pools with divergent fragment size 301 distributions (Supplementary Figure 4b) were merged based on equal molarity of molecules 302 between 42-70 kb according to FEMTOpulse, AATI genomic DNA quality check, and the resultant 303 DNA was used for replicates P1250R1 and P1250R2.

304 Size selection was performed on each DNA pool using the Sage Science BluePippin high 305 pass protocol (i.e., 0.75% Agarose Dye-Free / 0.75% DF Marker U1 high-pass 30-40 kb vs3) with 306 starting point at 40 kb to ensure that most molecules were over 40 kb (Supplementary Figure 1a-e). 307 For P200R1/2 and P1250R1/2, the size-selected DNA was evaluated using the Qubit fluorometer 308 and TapeStation analyzer. After selection, for each library 0.75 ng DNA were loaded into the 10X 309 Chromium controller (Supplementary Table 1). P50 was subjected to the same quality control 310 measures. For this library three different amounts of DNA (0.25 ng, 0.40 ng and 0.75 ng) were 311 loaded into the 10X Chromium Controller, generating the P50L25, P50L40, P50L75 libraries.

#### 312 **Pollen DNA extraction and library preparation**

313 Col CEN3 grt 420 and Ler seeds were stratified for 4-7 days at 4°C before sowing on soil in 314 18-pot trays and growing under 16 h light, 8 h dark cycles at 20°C until flowering. Pollen from a 315 single Ler plant was used to pollinate a single Col CEN3 qrt 420 stigma to generate F1 plants. To obtain pollen and extract DNA, we adapted a method from Drouaud and Mezard<sup>47</sup> as follows: 316 317 Inflorescences were collected from a single F<sub>1</sub> plant and ground in 1 mL of 10% sucrose using a 318 mortar and pestle. 9 mL of 10% sucrose were added to the mortar slurry and the resulting 319 homogenate was mixed by pipetting up and down with a wide bore 1 mL pipet tip and filtered 320 through an 80-µM nylon mesh before centrifugation at 350 x g for 10 min at 4°C. The supernatant 321 was discarded and the pollen pellet was washed two times with 10% sucrose. The pellet was

#### CO detection using linked-read sequencing

322 resuspended in four volumes of lysis buffer (100 mM NaCl, 50 mM Tris-HCl (pH 8)), 1 mM EDTA, 323 1% SDS and proteinase K was added to achieve a concentration of 20 µg/mL. Five to ten 2-mm 324 glass beads were added to the sample and it was vortexed at full speed for 30 s. To check for pollen 325 disruption, a 1-uL sample was removed and mixed with 10 uL of lysis buffer and examined under a 326 microscope. During this time, we verified that the pollen sample was free of large amounts of cell 327 debris. The sample was then vortexed for 30 seconds and checked for pollen disruption. An equal 328 volume of Tris-saturated phenol was added and the sample was placed on a rotating wheel at room 329 temperature for 30 min. After centrifuging at 15,000 x g for 10 min, the supernatant was transferred 330 to a new tube and mixed with an equal volume of 24:1 chloroform: isoamvlalcohol and homogenized 331 by shaking the tube. The tube was centrifuged again at 15,000 x g for 10 min and the supernatant 332 was mixed with 0.7 volumes isopropanol in a new tube and inverted gently. After another 333 centrifugation step at 15,000 x g for 10 min, the supernatant was discarded and the pellet was 334 washed with 1 mL of 70% EtOH. After a final centrifugation at 15,000 x g for 2 min, the supernatant 335 was discarded and the DNA pellet was allowed to dry at room temperature. The pellet was 336 resuspended in 50 µL of 10 mM Tris-HCl, pH 8 with 0.1 mM EDTA and stored at 4°C.

337 DNA was analyzed by field inversion gel electrophoresis to confirm that most molecules 338 were around 48 kb (Supplementary Figure 1f) before preparing the linked-read library (load: 1.00 339 ng) using the Chromium Genome Reagent Kit, the Chromium Genome Library Kit & Gel Bead Kit, 340 Chromium Genome Chip Kit v2, and Chromium i7 Multiplex Kit according to the manufacturer's 341 instructions. As 1.00 ng is equivalent to 10<sup>6</sup> Mb, the haploid genome size is ~135 Mb<sup>41</sup>, and each 342 pollen grain is composed of three cells, there were ~3,000 pollen grains expected with unique 343 recombinant genomes.

