

1 Genomic Rearrangements Considered as Quantitative Traits

2

3 Martha Imprialou¹, André Kahles², Joshua G. Steffen³, Edward J. Osborne³,
4 Xiangchao Gan⁴, Janne Lempe⁴, Amarjit Bhomra¹, Eric Belfield⁵, Anne
5 Visscher^{5,6}, Robert Greenhalgh³, Nicholas P Harberd⁵, Richard Goram⁷,
6 Jotun Hein⁸, Alexandre Robert-Seilaniantz⁹, Jonathan Jones¹⁰, Oliver
7 Stegle¹¹, Paula Kover¹², Miltos Tsiantis⁴, Magnus Nordborg¹³, Gunnar
8 Rättsch², Richard M. Clark^{3,14}, Richard Mott^{1,15}

9

10 ¹ Wellcome Trust Centre for Human Genetics, University of Oxford, OX3 7BN,
11 UK

12 ² Memorial Sloan-Kettering Cancer Center, New York City, NY 10065, USA

13 ³ Department of Biology, University of Utah, Salt Lake City, UT, 84112-0840,
14 USA

15 ⁴ Max Planck Institute for Plant Breeding Research, 50829 Köln, Germany

16 ⁵ Department of Plant Sciences, University of Oxford, Oxford, OX1 3RB, UK

17 ⁶ Department of Comparative Plant and Fungal Biology, Royal Botanic
18 Gardens Kew, Ardingly, RH17 6TN, UK

19 ⁷ John Innes Centre, Norwich, NR4 7UH, UK

20 ⁸ Department of Statistics, University of Oxford, OX1 3TG, Oxford, UK

21 ⁹ UMR INRA-Agrocampus Ouest-Université de Rennes 1, 35653 Le Rheu
22 Cedex, France

23 ¹⁰ The Sainsbury Laboratory, Norwich Research Park, Norwich NR4 7UH, UK.

24 ¹¹ European Bioinformatics Institute, Hinxton, Cambridge, CB10 1SD, UK

- 1 ¹² Dept of Biology and Biochemistry, University of Bath, Bath, BA2 7AY, UK
- 2 ¹³ Gregor Mendel Institute of Molecular Plant Biology, Vienna, 1030, Austria
- 3 ¹⁴ Center for Cell and Genome Science, University of Utah, Salt Lake City, UT,
- 4 84112-0840, USA
- 5 ¹⁵ UCL Genetics Institute, University College London, WC1 6BT, UK
- 6

1 Abstract

2

3 To understand the population genetics of structural variants (SVs), and their
4 effects on phenotypes, we developed an approach to mapping SVs,
5 particularly transpositions, segregating in a sequenced population, and which
6 avoids calling SVs directly. The evidence for a potential SV at a locus is
7 indicated by variation in the counts of short-reads that map anomalously to the
8 locus. These SV traits are treated as quantitative traits and mapped
9 genetically, analogously to a gene expression study. Association between an
10 SV trait at one locus and genotypes at a distant locus indicate the origin and
11 target of a transposition. Using ultra-low-coverage (0.3x) population
12 sequence data from 488 recombinant inbred *Arabidopsis* genomes, we
13 identified 6,502 segregating SVs. Remarkably, 25% of these were
14 transpositions. Whilst many SVs cannot be delineated precisely, PCR
15 validated 83% of 44 predicted transposition breakpoints. We show that
16 specific SVs may be causative for quantitative trait loci for germination, fungal
17 disease resistance and other phenotypes. Further we show that the
18 phenotypic heritability attributable to sequence anomalies differs from, and in
19 the case of time to germination and bolting, exceeds that due to standard
20 genetic variation. Gene expression within SVs is also more likely to be
21 silenced or dysregulated. This approach is generally applicable to large
22 populations sequenced at low-coverage, and complements the prevalent
23 strategy of SV discovery in fewer individuals sequenced at high coverage.

24

1 Introduction

2

3 Whilst genome resequencing has become cheap and ubiquitous, and can
4 readily determine variations such as SNPs and very small indels, the problem
5 of identifying structural variants (SVs) and rearrangements remains a
6 challenge, despite continual improvement in algorithms for calling SVs. The
7 current gold standard for determining SVs between individuals is by *de-novo*
8 assembly¹. This requires very high-coverage paired-end sequence over a
9 range of insert sizes, together with long-range information for scaffolding.
10 Advances in long-read technologies^{2,3} are beginning to aid this process, but
11 the relatively high cost and low throughput of this strategy limits its
12 applicability to smaller numbers of genomes, and leaves open two important
13 questions. First, whether an SV identified in an individual is unique, or is
14 frequent enough to contribute appreciably to phenotypic heritability in a
15 population. Second, whether an SV represents a local rearrangement, such
16 as a deletion, inversion or tandem copy-number variant (CNV), or is long-
17 range, such as a transposition^{4,5}.

18

19 SVs are often revealed by the anomalous alignment of short-reads to the
20 reference genome. Specific anomaly signatures characterize different types of
21 SVs (Table 1). Thus, same-strand pairs indicate inversion, high read coverage
22 duplications, abnormal insert sizes and unpaired reads are classified as
23 indels.

1 These signatures include excess read coverage (e.g., duplications, Copy
2 Number Variants (CNV)), discordant distances between read pairs (e.g.,
3 indels) and inconsistent read orientation (e.g. inversions). These anomalies
4 arise, often in combination, because the reads have been aligned to the
5 wrong genome – the anomalies should disappear if instead the reads were
6 aligned to the true genome. This idea is used by algorithms such as GATK⁶
7 and Platypus⁷ that identify small indels by local realignment, and in whole-
8 genome reassembly by iterative realignment⁸.

9

10 Many SV-calling algorithms utilize these read-anomaly signatures to identify
11 SVs segregating in individuals sequenced at high coverage^{9–16}. These
12 methods focus on short-range SVs because of the difficulties in distinguishing
13 long-range rearrangements from read mapping errors. They are also designed
14 to work best when calling SVs in individuals sequenced at intermediate to high
15 coverage; for example, two of the most recent SV-callers, LUMPY¹⁵ and
16 WHAM¹⁶ are most sensitive when sequence coverage is at least 10x. In other
17 applications, e.g cancer resequencing, typical coverage is even higher, at 30x
18 or above.

19

20 The problem of calling SVs from population sequence data presents additional
21 challenges. Population studies are generally conducted for the purpose of
22 genetic association, and consequently require large sample sizes. Population
23 sequencing provides an alternative to genotyping by SNP arrays,
24 simultaneously providing both haplotype reference panels for imputation¹⁷ and

1 cohorts for disease mapping^{18,19}. As the sample size increases, it becomes
 2 possible to reduce the coverage of each individual dramatically, yet still
 3 impute single nucleotide polymorphism (SNPs) accurately²⁰. Consequently
 4 one would want to be able to call SVs as well as SNPs and to test them for
 5 association. Although the information present in each sample is sparse, and
 6 therefore it would be difficult to call SVs (and SNPs) on an individual basis, by
 7 pooling information across samples it might be possible to determine common
 8 SVs analogously to the way SNPs are imputed.

9

10 A further challenge, which is not confined to low-coverage sequencing, is that
 11 presented by complex SVs. Unlike simple indels, inversions and
 12 transpositions, where a segment with well-defined breakpoints is affected,
 13 many SVs are composites of multiple events²¹, often driven by transposons
 14 and other repetitive mobile elements. Complex SVs resist simple
 15 classification, and it may be impossible to determine the precise sequence of
 16 mutations that occurred in the lineages separating the reference genome from
 17 that of the sequenced individual. Whilst current algorithms for calling SVs in
 18 simulated high-coverage human data can identify simple SVs with sensitivities
 19 of about 90% depending on the type of SV¹⁶, they are less accurate when
 20 applied to real data, and their performance on complex SVs is unreported.

21

22 Nonetheless, even though it may be difficult to delineate complex SVs, there
 23 can still be strong evidence from read-mapping anomalies that an SV of some
 24 sort exists at a locus. If the intensity of its anomaly signature can be used as a

1 proxy for the purpose of testing genetic association, then one need not
2 delineate the SV precisely. It then follows that the genomewide information
3 captured by these anomalies could be used to compute relationships between
4 individuals based on their structural similarities alone, and hence to estimate
5 the heritability attributable to this source of variation.

6

7 Here, we ask whether low-coverage population sequencing provides new
8 ways for mapping SVs and estimating heritability, complementing the
9 sequencing of fewer individuals at high coverage. As an illustration, we
10 investigate the architecture and phenotypic impact of structural variation in
11 *Arabidopsis thaliana*. Among natural accessions of *Arabidopsis*, structural
12 variation is plentiful ²². The extent of rDNA repeats ²³ and mobile transposable
13 elements ²⁴ vary between accessions, and variation in the overall amounts of
14 both classes of repetitive sequence elements are complex traits, partially
15 under genetic control. In this study we investigate all types of structural
16 variation in *Arabidopsis*, including those not mediated by mobile elements.
17 We show that long-range transpositions are common, and that structural
18 variation has a significant impact on particular quantitative trait loci (QTLs)
19 and on trait heritability, distinct from that explained by other types of sequence
20 variation.

21

22

23 **Results**

24

1 **Structural Variants as Quantitative Traits**

2

3 We combined established ideas from signature-based SV identification with
4 quantitative genetics to analyse structural variation in population sequence
5 data. The following scenario motivates our reasoning: suppose an SV arose in
6 a certain population ancestor, α , transposing a genomic segment s originating
7 at a “source” locus L and targeting to a “sink” locus M . Source and sink can be
8 coincident or unlinked, but for the moment, suppose they are unlinked. If the
9 event is transposon-mediated, then the segment s is duplicated to s' at M , and
10 possibly altered, leaving the original s at L . Once random chromosomal
11 assortment and recombination has occurred, in the present-day population
12 there will be a mix of individuals carrying the segment at neither, one or both
13 loci.

14

15 In the descendent population, one individual is sequenced and becomes the
16 reference genome. Depending on the choice of reference individual and the
17 mechanism of transposition, the reference might carry zero or one copies of s
18 at the source and of s' at the sink.

19

20 Assume the reference has one copy of s and zero copies of s' . A population
21 sample will contain individuals with mix of all possible configurations at source
22 and sink. But only individuals that inherited the haplotype descended from α
23 at the sink carry the transposed segment, regardless of their haplotype at the
24 source. The individuals are sequenced with short-reads, and the reads are

1 mapped to the reference genome. Individuals carrying the transposition s' at
 2 the sink will have reads spanning the breakpoint that split between source and
 3 sink. Hence read mapping anomalies apparently originating at the source will
 4 be enriched in those individuals carrying the sink haplotype α : genotypes that
 5 tag α at the sink will be associated with anomalies at the source.

6

7 If on the other hand the reference contains both s at the source and s' at the
 8 sink then those individuals that did not inherit the haplotype α at the sink will
 9 appear to carry a deletion there. Reads with anomalously large insert sizes
 10 will map to the sink and will be associated with genotypes tagging the
 11 haplotype α at the sink – the generative role played by the source will be
 12 invisible.

13

14 Similarly, by considering situations where the source and sink are coincident –
 15 for example tandem duplications – in a population we would expect to
 16 encounter a mix of short-range *cis* and long-range *trans* associations between
 17 different classes of read-mapping anomalies and genotypes, depending on
 18 the diverse histories of each structural variants.