In total, eight 10X and 50 standard Illumina libraries were sequenced on HiSeq 3000/4000 with 151 bp paired-end reads, where the individual  $F_2s$  were at ~5x and the pools of  $F_2s$  and pollen nuclei were at 173-659x (Supplementary Table 2).

#### Sun and Rowan et al.

### 347 Molecule recovery using linked-read alignments

348 Col-0 and Ler reference genomes were indexed using the function mkref of longranger 349 (v2.2.2 10X Genomics). Linked-reads of each sample were aligned against the Col-0 and Ler 350 reference genomes separately using longranger wgs with options: --id=ALIGN-ID --351 reference=MKREF-INDX-FOLDER --fastqs=READ-PATH --sample=READ-ID --localcores=40 --352 localmem=192 --noloupe --sex=male --vcmode=freebayes --library=LIBRARY-ID and all other 353 options under default settings (note that option --sex was set as "male", which is required by the tool 354 but not affecting the alignment here). Read alignments with the same read barcodes were clustered 355 to molecules according to their genomic location ensuring that neighboring read alignments were 356 not more than 25 kb away from each other. The resultant molecules (over 1 kb) were further filtered 357 for falsely merged molecules by removing very long or densely covered molecules (Supplementary 358 Note 1).

359 Generating the CO benchmark set

360 For each individually sequenced F<sub>2</sub>, read alignments against the Arabidopsis Col-0 reference sequence<sup>36</sup> and variant calling were performed using *Bowtie2*<sup>43</sup> (version 2.2.8) and 361 SAMtools/BCFtools<sup>44</sup> (version 1.3.1). TIGER<sup>31</sup> was used to reconstruct the parental haplotypes 362 using SNPs in co-linear regions between the Col-0 and Ler<sup>37</sup> genomes. We determined an average 363 of 8.3 COs per diploid genome, totaling 415 COs with a median breakpoint resolution of 0.5 kb 364 365 (Supplementary Table 2). Of those, 15 COs were located in regions with an inter-marker distance of 366 more than 10 kb. As local CO breakpoint identification relies on markers near the breakpoint, we 367 excluded these COs, generating a final benchmark set of 400 COs.

CO detection using linked-read sequencing

#### 368 Estimating chromosomal CO distributions (recombination landscape)

Recombination landscapes were estimated using sliding windows along each chromosome (window size 1 Mb, step size 50 kb). Within a given window, the CO frequency *f* was calculated by *f*  $= n / t_i$ , where *n* is the number of COs in the window and  $t_i$  is the total CO number within the respective chromosome *i* across each of the focal libraries. This ensures that different CO landscapes can be compared to each other within the given chromosome.

### 374 **Relative CO frequency estimation**

375 For frequency estimation, we first subsampled molecules from the pools of interest by 376 intersecting their molecule characteristic distributions. Using the molecule length distributions as the 377 basis (Supplementary Figure 6a), the molecules were randomly sampled down within 1 kb large 378 bins. For each bin, the number of randomly selected molecules per pool was equal to the lowest 379 number of molecules across all pools within this specific bin. To increase the subsampling rate, 80% 380 of the sampled molecules were further randomly selected within each bin. After this second 381 subsampling, the read and base coverage distributions were also highly similar (Supplementary 382 Figure 6b-c). Final molecule filtering was applied following the same strategy as used for the complete sets of molecules, by applying 30 kb as size and 24 as read number thresholds. 383

Then, COs per million molecules ( $C^{M}$ ) was calculated for each pool. The process of molecule sub-sampling, CO identification and  $C^{M}$  calculation was repeated fifty times for each pool. The means ( $\mu$ ) of all 50  $C^{M}$  values and their 95% confidence intervals ( $\mu \pm s$ , where *s* is 2.0096\* $\sigma$ /50^0.5 with  $\sigma$  being the standard deviation) were used to assess significant differences between samples.

### 388 **DNA methylation level estimation**

389 Within a recent study, we have assessed DNA methylation for *A. thaliana* Col-0<sup>45</sup>. Methylation 390 level *M* for a CO or a random interval was calculated by  $M = N_{met} / N$ , where  $N_{met}$  is the number of

#### Sun and Rowan et al.

reads supporting methylated cytosines at all C and G sites, while *N* is the total number of reads at these sites.