19

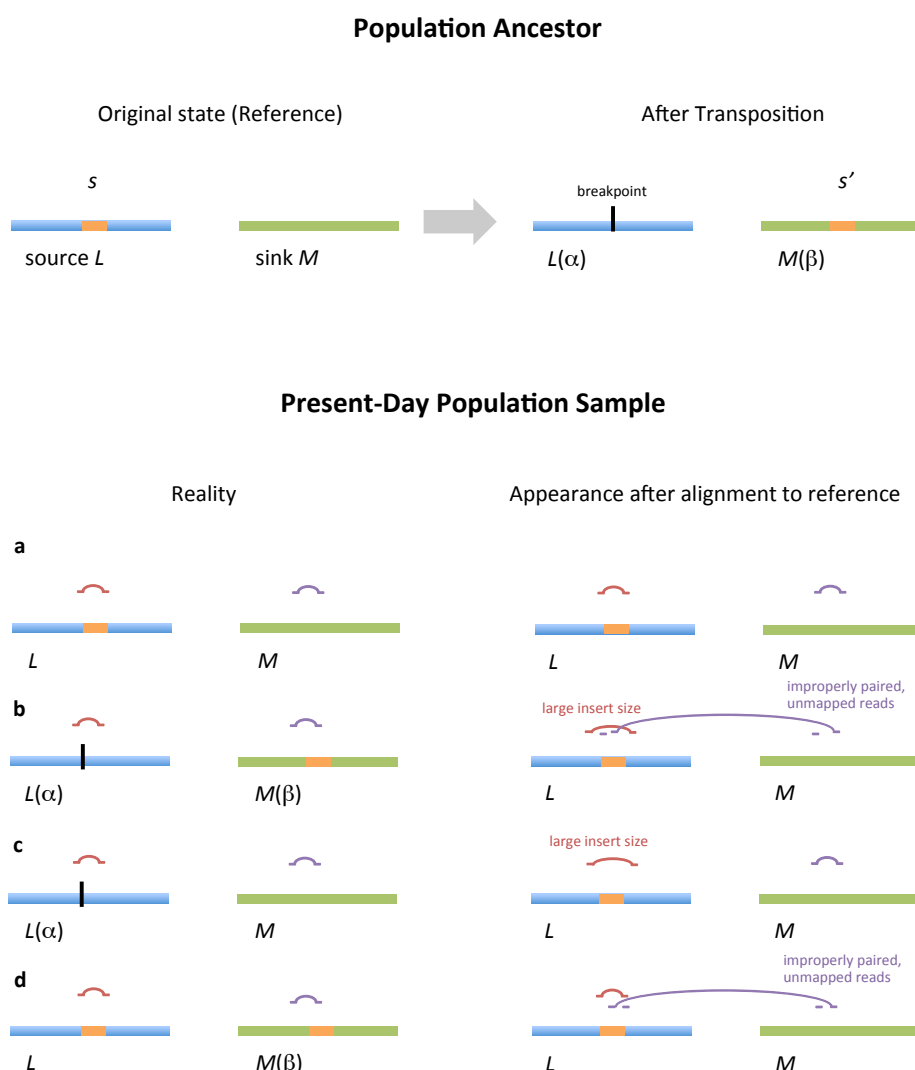


Figure 1 Effects of a transposition on short-read mapping. Chromosomes are horizontal bars and read pairs are pairs of horizontal lines linked by curves. Upper shows a population ancestor corresponding to the reference genome (left) undergoing a transposition (right), in which a segment s at source locus L with haplotype context α is copied to s' at recipient sink locus M with haplotype context β . Lower shows all four possible combinations (a-d) of source L and sink M haplotype in descendants. On left are shown the alignment of reads to the true haplotypes, where there are no read-mapping anomalies. On right are shown the various read-mapping anomalies that arise, depending on the true haplotype backgrounds at source and sink, upon alignment to the reference genome.

1 To apply these ideas in practice, we count the numbers of anomalous reads
 2 mapping to each source L in a population sample, treat it as a quantitative
 3 trait, and proceed to identify genetic loci containing variation that correlate
 4 with variation in the trait. This procedure defines a SV quantitative trait locus
 5 (SV-QTL) linking T_{Li} , the number of anomalous reads mapping to locus L in
 6 individual i and the haplotype H_{Mi} at sink locus M in individual i (**Figure 1**,
 7 **Methods**). *cis* SV-QTLs where the source and sink overlap indicate local
 8 structural variants such as CNVs, deletions and inversions; *trans* SV-QTLs
 9 indicate transpositions (insertional translocations) or larger scale
 10 rearrangements. In this way we can determine whether an SV is in *trans*, its
 11 originating haplotype, which individuals now carry it (**Figure S1**), and its
 12 frequency (**Figure S2**).

13

14 We interpret the matrix of SV-traits across all loci as a Euclidean
 15 representation of haplotype space, in the sense that, if two individuals are
 16 genetically similar then their SV-trait vectors should be close together.
 17 Consequently we define a genome-wide similarity between individuals based
 18 on the similarity of their anomalies, as a weighted average of their locus-
 19 specific similarities. Taken across all individuals, these generate a structural
 20 variation similarity matrix, analogous to a SNP-based genetic relationship
 21 matrix. This matrix was used to estimate the heritability of a phenotype with
 22 respect to structural variation, and compared to the heritability associated with
 23 SNP variation.

24

1 **Structural Variation in Arabidopsis**

2

3 We used our strategy to map *cis* and *trans* SVs in the 120Mb genome of the
 4 plant *Arabidopsis thaliana*. We sequenced 488 of the *Arabidopsis* Multiparent
 5 Advanced Generation Inter-Cross (MAGIC) recombinant inbred lines²⁵ at
 6 ~0.3x coverage using 51bp paired-end Illumina reads. The MAGIC lines
 7 descend from 19 ancestral founder accessions that have been sequenced at
 8 high coverage⁸ (**Table S1**) such that each ~120Mb genome is a mosaic of the
 9 19 founder haplotypes. Consequently we expect most SVs segregating in
 10 MAGIC to also segregate in the founders, thereby providing a means to verify
 11 any SVs we detect. The choice of MAGIC lines rather than natural accessions
 12 means that the confounding effects of population structure and of selection
 13 are largely absent from the population. Very rare alleles with frequency below
 14 1/19=4.5% are uncommon, increasing the power to detect QTLs. However,
 15 MAGIC QTL mapping resolution is also poorer, at ~200kb, compared to ~10kb
 16 in natural accessions.

17

18 We mapped the reads to the TAIR10 reference using Stampy^{26,27} and
 19 inferred the mosaic of each line using a hidden Markov model (HMM)
 20 implemented in the software ‘reconstruction’ available on request from the
 21 authors. The algorithm uses as input SNP calls for each MAGIC genome, and
 22 the set of 1.2M biallelic variants in the 19 founders (excluding loci tagged as
 23 within transposons, and those sites called as heterozygous or multi-allelic in
 24 the founders)¹⁷, and finds the most likely sequence of haplotype assignments

1 for each chromosome. Because the lines were called at low coverage, most
2 SNP sites were not covered by reads in an given; consequently we called on
3 average 301k SNPs per line (using GATK⁶) (ie a randomly sampled of ~25%
4 of the 1.2M sites). However, this amount of data is sufficient for the HMM to
5 determine the founder mosaic accurately; we estimated by simulation that the
6 algorithm can delineate the mosaic breakpoints (which correspond to
7 recombination events) to within ~2kb (data not shown).

8

9 Using this procedure, we reconstructed each MAGIC genome into ~34
10 haplotype blocks on average with mean size 3.48Mb, representing
11 contributions from about 11 founder haplotypes (**Table S2**), and imputed the
12 full variant catalogue into each lines. Comparison of imputed SNPs with 782
13 GoldenGate SNP genotypes measured in 370 of the MAGIC lines²⁵ showed
14 98% concordance.

15

16 To map SVs, we divided the reference genome into 11,915 abutting source
17 loci, each 10kb wide, and computed six measures of anomalous read
18 mapping in each locus ($6 \times 11,915 = 71,490$ SV trait vectors) (**Methods, Table**
19 **1a, Table S3**). Four of these measures address different types of anomalous
20 read mapping that provide evidence of specific anomalies, namely high read
21 coverage for duplications, strandedness of reads for inversions, anomalously
22 large insert size for translocations and unpaired reads for deletions. The
23 remaining two measures are linear combinations of other measures that could
24 co-exist.

(a)

trait type	SV-QTLs	Unique	<i>cis</i>	<i>trans</i>
IP	1997	833	1617	380
ER	184	165	112	72
LIS	2051	585	1677	374
SS	1950	1887	1358	592
U	2060	1998	1530	530
U+LIS	2033	431	1661	372
Total	10275	5899	7955	2320

(b)

SV type	SV-QTLs	<i>cis</i>	<i>trans</i>
duplication	175	109	66
indel	3035	3035	0
inversion	1976	1373	603
other	1316	381	935
Total	6502	4898	1604

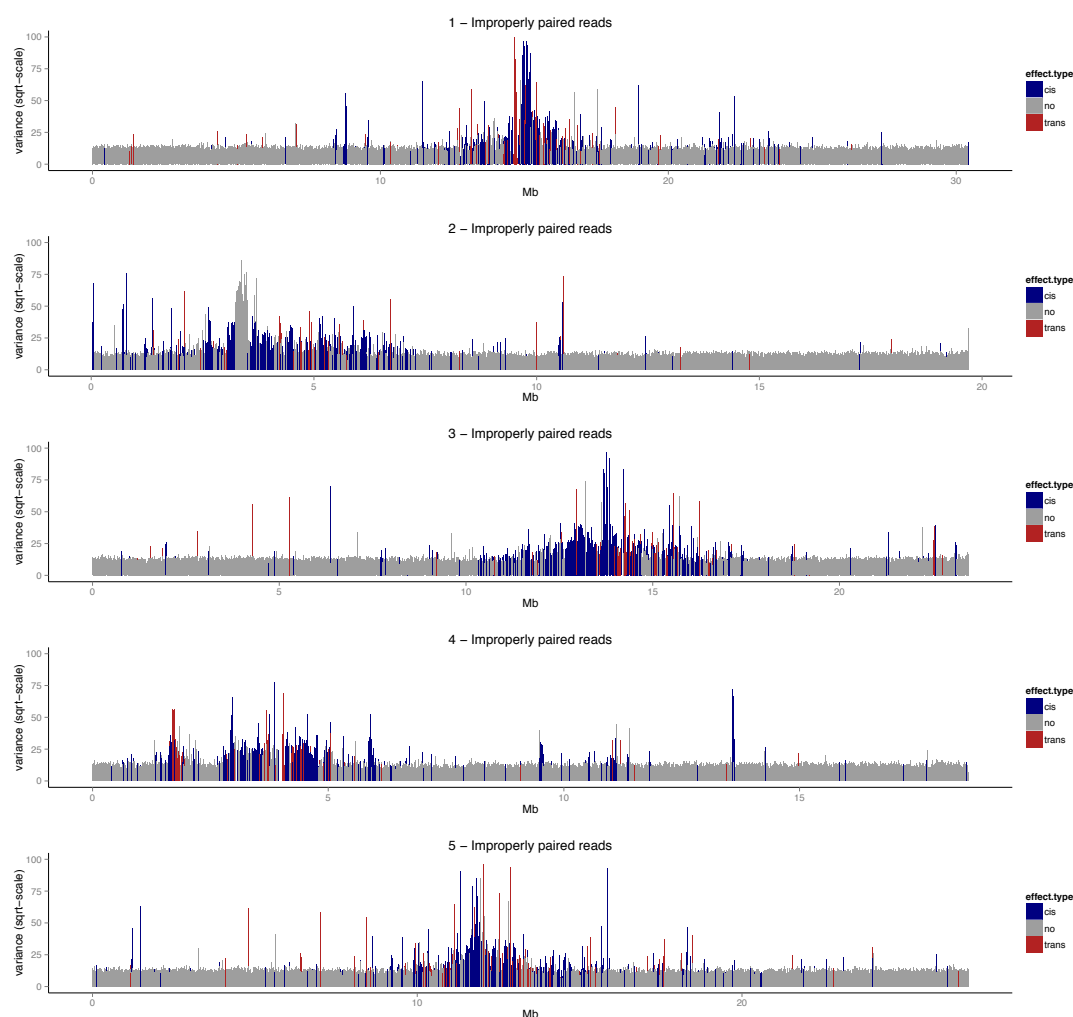
Table 1 (a) MAGIC SV-QTLs classified by read pair anomaly type. SV-QTLs: total number of QTLs detected using each anomaly type (if the same QTL was detected by multiple anomalies then it is counted multiple times in this column), Unique: number of QTLs detected only by a single anomaly category, *cis*: number of *cis* SV-QTLs, *trans*: number of *trans* SV-QTLs. **(b)** MAGIC SV-QTLs classified by QTL type, after removing duplicates. SV-QTLs: number of structural variants of each type, *cis*: number of *cis* SV-QTLs (source and sink within 2Mb from each other), *trans*: number of *trans* SV-QTLs [Note that the total number of SV-QTLs is 10,275, of which 6,502 are distinct after removing overlapping events, and 5,899 unique to a single anomaly type.]