# 393 Data availability

Read data of all eight 10X linked-read libraries (ERS2851779-ERS2851786) and 50 whole genome sequencing libraries (ERS2851943-ERS2851992) are available in BAM format from European Nucleotide Archive under accession number ERP111558.

# 397 Acknowledgements

398 The authors would like to thank Ian R. Henderson (Department of Plant Sciences, University 399 of Cambridge) for providing the CO breakpoint lists, Erik Wijnker (Wageningen University) for 400 providing seeds and Ulrike Hümann, Manish Goel, Wen-Biao Jiao, Vidya Oruganti, and Onur Dogan 401 (Max Planck Institute for Plant Breeding Research) for their help in the greenhouse. We also would 402 like to thank 10X Genomics for their help on setting up the longranger software, advice on DNA extraction, and their kind donation of library reagents to support the development of recombination 403 404 identification. We thank Lutz Froenicke and the DNA Technologies Core at UC Davis for 10X library 405 sequencing support and discussions. We also acknowledge helpful discussions with Detlef Weigel 406 at the Max Planck Institute for Developmental Biology and Kyle Fletcher, William Palmer, and 407 Sebastian Reves-Chin-Wo at UC Davis. We thank Felicity Jones and Frank Chan (Friedrich 408 Miescher Laboratory of the Max Planck Society) for inspiring the extension of this work to gametes.

# 409 **Author contributions**

H.S., B.A.R. and K.S. designed the project. H.S. and B.A.R. performed all analysis. B.A.R.,
H.S. and P.J.F. prepared the samples. B.A.R., R.B., J.F. and B.H. performed DNA extraction, quality

CO detection using linked-read sequencing

- 412 control, library preparation and sequencing. K.S., B.A.R., A.M.H. and R.M.W. supervised the project.
- 413 H.S., B.A.R. and K.S. wrote the manuscript. All authors read and approved the final manuscript.

# 414 Funding

- 415 This work was supported by a Max Planck Society postdoctoral fellowship and a UC Davis
- 416 Genome Center Pilot Project grant.

# 417 **Competing Interests**

418 None.

## 419 **Code availability**

- 420 Custom code used for identification of recombinant molecules and frequency calculation can
- 421 be found online at *https://github.com/schneebergerlab/Scripts\_supporting\_Sun\_Rowan\_et\_al\_2019*.

# 422 **References**

- 423 1. Barton, N. H. & Charlesworth, B. Why sex and recombination? *Science* **281**, 1986-1990 (1998).
- 424 2. Rice, W. R. & Chippindale, A. K. Sexual recombination and the power of natural selection.
  425 Science **294**, 555-559 (2001).
- 426 3. McDonald, M. J., Rice, D. P. & Desai, M. M. Sex speeds adaptation by altering the dynamics of
  427 molecular evolution. *Nature* 531, 233-236 (2016).
- 428 4. Mancera, E., Bourgon, R., Brozzi, A., Huber, W. & Steinmetz, L. M. High-resolution mapping of 429 meiotic crossovers and non-crossovers in *yeast. Nature* **454**, 479-485 (2008).
- Kauppi, L., Jeffreys, A. J. & Keeney, S. Where the crossovers are: recombination distributions in
  mammals. *Nat. Rev. Genet.* 5, 413-424 (2004).
- 432 6. Yamada, S. et al. Genomic and chromatin features shaping meiotic double-strand break

Sun and Rowan et al.