12

1 Genetic association between each of the SV-trait vectors and the local
2 haplotype space was determined using a one-way ANOVA. We chose to
3 determine association at the level of haplotypes rather than SNPs for two
4 reasons. First, the founder haplotype space in the MAGIC lines is well-
5 defined, and measuring association with haplotypes can capture relationships
6 invisible at the level of SNPs. Second, the set of haplotype tests - defined by
7 the union of all the breakpoints, comprising 16,700 haplotype blocks, such
8 that the ancestral haplotype of all lines is unchanged within each block –
9 means about 75 times fewer tests are performed, thereby speeding up the
10 procedure (**Methods**). To determine genome-wide significance thresholds for
11 SV-QTLs we performed 100 phenotype permutations for each trait and then
12 fitted extreme value distributions (evd) to the genome-wide maxima of the
13 permutations (**Methods**). We merged together probable duplicate overlapping
14 SV-QTLs identified by multiple anomaly types.

15

16 After removing duplicates we identified 6,502 SV-QTLs at 1% study-wide false
17 discovery rate (evd $P < 0.001$) (**Table S3**). Of these, 1,604 (25%) were *trans*,
18 defined as mapping over 2Mb from the source. Overall, 4,073/11,915 (34.2%)
19 source loci harboured structural variants. Whilst we have greater power to
20 detect larger SVs, 2,379 overlapped annotated indels shorter than 2kb⁸.



1

2 **Figure 2** Genome-wide distribution of the variance for the trait “improperly-
3 paired reads” (the number of reads mapping to a locus with mapping
4 anomalies), computed in 10-kb windows. The x-axis shows genomic position
5 and the y-axis the variance of each trait vector scaled by its mean. Each
6 vertical line corresponds to a window. Those with SV-QTLs are blue (*cis*) and
7 red (*trans*). Centromeres are marked by pink bars.

8

9 The likelihood that a structural trait vector has an SV-QTL increases with its
10 variance (**Figure 2**). SV-QTLs are enriched around centromeres, as expected.

11 Away from the centromeres, **Figure 2** also shows that bins with variable SV
12 traits are isolated, rather than in clusters. **Figure 3a** shows the genome-wide

1 distribution of SVs segregating in one MAGIC founder, Ler-0. **Figure 3** and
2 **Table S3** show trans SV-QTLs link all five chromosomes.

3

4 In 319 SVs we were able to pinpoint both breakpoints, using contigs from *de-*
5 *novo* assemblies of the 19 founder genomes⁸ (see validation section below).

6 Mean SV size was 53kb in these SVs, and the largest was 189kb. Thus the
7 many of the SVs we discovered are too large to be due to insertions of small
8 transposable elements. This probably reflects our lack of power to detect very
9 small events, but also emphasizes that not all SVs are driven by mobile
10 elements.

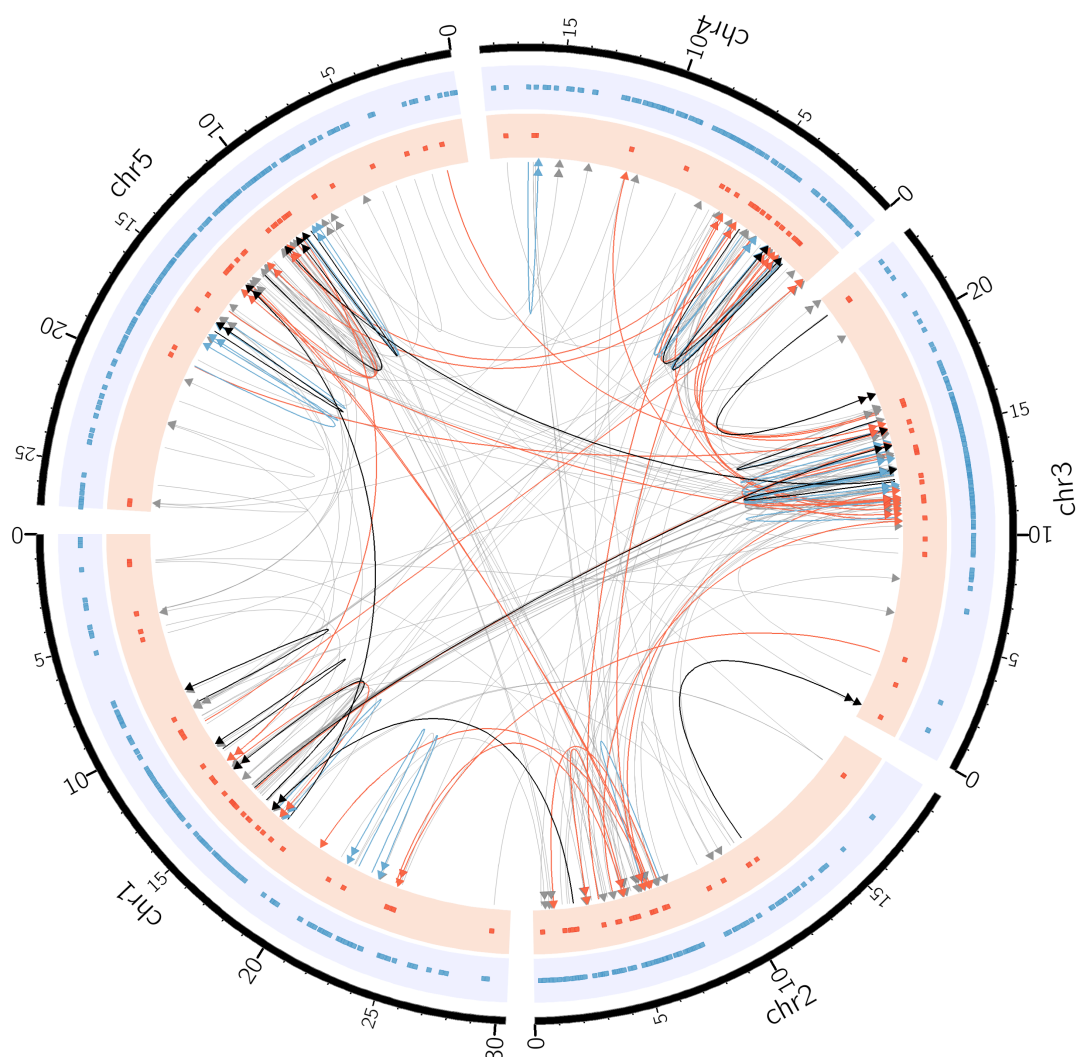
11

12 **Validation**

13

14 Genome-wide confirmation of SVs using short-read sequence is challenging
15 because SV breakpoints often associate with transposons and repeats that
16 hinder read-mapping and reassembly. However, among our SV-QTLs are
17 several known rearrangements. These include trans SV-QTLs linking a cluster
18 of rDNA repeats at ~14.2Mb on chromosome three to clusters at the ends of
19 chromosome two. Polymorphisms in these clusters are implicated in massive
20 genome size variations among *Arabidopsis* accessions²⁸. We also identified
21 the known knob inversion on chromosome 4 as reciprocal transpositions
22 linking 1.61Mb and 2.65Mb²⁹, and a 93kb inverted transposition identified
23 previously in a cross between Ler-0 and Col-0³⁰, and found it was present in
24 12 MAGIC founders.

1

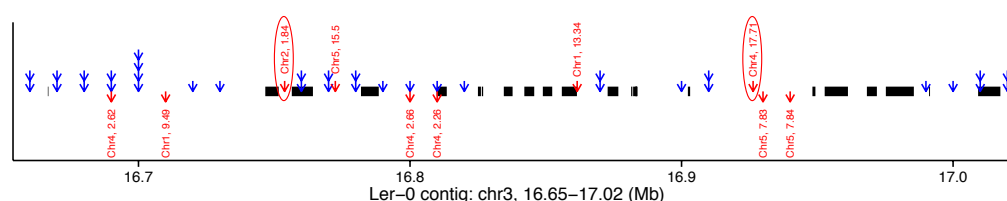


2

3 **Figure 3** Structural variants segregating in the accession Ler-0. The grey
4 directed lines show SV-QTLs with the arrows pointing towards the sink locus.
5 Red and blue links indicate 37 *trans* and 30 *cis* SV-QTLs confirmed by *de*
6 *novo* contigs. The black links show 16 SVs confirmed by PCR (7 *cis*, 9 *trans*).
7 Double arrows in links indicate inversions. The dots in the red and blue tracks
8 mark the sources (*trans* and *cis*, respectively) of all SVs associated with the
9 Ler-0 haplotype.

10 To validate further SVs we compared our SV calls for the founder accession
11 Ler-0 against two Ler-0 contigs (chr3:16.65-17.02Mb, chr5:25.06-25.23Mb)
12 that were independently re-sequenced and manually reassembled³¹, thereby
13 constituting a gold standard for comparison. The chromosome 3 contig

(**Figure 4**) is enriched in SVs (83 indels, 31 larger than 100bp), consistent with our analysis: 42 SV-QTL sources (36 *cis* and 6 *trans*) are in this region and 4 *trans* SV-QTLs map into it. As would be expected, the sources of these SV-QTLs are within gaps in the contig. Furthermore, alignment revealed two long-range SVs within the contig (a transposition and a duplication which align to chromosomes four and two, respectively), which coincide with the source and sink of two *trans* SV-QTLs mapped within the contig. Similarly, in the chromosome 5 contig, 6 *cis* SV-QTLs correspond to deletions (**Figure S3**).



10

Figure 4 Alignment of a manually assembled contig from Ler-0, chr3:16.65-17.02Mb to the reference genome annotated with SV-QTLs. Thick black lines show alignments to reference genome. Blue arrows show the sources of *cis* SV-QTLs; stacked arrows mean multiple read anomaly traits had SV-QTLs. Red arrows display *trans* QTLs with arrows starting from the source and pointing towards the sink. Gaps in the contig alignment indicate loci where Ler-0 did not align to the reference, with the exception of two transposed segments that mapped to chromosomes 2 and 4 at positions concordant with the sources of two *trans* SV-QTLs (circled).