- 433 formation and repair in mice. *Cell Cycle* **16**, 1870-1884 (2017).
- 434 7. Singhal, S. *et al.* Stable recombination hotspots in birds. *Science* **350**, 928-932 (2015).
- 435 8. Miller, D. E. *et al.* Whole-Genome Analysis of Individual Meiotic Events in Drosophila
  436 melanogaster Reveals That Noncrossover Gene Conversions Are Insensitive to Interference and
- 437 the Centromere Effect. *Genetics* **203**, 159-171 (2016).
- 438 9. Kianian, P. M. A. *et al.* High-resolution crossover mapping reveals similarities and differences of
  439 male and female recombination in *maize*. *Nat. Commun.* 9, 2370 (2018).
- 440 10. Salomé, P. A. *et al.* The recombination landscape in *Arabidopsis thaliana* F2 populations.
  441 *Heredity* 108, 447-455 (2012).
- 442 11. Higgins, J. D., Osman, K., Jones, G. H. & Franklin, F. C. H. Factors underlying restricted
  443 crossover localization in barley meiosis. *Annu. Rev. Genet.* 48, 29-47 (2014).
- 444 12. Sturtevant, A. H. A crossover reducer in *Drosophila melanogaster* due to inversion of a section
  445 of the third chromosome. *Biol. Zent. Bl.* 46, 697-702 (1926).
- 446 13. Dobzhansky, T. The Decrease of Crossing-Over Observed in Translocations, and Its Probable
  447 Explanation. *Am. Nat.* 65, 214–232 (1931).
- 448 14. Chakraborty, U. & Alani, E. Understanding how mismatch repair proteins participate in the
  449 repair/anti-recombination decision. *FEMS Yeast Res.* 16, (2016).
- 450 15. Choi, K. *et al. Arabidopsis* meiotic crossover hot spots overlap with *H2A.Z* nucleosomes at gene
  451 promoters. *Nat. Genet.* 45, 1327-1336 (2013).
- 452 16. Underwood, C. J. *et al.* Epigenetic activation of meiotic recombination near *Arabidopsis thaliana*
- 453 centromeres via loss of H3K9me2 and non-CG DNA methylation. *Genome Research* (2018).
- 454 doi:10.1101/gr.227116.117.
- 455 17. Yelina, N. E. *et al.* Epigenetic remodeling of meiotic crossover frequency in *Arabidopsis thaliana*456 DNA methyltransferase mutants. *PLoS Genet.* 8, e1002844 (2012).
- 457 18. Marand, A. P. et al. Meiotic crossovers are associated with open chromatin and enriched with

CO detection using linked-read sequencing

- 458 Stowaway transposons in potato. *Genome Biol.* **18**, 203 (2017).
- 459 19. Ziolkowski, P. A. *et al.* Natural variation and dosage of the HEI10 meiotic E3 ligase control
   460 *Arabidopsis* crossover recombination. *Genes Dev.* **31**, 306-317 (2017).
- 461 20. Serra, H., Lambing, C., Griffin, C.H., Topp, S.D., Nageswaran D.C., Underwood C.J., Ziolkowski,
- 462 P.A., Séguéla-Arnaud, M., Fernandes, J.B., Mercier R. and Henderson I. Massive crossover
- 463 elevation via combination of *HEI10* and *recq4a recq4b* during *Arabidopsis* meiosis. *PNAS* 115,
  464 2437-2442 (2018).
- 465 21. Fernandes, J. B., Seguéla-Arnaud, M., Larchevêque, C., Lloyd, A. H. & Mercier, R. Unleashing
  466 meiotic crossovers in hybrid plants. *PNAS* **115**, 2431-2436 (2017).
- 467 22. Girard, C. *et al.* FANCM-associated proteins *MHF1* and *MHF2*, but not the other Fanconi anemia
  468 factors, limit meiotic crossovers. *Nucleic Acids Res.* 42, 9087-9095 (2014).
- 469 23. Crismani, W. et al. FANCM limits meiotic crossovers. Science 336, 1588–1590 (2012).
- 470 24. Higgins, J. D. *et al. AtMSH5* partners *AtMSH4* in the class I meiotic crossover pathway in
  471 *Arabidopsis thaliana*, but is not required for synapsis. *Plant J.* **55**, 28-39 (2008).
- 472 25. Berchowitz, L. E., Francis, K. E., Bey, A. L. & Copenhaver, G. P. The role of *AtMUS81* in
  473 interference-insensitive crossovers in *A. thaliana*. *PLoS Genet.* **3**, e132 (2007).
- 474 26. De Storme, N. & Geelen, D. The impact of environmental stress on male reproductive
  475 development in plants: biological processes and molecular mechanisms. *Plant Cell Environ.* 37,
  476 1-18 (2014).
- 477 27. Jackson, S., Nielsen, D. M. & Singh, N. D. Increased exposure to acute thermal stress is
  478 associated with a non-linear increase in recombination frequency and an independent linear
  479 decrease in fitness in *Drosophila*. *BMC Evol. Biol.* **15**, 175 (2015).
- 480 28. Konieczny, A., and F. M. Ausubel. A procedure for mapping Arabidopsis mutations using co481 dominant ecotype-specific PCR-based markers. *Plant J.* 4, 403-410 (1993).