20

We also used an independent *de-novo* assembly of Ler-0 built from long PAC-BIO reads, Genbank accession GCA_000835945.1³² to validate our *trans* SV predictions. This assembly was constructed algorithmically without manual revisions, and so is not guaranteed to be correct. Further, the Ler-0 individual

1 sequenced in the PAC-BIO assembly was different from the individual that
 2 founded the MAGIC population and therefore might carry private structural
 3 variations. Nonetheless, we expect it to be more accurate and contiguous
 4 than a Ler-0 assembly built from short Illumina reads alone. We took those
 5 3080 Illumina paired-end reads for Ler-0 from ⁸ that carried large insert size
 6 mapping anomalies when mapped to TAIR10 and that mapped to the sources
 7 of our predicted Ler-0 *trans* SV-QTLS, and then mapped them to the PAC-BIO
 8 assembly using bwa²⁶. These Illumina reads are from an individual grown
 9 from the same batch of seeds used to found the MAGIC population in ~2007,
 10 and should therefore share the same structural variants. Read anomalies that
 11 gave rise to correct SV predictions should map contiguously to the PAC-BIO
 12 assembly, under the assumption that the latter assembly is a more accurate
 13 representation of the Ler-0 genome. We found 2460 (80%) of these formerly
 14 split Illumina read pairs now mapped contiguously, defined as both members
 15 of a read-pair mapping to the PAC-BIO assembly with an insert size below
 16 600 bp.

17

18 With the exception of these manually assembled Ler-0 contigs and the
 19 provisional Ler-0 PAC-BIO assembly, the MAGIC founders are not
 20 contiguously reassembled into a genome-wide gold standard reference panel.
 21 Nevertheless, they provide some information to test our SV predictions. To do
 22 this, at each SV-QTL we predicted which founder haplotypes carried SVs at
 23 the origination of the population. Using the low coverage data for the 488
 24 MAGIC lines, at each SV-QTL we predicted which group of founders carried

1 the SV allele vs the reference allele based on correspondence between their
 2 SV-trait value and predicted founder allele, using the fact that SV haplotypes
 3 will have elevated anomalous reads at the source. We were able to do this
 4 confidently at 2,391 SVs where the founders divided into two groups, the
 5 remainder having complex multi-allelic SV predictions (**Methods**). We then
 6 examined the independently-collected high-coverage reads in each the 19
 7 MAGIC founders⁸ for read-mapping signatures that supported the predicted
 8 grouping of founders at each SV. We counted the read pairs linking source
 9 and sink at each of the 2,391 SVs in the 19 high coverage founders. At
 10 1,585/2,391 (66.3%, FDR 7.5%) SVs we observed significant differences in
 11 anomalies between the predicted groupings of founders (**Figure S4**, which
 12 also shows that the majority of SVs were mapped within 50kb). In the
 13 founders, the mean SV allele frequency was 6/19=31%. Only 387 (12%) were
 14 private to a single founder (**Figure S2**), in contrast to the fraction of SNPs
 15 (45%) that are private to a single founder⁸.

16

17 A related analysis using low-coverage reads from the 488 MAGIC genomes,
 18 but independent of founder predictions, (**Methods**) supported 1228/2391
 19 (51.3%, most also supported by the founder genomes) and 1631/4111
 20 (39.7%) of those remaining SVs without founder predictions. In total
 21 2,965/6,502 (45.6%) SVs were supported by either method.

22 ***Breakpoint Prediction and Confirmation***

23

1 In order to estimate SV sizes and identify SV breakpoints that could be tested
 2 experimentally by PCR, we next assembled the high-coverage sequence data
 3 for the MAGIC founders into *de-novo* contigs. No scaffolding was attempted in
 4 order to produce conservative high-quality short contigs, each up to a few
 5 kilobases long. We aligned these contigs to the reference to find alignments
 6 split between sources and sinks. We mapped breakpoints for 420 SV-QTLs
 7 (**Methods, Table S4**): in 319 SV-QTLs both breakpoints were identified. We
 8 designed PCR primers around 77 breakpoints from 45 SVs (both breakpoints
 9 in 7 SVs, and one in each of the remaining 38). We validated 37 SVs (83%),
 10 comprising 61 (79%) breakpoints, in 14 *cis* (6 inversions, 7 transpositions, 1
 11 indel) and 23 *trans* (23 transpositions, 13 with inversions) SV-QTLs (**Table**
 12 **S5**).

13

14 Consistent with our difficulties in predicting biallelic founder alleles, in 11 SVs
 15 the breakpoints were polymorphic among the founders carrying the SV, and in
 16 5 transpositions the orientation of the SV differed between founders. These
 17 results emphasise the difficulties frequently reported when delineating SVs.

18

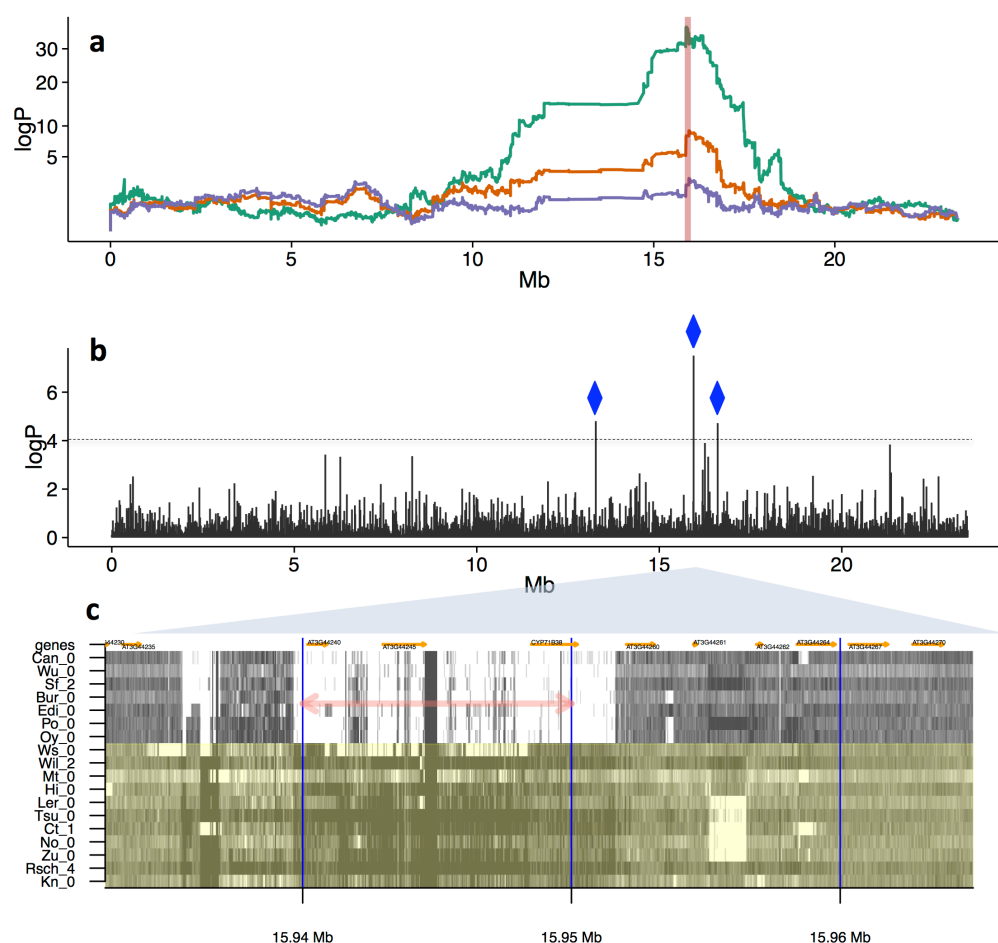


Figure 5 Association of haplotypes and SVs **(a)** Genome scans over chromosome 3 (x-axis: genomic position, y-axis: logP of association). Orange: association of local haplotype with germination time (days), peaking at 15.93Mb. Green: association of local haplotype with the SV trait unpaired reads at the source locus 15.94-15.95Mb (indicated by the vertical red line), explaining 8.13% of the variance in germination time, with an SV-QTL mapped at the same position as the germination QTL. Purple: residuals of germination time after regressing out the SV trait, ablating the QTL. **(b)** Chromosome-wide Pearson correlations between germination time and the numbers of unpaired reads measured at each 10kb source locus (x-axis: genomic position, y-axis: $-\log_{10}$ P-value of test that the correlation is zero). Three source loci correlate strongly with germination ($\log P > 4$), all with *cis* SV-QTLs (blue diamonds). **(c)** Structural variation in the MAGIC founders. Shown is the read coverage in 18 accessions (labelled on y-axis), over ~30kb surrounding around 15.94Mb (x-axis). Dark shades indicate high coverage, light shades low coverage. The 10kb intervals used to define source loci are delineated by vertical blue lines. The source locus giving rise to the SV-QTL in (a), (b) is marked with a pink double-arrow. Those founder accessions predicted to carry the reference allele (No-0, Ct-1, Mt-0, Wil-2, Ler-0, Tsu-0, Rsch-4, Kn-0, Zu-0, Hi-0, Ws-0) are in green, those predicted to carry the SV are in grey. Genes are annotated in orange.

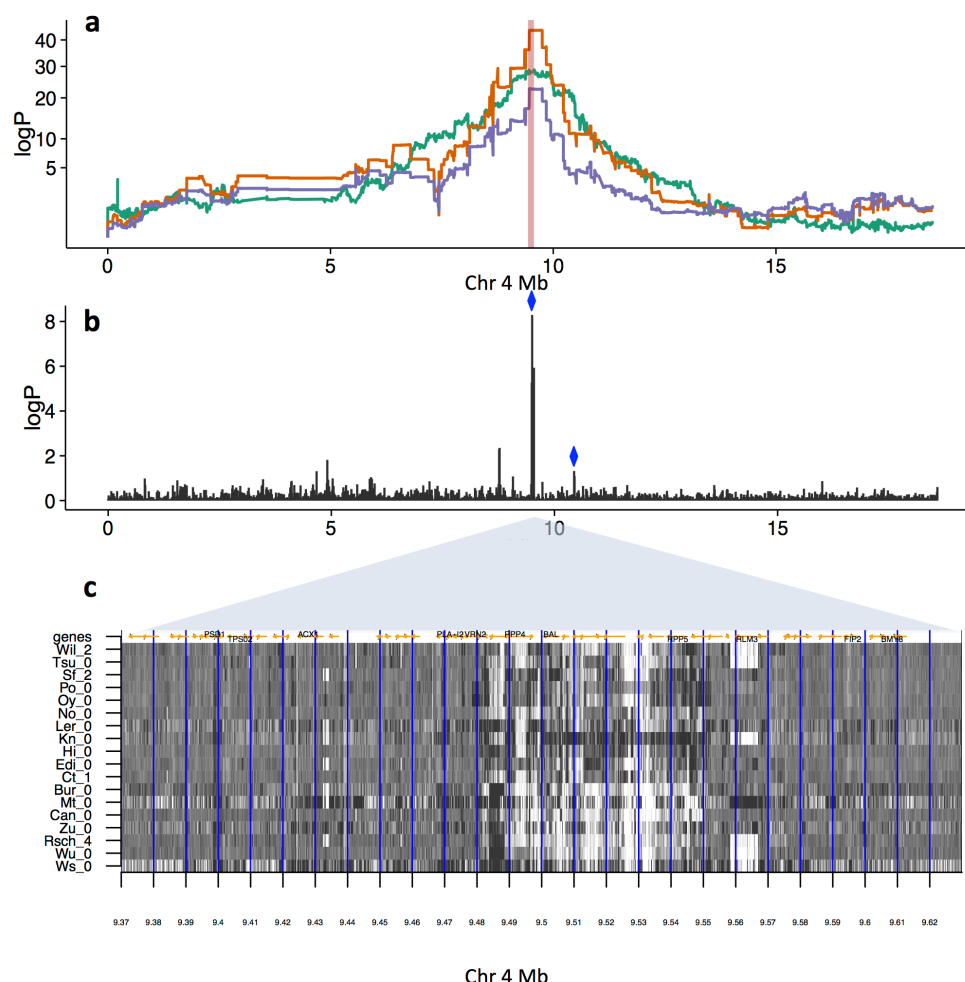


Figure 6 Effects of SVs on resistance to *Albugo laibachii* infection, (a) Genome scans on chromosome 4. Orange: Association with resistance. The peak of association for is at 9.50Mb. Green: Association with SV-trait improperly paired reads at source 9.50-9.51. Purple: Resistance after two SV traits have been regressed out measuring improperly paired reads (sources chr4, 9.50-9.51Mb (green line) and chr4, 10.44-10.45Mb (not shown), both marked with blue diamonds in fig (b)) that together explain 24.7% of the phenotypic variance. (b) logP of association between SV traits for improperly-paired reads and the resistance trait. There is a cluster of associated traits near 9.50Mb, in addition to the more weakly associated trait at 10.44-10.45Mb. (c) Structural variation in high-coverage sequence in the MAGIC founders around 9.50Mb. Shown is the number of improperly-paired reads (dark: high values, light: low values) in 18 accessions (labeled on y-axis), between 9.37-9.63Mb (x-axis). The 10kb intervals used to define source loci are delineated by vertical blue lines. There is a region of complex structural variation spanning 9.48-9.55Mb approximately, with considerable variation between the founder accessions. Genes are marked by orange arrows, and selected genes, some implicated in disease resistance at this locus, are labelled.

1 ***Effects of SVs on phenotypic QTLs and gene expression***

2 We next investigated associations between SVs and 9 physiological
3 phenotypes, either previously published^{25,33} or new to this study (**Table S6**).

4 We found 16 distinct SV-QTLs (8 in *trans*, **Table S7**) that overlap
5 physiological QTLs. In some cases, regressing the SV-trait from the
6 physiological trait ablated the physiological QTL, consistent with, albeit not
7 proof, that the SV is causal. This is illustrated by a QTL for germination time²⁵
8 on chromosome 3, which is ablated by a *cis* SV-QTL for unpaired reads at
9 around 15,936,650-15,951,640bp (**Figure 5a,b**). Our analysis predicted that 7
10 founders would carry a deletion at this locus, which was confirmed by the
11 independent founder sequences (**Figure 5c**), revealing a 15kb deletion of
12 three genes, *AT3G44240* (Polynucleotidyl transferase, ribonuclease H-like
13 superfamily protein), *AT3G44245* (pseudogene of cytochrome P450, family
14 71, subfamily B, polypeptide 21), and *CYP71B38* (*AT3G44250*, cytochrome
15 P450, family 71, subfamily B, polypeptide 38). Other SVs segregate nearby,
16 but with allelic patterns inconsistent with the trait and therefore unlikely to be
17 causal. It is therefore probable that the causal variant(s) lies within the deleted
18 region. The three genes are not known to affect germination, although a
19 mutant of another Polynucleotidyl transferase, *AHG2* (*AT1G55870*) does³⁴.

20

21 We found similar effects on the chromosome 4 QTL for resistance to the
22 fungal pathogen *Albugo laibachii*, isolate Nc14³⁵ (**Figure 6, Table S7**).

23 Variation in the number of unpaired reads at 9.50-9.51Mb explains 18.3% of
24 the variance in resistance, and is adjacent to a cluster of Leucine-rich repeat

1 genes, and the genes RPP4³⁶, BAL³⁷ and RPP5. This locus is rearranged in
2 some Arabidopsis accessions and known to be involved in disease
3 resistance³⁷; **Figure 6** confirms the founder genomes have complex,
4 polymorphic SVs in this region. Since the resistance QTL is not completely
5 ablated by the SV traits associated with it, additional non-structural variants
6 likely contribute to it.

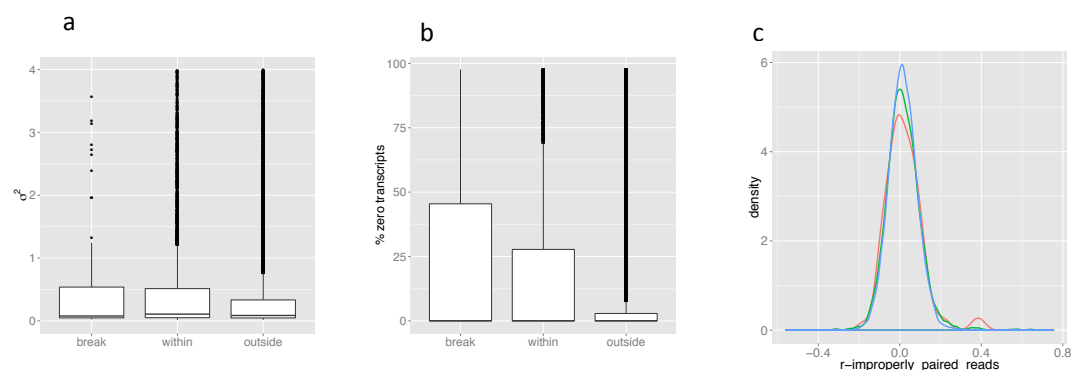
7

8 Importantly, **Figures 5b,6b** show that correlations between SV traits and
9 phenotypes are tightly localized, generally to width of a single SV trait window,
10 in contrast with wider linkage disequilibrium decay seen in QTL genetic
11 mapping (**Figure 5a**). Consequently correlations between SV traits and
12 physiological traits can sometimes pinpoint causal variants within
13 physiological QTLs which are otherwise too broad to localize (mapping
14 resolution in the MAGIC population is about 200kb²⁵).

15

16 We also corroborated studies^{21,24} showing SVs associate with gene
17 dysregulation, even when the gene sequence is undisturbed. Within those
18 SVs with mapped breakpoints, 119 genes spanned the breakpoints, 6,909 lay
19 inside the SVs (**Table S8**) and 21,747 outside. Using RNA-seq from 200
20 MAGIC aerial seedlings, scaled expression variance increased among genes
21 spanning breakpoints (t-test: $P < 9 \times 10^{-3}$) and within SVs ($P < 1 \times 10^{-13}$)
22 (**Figure 7a**). Similarly larger fractions of lines had silenced transcripts for
23 genes spanning breakpoints (t-test $P < 1.2 \times 10^{-2}$) and within SVs ($P <$

- 1 2×10^{-52}) (**Figure 7b**). Expression within SVs was more correlated with local
- 2 SV traits than outside SVs (F-test $P < 2.1 \times 10^{-6}$) (**Figure 7c**).



3

- 4 **Figure 7** Variation of expression in 200 MAGIC leaf transcriptomes, in genes
- 5 spanning SV breakpoints, within SVs or outside SVs. **(a)** Boxplots of transcript
- 6 variance (scaled by the mean). **(b)** Boxplots of the fractions of silenced genes
- 7 **(c)** Distributions of the Pearson correlations between gene expression and
- 8 number of abnormal insert size reads in the locus containing the gene (red:
- 9 spanning breakpoints, green: within SVs, blue: outside SVs).

10

Phenotype	h_H^2	h_{SNP}^2	h_{SV}^2				
			IP	LIS	SS	U	U+LIS
Resistance (resistance to <i>Albugo laibachii</i>)	0.000 (0.139)	0.258 (0.085)	0.490 (0.335)	0.511 (0.307)	0.000 (NA)	0.673 (0.504)	0.503 (0.314)
RosetteLeafNumber.LongDay (number of leaves in a rosette for plants grown under long daylight)	0.228 (0.081)	0.322 (0.076)	0.463 (0.148)	0.456 (0.146)	1.000 (NA)	1.000 (0.377)	0.447 (0.146)
RosetteLeafNumber.ShortDay (number of leaves in a rosette for plants grown under short daylight)	0.038 (0.060)	0.047 (0.062)	0.000 (NA)	0.000 (NA)	0.000 (NA)	0.000 (NA)	0.000 (NA)
bolting.Bath (bolting time in a greenhouse)	0.426 (0.064)	0.476 (0.048)	0.783 (0.093)	0.783 (0.093)	0.952 (0.047)	0.989 (0.025)	0.785 (0.092)
days.to.germ.x (germination time)	0.220 (0.068)	0.149 (0.063)	0.385 (0.116)	0.357 (0.113)	0.598 (0.165)	0.835 (0.146)	0.365 (0.114)
fieldFT.pl (flowering time in the field)	0.000 (0.068)	0.095 (0.076)	0.000 (0.179)	0.000 (0.130)	0.000 (0.913)	0.000 (NA)	0.000 (0.145)
fieldRD.pl (rosette diameter plasticity)	0.000 (NA)	0.000 (0.063)	0.000 (0.085)	0.000 (0.084)	0.000 (0.239)	0.166 (0.220)	0.000 (0.085)
leaves.day.28.given.days.to.germ (residuals for number of leaves at day28 regressed on germination)	0.193 (0.081)	0.299 (0.066)	0.391 (0.146)	0.362 (0.140)	0.836 (0.189)	0.675 (0.272)	0.366 (0.142)
tll_branch.BATH (total number of branches of plants)	0.106 (0.048)	0.196 (0.054)	0.276 (0.104)	0.275 (0.100)	0.419 (0.193)	0.616 (0.214)	0.279 (0.102)

1

2 **Table 2 Estimates of heritability.** h_H^2 is haplotype-based heritability. h_{SNP}^2 is
3 SNP-based heritability. h_{SV}^2 is the heritability estimated from structural variant
4 anomaly traits. Numbers in brackets are the standard errors of the heritability
5 estimates above. ER: Excess Reads, IP: Improperly-paired, LIS: Large Insert
6 Size, SS: Same Strand, U: Unpaired, U+LIS: Unpaired or Large Insert Size.
7 Heritabilities for excess reads are not reported because the fraction of bins in
8 any individual containing non-zero entries was too small.

9

10 Effects of SV-traits on Heritability

11

12 Finally, we treated the SV traits as if they were quantitative noisy genotypes to
13 define pair-wise correlations between MAGIC lines, as weighted correlations
14 of their SV traits (**Methods**). We constructed SV genetic relationship matrices
15 (GRMs) K_{SV} , which we used to compute the SV-heritability h_{SV}^2 of each of the
16 physiological traits mapped above by analogy with the mixed models used for

1 estimating SNP-based heritability³⁸. This idea is similar to the use of gene
 2 expression data to model intersample relationships³⁹. We also compared
 3 these SV-heritabilities with those obtained from “classical” haplotype K_H or
 4 SNP-based K_{SNP} GRMs (**Table 2**). K_H was computed from the identity
 5 between haplotype mosaics (and so measures identity by descent), K_{SNP} and
 6 K_{SV} were computed from the correlations of 1.2M imputed SNPs or 12k SV-
 7 traits respectively (Methods). We also computed SV heritability when only the
 8 most variable 50% or 25% of SV-traits were included, to investigate if
 9 heritability was concentrated at the most structurally variable loci.