Sun and Rowan et al.

- 482 29. Baird, N.A., Etter, P.D., Atwood, T.S. *et al.* Rapid SNP discovery and genetic mapping using
  483 sequenced RAD markers. *PLoS One* **3**, e3376 (2008).
- 30. Rowan, B.A., Danelle, K.S., Eunyoung, C., Derek. S.L. and D. Weigel. Methods for Genotypingby-Sequencing. *Methods in Molecular Biology* **1492**, 221-42 (2017).
- 31. Rowan, B. A., Patel, V., Weigel, D. & Schneeberger, K. Rapid and inexpensive whole-genome
  genotyping-by-sequencing for crossover localization and fine-scale genetic mapping. *G3* 5, 385-
- 488 398 (2015).
- 32. Melamed-Bessudo, C., Yehuda, E., Stuitje, A. R. & Levy, A. A. A new seed-based assay for
  meiotic recombination in *Arabidopsis thaliana*. *Plant J.* **43**, 458-466 (2005).
- 33. Francis, K. E. *et al.* Pollen tetrad-based visual assay for meiotic recombination in *Arabidopsis*. *PNAS* **104**, 3913-3918 (2007).
- 34. Choi, K., Yelina, N. E., Serra, H. & Henderson, I. R. Quantification and Sequencing of Crossover
   Recombinant Molecules from *Arabidopsis* Pollen DNA. in *Haplotyping: Methods and Protocols*

495 (eds. Tiemann-Boege, I. & Betancourt, A.) **1551**, 23-57 (Springer New York, 2017).

- 496 35. Mayjonade, B., Gouzy, J., Donnadieu, C., Pouilly, N., Marande, W., Callot, C., Langlade, N.,
- 497 Muños, S. Extraction of high-molecular-weight genomic DNA for long-read sequencing of single
  498 molecules. *Biotechniques* 61, 203-205 (2016).
- 36. The Arabidopsis Genome Initiative. Analysis of the genome sequence of the flowering plant
  Arabidopsis thaliana. *Nature* 408, 796-815 (2000).
- 37. Zapata, L., Ding, J., Willing, E.M., Hartwig, B., Bezdan, D., Jiao, W.B., Patel, V., Velikkakam,
  J.G., Koornneef, M., Ossowski, S., Schneeberger, K. Chromosome-level assembly of *Arabidopsis thaliana* Ler reveals the extent of translocation and inversion polymorphisms. *PNAS*113, E4052-60 (2016).

CO detection using linked-read sequencing

- 38. Wijnker, E., James, G.V., Ding, J., Becker, F., Klasen, J.K., Rawat, V., Rowan, B.A., de Jong,
  D.F., de Snoo, C.B., and Zapata, L. et al. The Genomic Landscape of Meiotic Crossovers and
  Gene Conversions in *Arabidopsis Thaliana*. *Elife* 2, e01426 (2013).
- 39. Shilo, S., Melamed-Bessudo, C., Dorone, Y., Barkai, N., Levy, AA. DNA Crossover Motifs
  Associated with Epigenetic Modifications Delineate Open Chromatin Regions in *Arabidopsis*. *Plant Cell* 27, 2427-36 (2015).
- 511 40. Choi, K., Zhao, X., Tock, A. J., Lambing, C., Underwood, C.J., Hardcastle, T.J., Serra, H., Kim,
- 512 J., Cho, H. S., Kim, J. et al. Nucleosomes and DNA Methylation Shape Meiotic, DSB Frequency
- in *Arabidopsis Thaliana* Transposons and Gene Regulatory Regions. *Genome Research* 28,
  532-546 (2018).
- 41. Sun, H., Ding, J., Piednoël, M., Schneeberger, K. *findGSE*: estimating genome size variation
  within human and Arabidopsis using *k*-mer frequencies. *Bioinformatics* **34**, 550-557 (2018).
- 42. Giraut, L., Falque, M., Drouaud, J., Pereira, L., Martin, O.C., Mézard, C. Genome-wide
  crossover distribution in *Arabidopsis thaliana* meiosis reveals sex-specific patterns along
  chromosomes. *PLoS Genet.* 7, e1002354 (2011).
- 43. Langmead, B. and Salzberg, S. L. Fast gapped-read alignment with *Bowtie* 2. *Nat Methods* 9,
  357-9, (2012).
- 44. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, H., Marth, G., Abecasis, G.,
  Durbin, R., 1000 Genome Project Data Processing Subgroup. The sequence alignment/map
  format and *SAMtools. Bioinformatics* 25, 2078-9 (2009).
- 45. Willing, E. M., Rawat, V., Mandáková, T., Maumus, F., James, J. V., Nordström, K. J., Becker,
  C., Warthmann, N., Chica, C., Szarzynska, B. et al. Genome expansion of *Arabis alpina* linked
  with retrotransposition and reduced symmetric DNA methylation. *Nature plants* 1, 14023 (2015).
  doi: 10.1038/nplants.2014.23.

Sun and Rowan et al.

- 529 46. Dreissig, S., Fuchs, J., Himmelbach, A., Mascher, M. and Houben, A. Sequencing of Single
- 530 Pollen Nuclei Reveals Meiotic Recombination Events at Megabase Resolution and Circumvents
- 531 Segregation Distortion Caused by Postmeiotic Processes. *Front Plant Sci.* **8**, 1620 (2017).
- 532 47. Jan, D. and Mézard, C. Characterization of Meiotic Crossovers in Pollen from Arabidopsis
- 533 Thaliana. *Methods in Molecular Biology* 745, 223-49 (2011).

CO detection using linked-read sequencing

## 535 48. Table 1. Effect of increasing the molecule number of increased DNA loading

| Pool  | P50L25    | P50L40    | P50L75     |
|---|-----------|-----------|------------|
| Raw molecules (≥1kb) <sup>C</sup>           | 3,577,104 | 5,804,559 | 12,551,833 |
| Raw molecules (≥1kb) <sup>L</sup>           | 3,462,796 | 5,631,294 | 12,240,762 |
| Filtered molecules <sup>C</sup>             | 2,664,977 | 4,783,828 | 10,651,147 |
| Filtered molecules <sup>L</sup>             | 2,565,403 | 4,642,515 | 10,417,489 |
| Total CO molecules <sup>CL</sup>            | 558       | 1,012     | 2,519      |
| TP CO molecules <sup>CL</sup>               | 475       | 804       | 1,874      |
|   | (85.1%)   | (79.4%)   | (74.4%)    |
| TP unique (non-redundant) COs <sup>CL</sup> | 254       | 325       | 363        |
|   | (63.5%)   | (81.3%)   | (90.8%)    |
| FP molecules <sup>CL</sup>                  | 83        | 208       | 645        |
|   | (14.9%)   | (20.6%)   | (25.6%)    |

536 \*TP, true positive; FP, false positive

537 \*C: number based on Col-0 reference genome<sup>36</sup>

538 \*L: number based on Ler reference genome<sup>37</sup>

\*CL: number based on intersection of two sets of predictions (of C and L)

Sun and Rowan et al.

|       | Pool    |                 | Molecules <sup>a</sup> | CO molecules <sup>b</sup> | $C^{M}$  |
|-------|---------|-----------------|------------------------|---------------------------|----------|
|       |         | P200R1          | 4,845,588±1,556        | 129±7                     | 26.7±0.4 |
| $F_2$ | F       | P200R2          | 4,856,886±1,371        | 130±7                     | 26.8±0.4 |
|       | P1250R1 | 4,650,420±1,662 | 126±9                  | 27.2±0.6                  |          |
|       |         | P1250R2         | 4,309,936±1,454        | 117±10                    | 27.1±0.7 |
|       | Pollen  | P8000           | 4,902,994±1,722        | 212±10                    | 43.2±0.6 |

# 540 Table 2 Relative recombination frequency in F<sub>2</sub> and pollen pools

541 \*Note: for **a** and **b**, the values are  $\mu \pm \sigma$ , where  $\mu$  is the mean and  $\sigma$  is the standard deviation of the molecule numbers in 50

random sub-samplings. For  $C^{M}$ , the values are  $\mu \pm s$  giving 95% confidence intervals for the mean  $\mu$  of  $C^{M}$ , where s is

543 2.0096\* $\sigma$ /50^0.5 with  $\sigma$  being the standard deviation.