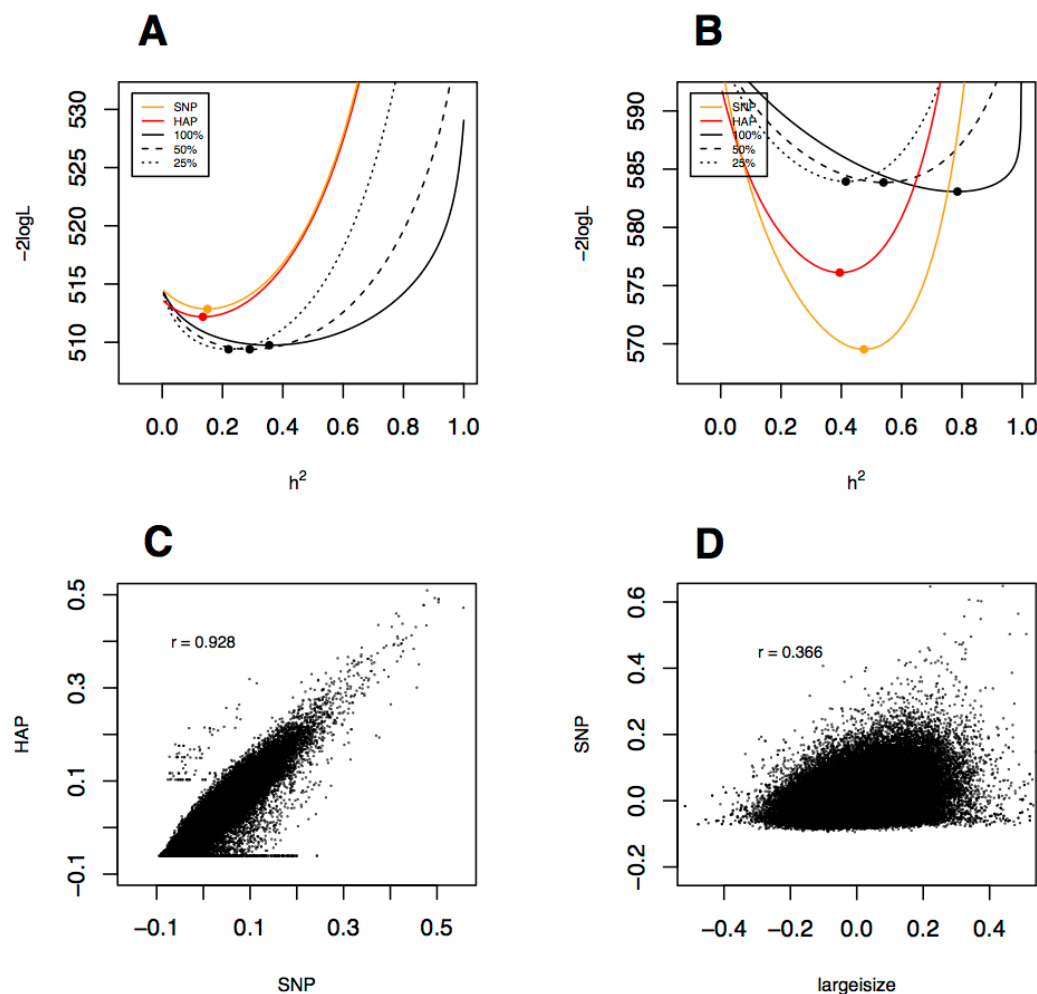
10

11 As expected, SNP-based heritability h_{SNP}^2 is generally similar to haplotype-
 12 based heritability h_H^2 for all phenotypes tested. However, the heritability h_{SV}^2
 13 captured by the six measures of SV anomaly is more variable, sometimes
 14 being close to zero, but sometimes exceeding classical heritability by a
 15 considerable margin (**Table 2**). The standard error of h_{SV}^2 was typically about
 16 twice that of h_{SNP}^2 or h_H^2 , (approximately 0.1 compared to 0.05), presumably
 17 reflecting greater uncertainty in SV-traits than in SNPs or haplotypes.
 18 Therefore the larger heritability estimates should be treated with caution.
 19 Nonetheless, for phenotypes such as times for germination or bolting, the
 20 standard errors of all estimates are comparable at ~0.05 and is possible to
 21 make meaningful comparisons. **Figure 8A,B** illustrates likelihood curves the
 22 times to germination (A) bolting (B), for SNP, haplotype and large insert-size
 23 anomalies. Visualising the entire curves gives a better sense of the
 24 uncertainty of the maximum likelihood estimates at the curves’ minima (the

1 standard errors in **Table 2** are asymptotic estimates based on the curvature at
 2 these minima). The Figure 8B shows that for bolting time the heritability
 3 attributable to all large-size SV-traits, $h^2_{large\ size}$, is close to 80%, compared to
 4 40-50% for haplotype or SNP-based estimates. As the fraction of SV traits is
 5 reduced by progressively removing those traits with lower variance, $h^2_{large\ size}$
 6 reduces to that of SNPs or haplotypes. This suggests that there is genome-
 7 wide structural variation that is not tagged by standard genetic variation, and
 8 which has important effects on specific phenotypes. These effects are not
 9 universal, as Figure 8A shows for germination time, where heritability is
 10 similar for all estimates.

11

12 The relative independence of the heritability estimates borne out by low
 13 correlations between the corresponding elements of SNP and SV-based
 14 GRMs, which range around 0.3 depending on the anomaly type (Figure 8D
 15 shows the relationship between GRMs computed from SNPs vs large
 16 insertsize anomalies), compared to the correlation of 0.93 between SNP and
 17 haplotype based GRMs (Figure 8C).



1

2 **Figure 8** (a,b) log-likelihood curves for two phenotypes bolting.BATH (large
3 insert size read anomalies) and RosetteLeafNumber.ShortDay (unmapped
4 read anomalies), illustrating contrasting behavior of heritability estimates
5 based on structural variants, SNPs and haplotypes. Log-likelihood curves as
6 functions of heritability are plotted for the GRMs estimated from SNPs,
7 haplotypes and various fractions of anomalies. The maximum likelihood
8 estimates of each heritability measure correspond to the minima of the
9 corresponding curves, and are marked with dots. (c,d) Scatter plots
10 comparing the off-diagonal elements of genetic relationship matrices. (c) K_{SNP}
11 vs K_H ; (d) K_{SNP} vs $K_{largeize}$.

12

1 Discussion

2

3 Our aim has been to understand better the architecture and impact of
4 structural variation in populations sequenced at low coverage. We used a
5 strategy that combines analysis of read-mapping signatures commonly used
6 to detect SVs in individuals sequenced at high coverage, with association
7 mapping in populations⁴⁰. A somewhat related concept was used for mapping
8 un-localised contigs into reference assemblies based on linkage
9 disequilibrium⁴¹.

10

11 In doing so, we have generated a partial catalog of SVs in Arabidopsis,
12 although our purpose is not to call SVs systematically, a task that remains
13 challenging with short reads. Rather, we have shown how SVs' impact can be
14 assayed without necessarily calling them or mapping their breakpoints. In this
15 way, we can distinguish transpositions from local SVs, and determine the
16 approximate locations of transpositions. The privileged role of the reference
17 genome in the analysis means that some transpositions appear as deletions,
18 so we probably have underestimated their true frequency. Despite this, a
19 quarter of the SVs we detected are transpositions. Given the large numbers of
20 transposable elements in Arabidopsis - over 11,000 from over 300 families
21 are annotated in the reference²⁴ – this is unsurprising. However, many of the
22 SVs we mapped are too large, covering tens of kilobases, to be single
23 transposon-mediated events.

24

1 In the minority of cases where we delineated breakpoints exactly, we often
2 found SVs are complex combinations of different SV types. But often
3 breakpoints were not simple cut-and-paste transformations of the reference
4 genome, as illustrated in Figure 6c. Indeed, it is impossible to determine
5 precisely the changes that led to many observed structural variants.

6

7 Because we used ultra-low-coverage 0.3x sequence data, we divided the
8 Arabidopsis genome into 10kb bins when counting read-mapping anomalies.
9 With higher coverage and a larger sample size it would be possible to use a
10 larger number of narrower bins, thereby improving resolution. The public
11 release of over 3000 rice genomes sequenced at $\sim 14\times$ ⁴² and over 1000
12 Arabidopsis accessions sequenced at over $\sim 20\times$ ⁴³ means that there are now
13 large collections of inbred plant genomes available for analysis. Both of these
14 sets are worldwide surveys of germplasm, in which we expect SVs to
15 contribute significantly to, and be confounded with, their extensive population
16 structure, in contrast to the MAGIC population used here. Disentangling these
17 effects will be a challenging but important task.

18

19 Mapping SVs in a population brings new insights to the problem of QTL
20 analysis. First, an SV trait inside a QTL may entirely explain the genetic effect
21 at the QTL, and hence provide support for being the causal variant (e.g.
22 **Figure 6**). Second, SV traits are much more tightly localized than are QTLs:
23 there is little or no correlation between neighbouring SV traits so there are no
24 effects of linkage disequilibrium. Our analysis also shows that expression of

1 genes is often dysregulated or even silenced within large SVs, raising the
2 possibility that an SV causes multiple regulatory and phenotypic effects.

3

4 Finally we have shown that even in a population like Arabidopsis MAGIC
5 where the local haplotype space is known, structural variation has an impact
6 on heritability that cannot be explained by standard genetic variation. This is
7 unexpected given the breeding history and genetic architecture of the MAGIC
8 lines. For if an SV segregated among the founders of the MAGIC lines, then it
9 should be tagged by the local haplotype context, and therefore contribute to
10 both h_H^2 and h_{SV}^2 .

11

12 One possible explanation is that structural variation at loci rich in mobile
13 elements accumulates within each lineage, leading to SVs that are private to
14 each MAGIC line but tend to occur at the same loci, thereby creating similar
15 phenotypic effects. Supporting this, in our analysis the SV-relationship matrix
16 is calculated empirically, without regard to the ancestry of the MAGIC lines,
17 being solely a function of the counts of read-mapping anomalies. Therefore,
18 recalling that the history of each MAGIC line includes a private lineage of at
19 least five generations of selfing, should SVs accumulate recurrently but
20 independently in different lineages, then these could generate phenotypic
21 associations invisible to SNP or haplotype variation. In Arabidopsis, it is
22 known that some mobile elements are methylated, often in response to
23 environmental cues, and that such methylation plays a role in the epigenetic
24 control of certain phenotypes⁴⁴. Testing this hypothesis in Arabidopsis MAGIC

1 lines would require complete and precise reassembly of each genome using
2 long reads, annotation of mobile elements and determination of their
3 methylation status.

4

5 The general role that recurrent, but independent, genomic rearrangements
6 might play in Arabidopsis and in other species remains to be seen, but there
7 is no *a priori* reason why it should not be a driver of phenotypic variation. The
8 approach used here may therefore have wider application to other
9 populations, both to characterize the extent of transpositions and the impact of
10 cryptic structural variation on phenotypes.

11

12

1 **Methods**

2

3 **DNA extraction and sequencing** MAGIC lines were grown at Bath (lab of P.
4 Kover) or Oxford (lab of N.P. Harberd) in greenhouses or growth chambers
5 respectively. Leaves were harvested for DNA extraction. DNA isolation was
6 performed at the John Innes Centre, in 96 well plates using the DNeasy 96
7 Plant Kit and DNeasy 96 Protocol (www.qiagen.com). Sequencing was
8 performed by the Oxford Genomics Centre.

9

10 **Genomic DNA library construction and multiplexing** Samples were
11 quantified using the Quant-iT™ PicoGreen® dsDNA Kits (Invitrogen) and a
12 Genios plate scanner (Tecan) according to manufacturer specifications.
13 Sample integrity was assessed using 1% agarose gel. Approximately 300ng
14 of DNA were fragmented using a Covaris S2 system with the following
15 settings: Intensity: 5, Duty Cycle: 20, Cycles per Burst: 200, Time: 60 sec.
16 Distribution of fragments after shearing was determined using a Tapestation
17 D1200 system (Agilent/Lab901). DNA Libraries were constructed using the
18 NEBNext DNA Sample Prep Master Mix Set 1 Kit (NEB) with minor
19 modifications and a custom automated protocol on a Biomek FX (Beckman).
20 Ligation of adapters was performed using Illumina Adapters (Multiplexing
21 Sample Preparation Oligonucleotide Kit). Ligated libraries were size selected
22 using Ampure magnetic beads (Agencourt). Each library was PCR enriched
23 with 25 μ M each of the following custom primers:

24 Multiplex PCR primer 1.0

1 5'-
2 AATGATACGGCGACCAACGAGATCTACACTCTTTCCCTACACGACGCTCT
3 TCCGATCT-3'

4 Index primer

5 5'-
6 CAAGCAGAAGACGGCATACGAGAT[INDEX]CAGTGACTGGAGTTCAGACG
7 TGTGCTCTTCCGATCT-3'

8 Indexes used were 8bp long (manuscript in preparation). Enrichment and
9 adapter extension of each preparation was obtained using 5µl of size selected
10 library in a 50 µl PCR reaction. After 10 cycles of amplification (cycling
11 conditions as per Illumina recommendations) the reactions were purified with
12 Ampure beads (Agencourt/Beckman). The final size distribution was
13 determined using a Tapestation 1DK system (Agilent/Lab901). The
14 concentrations used to generate the multiplex pool were determined by
15 Picogreen. The library resulting from the pooling was quantified using the
16 Agilent qPCR Library Quantification Kit and a MX3005P instrument (Agilent)
17 before sequencing on an Illumina GAIIx as 50bp or 100bp paired end reads.
18 All steps for library construction, including the setup of the PCR reaction were
19 performed on a Biomek FX (Beckman). Post PCR cleanup was carried out on
20 a Biomek NXp (Beckman) whereas a Biomek 3000 (Beckman) was used to
21 generate the pools of 96 indexed libraries.

22

23 **Processing Sequence Reads and SNP Calling** The Illumina reads were
24 mapped to the *A. thaliana* reference genome (TAIR10) using Stampy version

1 v1.0.20²⁵. Alignments were stored in a separate BAM file for each MAGIC
 2 line. Previous sequencing for the 18 MAGIC line progenitors had produced a
 3 catalogue of 3,316,270 segregating SNPs¹⁷. We ran GATK v2.6²⁶ on the
 4 segregating SNPs to call variants for the 19 founders, setting the following
 5 read filters: Allele Balance, BaseQualityRankSumTest, Clipping
 6 RankSumTest, Coverage, DepthPerAlleleBySample, FisherStrand,
 7 GCContent, HaplotypeScore, LowMQ, MappingQualityRankSumTest,
 8 MappingQualityZero, MappingQualityZeroBySample, RMSMappingQuality,
 9 ReadPosRankSumTest. We filtered out SNPs that were triallelic, within
 10 transposons, or heterozygous for any founders.

11

12

13 **Definition of Structural Variant Traits** We divided the TAIR10 (Col-0)
 14 reference genome into 11,915 abutting 10 kb segments. Within each segment
 15 we computed six measures of anomalously mapped reads that are signatures
 16 of SVs. Let R be the set of all reads mapped to a genome of length L ; ρ is the
 17 number of reads in R , and ρ_l the number of reads mapped to a segment l of
 18 length 10kb. The read anomaly measures computed in each segment are:

19 1. **High read coverage:** $\rho_{hc} = \rho_l - 1.5E[\rho_l]$, where $E[\rho_l] = \frac{\rho \times l}{L}$ is the
 20 expected read coverage of the segment

21 2. **Unpaired reads:** ρ_u number of reads mapping to the segment whose
 22 pair is not mapped

23 3. **Pairs on the same strand:** ρ_s number of reads with pair on the same
 24 strand

1 **4. Reads with large insert size:** ρ_{is} number of read pairs with insert size
2 outside the range $m_s \pm IQR_s$ or mapped to different chromosomes,
3 m_s, IQR_s being the median and interquartile range of insert sizes of all the
4 reads in the sample.

5 **5. Unpaired reads or with large insert size:** $\rho_{ui} = \rho_u + \rho_i$

6 **6. Improperly paired reads:** $\rho_{uis} = \rho_u + \rho_i + \rho_s$

7 The last two traits are combination of others – certain SV types can cause
8 multiple anomaly signatures, so merging them may increase power. Each type
9 of read pair anomaly was measured in each of the 11,915 10kb segments,
10 determining 71,490 traits in total.

11 **Genome scan** We treated the SV traits like a gene expression eQTL study,
12 performing a genome scan for each one. Association was tested by fitting trait
13 vectors to the imputed ancestral haplotype at each locus in the 488 genome
14 mosaics. In combination, the mosaics partitioned the genome into 16,700
15 haplotype blocks, with the ancestral haplotype of all lines unchanged in each
16 one. Let y_{Ai} be the number of anomalous reads of a certain type at source
17 segment A in line i . At every haplotype block p we fit the linear model:

$$y_{Ai} = \mu_A + \sum_{s \in S} X_{pi}(s) \beta_{Ap}(s) + e_i$$

18 μ_A is the average trait value at A , $X_{pi}(s)$ is a binary indicator of whether line i
19 carries haplotype s at p , $\beta_{Ap}(s)$ is the effect of founder haplotype s and e_i the
20 (Normally distributed) error. The founder effects $\beta_{Ap}(s)$ were estimated by an
21 one-way ANOVA with null hypothesis:

$$H_0(p): \beta_{Ap}(s) = 0 \forall s$$

1 **Genome-wide significance** We denote the location of a sink locus
2 associated with an SV trait an SV-QTL. Each genome scan of a given SV trait
3 returns a p-value π_{Ap} for each of the 16,700 scanned blocks p . We selected
4 as candidate SV-QTLs for each mapped trait the locus with maximum
5 genome-wide negative logarithm of π_{Ap} , i.e.

$$\lambda_A = \max (-\log_{10}(\pi_{Ap}))$$

6 To control for the number of test in each scan (16,700) and correct for
7 associations driven by outliers we performed 100 permutations T_A of the trait
8 vector y_{Ai} and repeated the mapping for each one. We then fitted a
9 generalized extreme value distribution (GEV), using R evd package on the
10 $\lambda_A(t), t \in T_A$ values obtained by permutation, from which we obtained a
11 genomewide corrected p-value:

$$\gamma_A = -\log \left(1 - \exp \left(1 + \hat{s}_A \left(\frac{\lambda_A - \hat{\alpha}_A}{\hat{b}_A} \right)^{-\frac{1}{\hat{s}_A}} \right) \right)$$

12 $\hat{\alpha}_A, \hat{b}_A, \hat{s}_A$ are the MLEs of α_A, b_A, s_A . Study-wide SV-QTLs are selected at
13 $FDR < 10^{-2}$ ($\gamma_A < 10^{-3}$). At this FDR we mapped 10,275 SV-QTLs in total.
14 **Table S3** shows mapped QTLs per read anomaly category. 3,773 SV-QTLs
15 had coincident sources and sinks, probably corresponding to the same SV,
16 and so were counted only once. SVs are tabulated in **Table S4**.

17 **Cis and trans SV-QTLs** In total we mapped 6,502 distinct SV-QTLs, each

1 corresponding to a unique SV. SV-QTLs with sources and sinks within 2Mb
2 from each other were classified as *cis*, and the rest as *trans*.

3 **Prediction of SV allele frequency** We predicted the founder haplotypes
4 carrying SVs at the origination of the population, using the fact that SV
5 haplotypes will have elevated anomalous reads at the source. For each SV-
6 QTL the founders' contributions were arranged as a 19×19 table T whose
7 cells (i, j) carry the sum of read anomalies (of a certain type) at the source for
8 all lines carrying haplotype i at the sink and haplotype j at the source. A
9 founder is classified as carrying the SV if its corresponding row has generally
10 higher values than the rest of the table (we note that in *cis* QTLs the matrix is
11 almost diagonal).

12 For each cell of T , let t_{ij} be the sum of trait values for all genomes carrying
13 haplotypes i, j at the sink and source, respectively. The contributions of
14 founder i are estimated as the "row" effect: $r_i = \sum_j t_{ij}$, which are reorder such
15 that $r_1 > \dots > r_{19}$, with $r_1 \approx \dots \approx r_{19}$ under null. We rejected the null hypothesis
16 if there is a set $\{r_1, \dots, r_k\}$ such that $r_1 > \dots > r_k > r_{k+1} > \dots > r_{19}$. There are 18
17 such possible sets. The z-score of each set k is:

$$z_k = \frac{\sum_{j=1}^k r_j - E[r_k]}{\sigma(r_k)}$$

18 $E[r_k]$ and $\sigma(r_k)$ are estimated by 1000 permutations of T , denoted as R_{z_k} . We
19 choose k such that z_k is maximised. k is significant, hence the corresponding

20 founders carry the SV if the permutation p-value $\pi_{z_k} = \frac{|r_{z_k} \in R_{z_k} : r_{z_k} \geq z_k|}{N} \leq 10^{-2}$.

1 Examples of *cis* and *trans* SV-QTL tables with detectable founders are shown
2 in **Figure S1**. The test predicts founder haplotype groups at 2,391 SV-QTLs.

3 **Heritability**

4 For a given phenotype (such as germination time) measured in the MAGIC
5 lines, the phenotypic variance matrix is represented by the mixed model
6 $V = K\sigma_g^2 + I\sigma_e^2$ where K is the genetic relationship matrix (GRM) and I the
7 identity matrix. The phenotypic heritability is

8

$$h^2 = \sigma_g^2 / (\sigma_g^2 + \sigma_e^2)$$

9

10

11 We computed genetic relationship matrices K between the inbred MAGIC
12 lines in three ways:

13

14 **Identity By Descent (haplotype-based) K_H** We used the representations of
15 chromosomes as homozygous mosaics of the 19 founders to determine
16 identity by descent (IBD). Across all N MAGIC lines, we identified the union of
17 the mosaic breakpoints, and then segmented the genome of each MAGIC line
18 according to these breakpoints. Thus by construction, the founder haplotype
19 for each line is constant within each segment. The founder haplotype in
20 segment L in line i is represented by an indicator matrix H_{iLf} which is 1 if the
21 founder is f and 0 otherwise. If w_L is the fraction of the genome covered by L ,

1 and $f_{ijL} = \sum_f H_{iLf} H_{jLf}$ is the indicator of whether lines i, j are IBD at L , then

2 the fraction of the genome that is IBD for lines i, j is

3

$$d_{ij} = \sum_L w_L f_{ijL}$$

4

5 This matrix is then standardised to take the form of a genetic relationship

6 matrix. Let P_L be the probability that, given the observed population-wide

7 founder haplotype fractions at L , two randomly-sampled lines are IBD, i.e.

8

$$P_L = 2 \sum_{i < j} f_{ijL} / N(N-1)$$

9

10

11 Define

12

$$E_{ij} = \sum_L w_L (f_{ijL} - P_L)$$

$$\mu_i = \sum_j E_{ij} / N$$

$$\sigma_i^2 = \sum_j (E_{ij}^2 / N) - \mu_i^2$$

13

14

15 in order to compute the standardised IBD matrix K_H

16

$$K_{Hij} = \frac{\sum_L w_L (f_{ijL} - \mu_L)}{\sigma_i \sigma_j}$$

1 which has main diagonal 1 and off diagonal elements in the range $[-1,1]$.

2

3 **Identity by State (SNP-based) K_S** SNPs were imputed in the MAGIC lines by

4 using the haplotype mosaics and the catalog of variants in the 19 founders.

5 We treated each MAGIC line as being homozygous. We investigated down-

6 sampling the number of SNPs (as the total is over 1 million). Subsamples of

7 between 1% and 10% of the total SNPs were used to define the GRM in the

8 usual way for a homozygous population. Thus if $S_{ip} \in \{0,1\}$ encodes the SNP

9 genotype in individual i and SNP p , and if π_p is the allele frequency at p , then

10 the normalized genotype is

$$T_{ip} = \frac{S_{ip} - \pi_p}{\sqrt{\pi_p(1 - \pi_p)}}$$

11

12 Since the MAGIC lines are almost fully inbred the normalization is different

13 from that in an outbred population under Hardy-Weinberg equilibrium. The

14 SNP-based GRM is the matrix K_S with elements

15

$$K_{Sij} = \sum_p T_{ip} T_{jp} / M$$

16

17 That is,

$$K_S = TT' / M$$

18

1 which is always positive semi-definite.