544

CO detection using linked-read sequencing

# 546 Figure legends

#### 547 Figure 1. Experimental design and CO detection using linked-reads.

F<sub>1</sub> and F<sub>2</sub> plants were derived from crosses of two divergent *A. thaliana* accessions. *Leaf sampling*: 548 549 Leaves of 50 selected  $F_2$ s were individually sampled ( $a_{1\sim50}$ ). In addition, the leaves of all other plants 550 were pooled in batches of 50 F<sub>2</sub>s (for further merging), where the 50 individually sampled plants 551 formed one of the pools. Also, pollen from a single  $F_1$  plant was sampled. DNA extraction: DNA from 552 individual the samples  $a_{1\sim50}$  were extracted for Illumina whole-genome sequencing, while the DNA of 553 the 50-F<sub>2</sub> and pollen samples were extracted with a protocol for high molecular weight DNA (see Materials and Methods and Mayjonade et al. (2016)<sup>35</sup>). Pooling: Four and twenty-five 50-F<sub>2</sub> samples 554 were merged leading to 200-F<sub>2</sub> and 1,250-F<sub>2</sub> pools (P200 and P1250), each with two replicates (R1 555 556 and R2). The 50-F<sub>2</sub> pool composed of the 50  $F_2$ s that were individually sequenced was labelled as 557 P50, while the pollen pool was labelled as P8000. Library preparation: 50 individual  $F_2$  DNA 558 samples were used for preparing Illumina DNA TruSeg libraries. The P50 DNA sample was loaded 559 into a 10X Chromium Controller (illustration modified from 10X Genomics) with three different 560 amounts of DNA including 0.25 ng, 0.40 ng and 0.75 ng (P50L25, P50L40 and P50L75). 561 P200R1/R2 and P1250R1/R2 were loaded using 0.75 ng and P8000 was loaded using 1.00 ng. 562 Sequencing: All libraries were sequenced on Illumina HiSeg3000/4000 sequencers.

**Figure 2. Molecule recovery and genotyping, and real molecule characteristics.** *a.* Reads with the same barcode (represented by short lines in matching colors) that were aligned within close proximity (less than 25 kb apart) have been connected to recover molecules (M-1 to M-4). The genotypes that can be assessed with the read alignments of each recovered molecule fall into three major categories, non-recombinant (M-1, M-2), recombinant (M-3) and undetermined (M-4). *b.* Length of recovered molecules for each of the 10X libraries. *c.* Number of reads per recovered molecule. *d.* Recovered molecule base coverage.

#### Sun and Rowan et al.

570 Figure 3. Genome-wide CO landscape formation and feature association. a. Sliding window-571 based (window size 1 Mb, step size 50 kb) recombination landscapes calculated for true positives 572 (TP), false positives (FP) and the combined set of TP+FP recombinant molecules of P50L75. b. 573 Correlation of regional CO frequencies comparing TP v.s. TP+FP, FP v.s. TP+FP and FP v.s. TP. c. 574 Association of the TP and TP+FP sets with different genomic features in contrast to the random 575 expectation. d. Correlation of regional GC-content with CO frequency within the TP and TP+FP 576 sets. e. Association of TP or TP+FP sets with DNA methylation, compared with a random 577 expectation (Materials and Methods). Error bars reflect the standard deviation with respective to the 578 mean of methylation of a set of CO intervals. Note: in c and e, the expectations were obtained 579 based on randomly sampled intervals within the reference genome excluding heterochromatic 580 regions. The middle position of a CO interval was used for associations with genomic features in c. 581 where the promoter of a gene was defined as the 1,000 bp region upstream of the transcription start 582 site, and the gene start/end as the first/last 200 bp of a gene. For each feature, a permutation test 583 (oneway test in R "coin" package) was performed between the TP set (or TP+FP set) with 1,000 584 sets of randomly sampled 3,000 intervals. The asterisks indicate the observed values (either from 585 the TP or the TP+FP set) were significantly different from a random expectation.

Figure 4. Genome-wide recombination landscape in pollen. *a*. Comparison of sliding windowbased (window size 1 Mb, step size 50 kb) recombination landscapes in pollen and  $F_2$  populations (P1250R1/R2). *b*. Correlation of genome-wide CO frequencies in pollen and  $F_2$ s (Pearson's *r* 0.77~0.83, both *p*-value<2.2e-16).

590







Relative CO frequency (scaled per chromosome)