2

3 **Read Anomalies K_R** We constructed read-anomaly GRMs by analogy to
 4 SNP-based GRMs. Let X_{iL} be the read anomaly trait for individual i at locus L .
 5 Let α_L and τ_L^2 be the sample means and variances:

$$\alpha_L = \sum_i X_{iL} / N$$

6

$$\tau_L^2 = \sum_i X_{iL}^2 / N - \alpha_L^2$$

7

8 Define the standardized trait matrix W with elements

9

$$W_{iL} = \frac{(X_{iL} - \alpha_L)}{\tau_L}$$

10

11 The genetic relationship between individuals i, j is

12

$$K_{Rij} = \sum_L W_{iL} W_{jL} / M$$

13

14 where M is the number of loci. Then the read anomaly GRM $K_R = WW' / M$.

15

16 This formulation guarantees that the GRM is positive. The choice of loci that
 17 contribute to the GRM can be varied. Loci at which there is no variation in
 18 read anomaly are superfluous and so are ignored. Similarly, loci in which only

1 a small fraction of individuals are anomalous (say <3%) are likely to carry too
2 much weight after normalization and so may optionally be dropped,
3 analogously to the calculation of SNP-based GRMs using only high-frequency
4 SNPs.

5
6 In the MAGIC population, each of the 19 founders should be present at a
7 given locus in about $1/19 = 5.5\%$ of lines. Thus an SV that is private to a
8 single founder should give rise to a trait which is null (ie its un-normalised trait
9 value is zero) in 94.5% of lines on average. Loci with a much smaller fraction
10 of non-null trait values might represent private structural mutations.

11
12 We computed a separate kinship matrix for each of the six measures of read
13 anomaly, and estimated heritability by maximum likelihood.

14
15 **Validation by paired-end data** We used high and low coverage paired-end
16 reads from the 19 founders⁸ and from the MAGIC lines to search for
17 enrichment of read pairs linking the source and sink. For the high-coverage
18 test we restricted attention to the 2,391 SV-QTLs in which founders carrying
19 the SV are predicted, and compared the number of read pairs with one read
20 mapped in the 10kb of the source and within a variably-sized window of W kb
21 from the sink (association peak) ($W \in \{5,20,30,40,50,100,150,200,400\}$) in the
22 founders carrying the SV to the remaining founders using two-sided Fisher's
23 exact tests (FET) at the 5% level of significance. Given the haplotype
24 structure of the MAGIC lines, mapping resolution is variable between QTLs in

1 MAGIC - (200kb being the average in MAGIC²⁵) and may depend on the
2 significance of the association. In the low coverage data we performed the
3 same test comparing the 100 lines with the highest read anomaly trait value to
4 the rest of the population.

5 **Validation by *denovo* contigs** We used BLAT ⁴⁵ to align 5,524,143 short
6 contigs (50-1000bp) from existing *denovo* assembly contigs of the 18 non-
7 reference founder genomes to the reference (Col-0 TAIR10) to identify contigs
8 split across the source and sink locus. After alignment we excluded genomic
9 regions with annotated repeats or transposons and alignments that mapped to
10 over 5 genomic loci. We found 2,619 contigs with alignments split into disjoint
11 pieces over 420 QTLs' sources and sinks, suggesting a cut-and-paste
12 mechanism. We also found 460,656 (8.3%) shared contigs whose alignments
13 overlapped between source and sink regions (duplications, transposons,
14 Microhomology-Mediated Break-Induced Replication (MMBIR) sites and Non-
15 Allelic Homologous Recombination (NAHR) being possible explanations). We
16 randomized the SV-QTLs by circular genome permutation⁴⁶ to determine
17 whether such split and shared contig alignments are overrepresented near
18 SV-QTLs. In particular, for each SV-QTL i , if $a(i), b(i)$ are the original position
19 of the source and sink respectively, then a permuted SV-QTL $a_k(i), b_k(i)$ is
20 defined as:

$$a_k = (a(i) + \theta_k) \bmod L$$

$$b_k = (b(i) + \theta_k) \bmod L$$

$$\theta_k \sim \text{Unif}(0, L)$$

1 with the requirement that $a_k(i), b_k(i)$ must be on the same chromosome for
2 *cis* SV-QTLs. We then computed one-sided p-values $\pi_{split}, \pi_{shared}$. At the 1%
3 level, *trans* SV-QTLs were enriched for both split and shared alignments and
4 *cis* only for split.

5 **Validation by PCR** We designed PCR primers for 77 breakpoints from 44 SV-
6 QTLs predicted from *denovo* contigs. We considered two types of
7 experiments: *type I* experiments had primer oligos corresponding to remote or
8 inverted reference loci so PCR should produce a product in SV genomes and
9 not in the reference; *type II* experiment is a control experiment with the
10 reverse outcome (product in the reference, but not in SV genomes). In total,
11 we designed 96 *type I* experiments, one for each of the 77 breakpoints, and
12 19 control (*type II*) experiments, wherever possible.

13 We designed 20-30bp primer oligos based on the reference (TAIR10), using
14 Primer3⁴⁷, after masking out repeats, transposons and known polymorphisms.
15 SVs tend to be near such sequence features, so we had to relax the default
16 Primer3 criteria to detect oligos, and in particular we required: (i) Maximum
17 allowed product 1.5kb (ii) Annealing temperature 10-90°C (iii) GC-content 10-
18 90% (iv) Self-complementarity 8bp. Primer specificity was tested by BLAT⁴⁵.
19 In 30 (66.6%) SV-QTLs (46 *type I* experiments) at least one breakpoint was
20 confirmed, i.e. there was at least one *type I* experiment which amplified in a
21 subset of founder genomes (those carrying the SV) while not producing a
22 product in the reference, as expected. In a further 7 SV-QTLs (15.6%) (15
23 *type I* experiments) the founders carrying an SV-QTL amplified successfully,

1 but the reference genome also amplified unexpectedly. This may be due to
 2 the presence of highly similar sequence nearby, causing unexpected binding
 3 of one of the primers – potentially in the presence of duplications. Indeed in 10
 4 of these experiments evidence of duplications (multiple bands produced by
 5 PCR) was detected in more than two founder genomes. However, in all 15
 6 experiments there were at least three founder genomes behaving differently
 7 than the reference, indicating that the region is probably structurally variant,
 8 although the type of variant may be different to the one predicted by the
 9 mapping. We conclude that these results, despite ambiguous, probably
 10 indicate SVs, albeit likely polymorphic or of different type than originally
 11 predicted. The remaining 16 *type 1* experiments failed amplify in any founder.
 12 Of the 19 *type 2* experiments, 16 succeeded (worked as expected), 2 were
 13 ambiguous and 1 failed.

14 In total we confirm at least 30 (66.6%) SV-QTLs with at least one breakpoint,
 15 while for a further 7 (15.6%) we have evidence of structural variation – in total
 16 up to 82.2% of the tested SV-QTLs are confirmed.

17 **Association with physiological phenotypes.** For each of the six read
 18 anomaly categories, we computed Pearson correlations and corresponding p-
 19 values between 9 physiological phenotypes and the 11,915 traits measured
 20 genome-wide. We selected significant correlations with $\log P > 4$. After filtering
 21 out correlations driven by outliers (i.e. in which removal of the three most
 22 extreme samples reduced the correlation below the significance threshold) we
 23 found 549 traits associated with 40 phenotypes. Each physiological

1 phenotype had on average 1.56 associated SV traits of the same anomaly
2 type.

3 The effect of SVs on each phenotype was measured by a heritability-like
4 measure, h_{SV}^2 , estimated by linear models. Let y be the vector of phenotypic
5 values for a physiological phenotype with k correlated SV traits (of the same
6 type): X_1, \dots, X_k , represented by the matrix X . The phenotype is modelled as:

$$y = Xa + e$$

7 The k parameters a were estimated using the R function `glm()` and we
8 computed the residual sum of squares RSS and the total sum of squares (i.e.
9 variance of y) TSS. We also computed the individual effect sizes of all traits
10 contributing to the heritability, by fitting simple linear regression models.
11 Based on this analysis SV traits can explain up to 33% of the total phenotypic
12 variance.

13

14 We mapped QTLs for the phenotype residuals after regressing all/each
15 associated SV traits of the same type and compared them to the phenotype
16 QTLs.

17

18 **Published phenotypes** We used flowering time and rosette diameter data
19 from a field experiment³³, as well as phenotypes described in previously²⁵.

20

21 **Phenotyping resistance** Three replicates of each MAGIC line were grown at
22 the University of Bath in 2.5 inch plastic pots. Plants were monitored daily and
23 germination and bolting day recorded. After plants senesced, the

1 inflorescence height and the total number of branches were measured. In a
 2 separate experiment, MAGIC lines were grown in growth chambers in P24
 3 plastic trays and sprayed with *A. laibachii* race Nc14³⁵ when plants were 21
 4 days old. Nc14 zoospores were suspended in water at a concentration of
 5 10⁵ spores per ml and incubated on ice for 30 min. The spore suspension was
 6 then sprayed on plants using a spray gun, and plants were incubated in a cold
 7 room in the dark overnight. Infected plants were kept in a growth chamber
 8 under 10-h light and 14-h dark cycles with a 20°C day and 16°C night
 9 temperature⁴⁸. Resistance was defined as absence of pustules on the leaves
 10 at 7 days after inoculation. To minimize errors in scoring, resistant plants were
 11 monitored up to 14 days after inoculation. The experiment was reproduced
 12 twice.

13

14 **Collection of RNA** We obtained a subset of MAGIC lines from the
 15 Nottingham Arabidopsis Stock Centre (NASC). We grew 209 of the MAGIC
 16 lines at 20°C in Percival environmental chambers (Perry, IA, USA) and
 17 prepared total RNA as previously described for an earlier study with the
 18 MAGIC parental founders⁸; briefly, twenty aerial rosettes from seedlings at the
 19 fourth true leaf stage were pooled⁸. RNAseq library construction and
 20 sequencing was performed at the Oxford Genomics Centre (Oxford, UK) to
 21 produce 2 x 100bp reads using the Illumina non-strand specific method. Per
 22 Illumina HiSeq lane, samples were barcoded and run in 13-plexes to give
 23 approximately 14 million reads per sample.

24

Alignment of RNAseq reads and expression quantification All libraries were aligned to the TAIR10 reference genome using PALMapper v0.6⁴⁹, following a variation-aware alignment approach. Genome variants collected from the 19 founder strains as well as variants reported for a diverse natural population²⁸ were integrated and provided to the aligner as set of known variants. Briefly, the mapper used this set of variants in alignments to prevent reference biases in RNAseq read mapping (see previous work⁸ for a rationale). To facilitate accurate alignments, we further provided splice junction information collected in an earlier study with the founder strains⁸ as well as junction information extracted from the TAIR10 genome annotation. The full alignment parameter set for PALMapper was: -M 3 -G 0 -E 3 -I 12 -L 14 -K 12 -C 14 -I 5000 -NI 1 -SA 5 -UA 50 -CT 50 -JA 15 -JI 1 -z 10 -S -seed-hit-truncate-threshold 100 -report-map-read -report-spliced-read -report-map-region -report-splice-sites 0.9 -filter-max-mismatches 0 -filter-max-gaps 0 -filter-splice-region 5 -min-spliced-segment-len 1 -qpalma-use-map-max-len 10 -f bam -qpalma-prb-offset-fix -junction-remapping <junction_file> -score-annotated-splice-sites <junction_file> -max-dp-deletions 2 -use-variants-editop-filter -use-variants <variant_file> -filter-variants-minuse 1 -merge-variant-source-ids -use-iupac-snp-variants -filter-variants-map-window 20 -iupac-genome -filter-variants-maxlen 100 -index-precache

21

Gene expression quantification. We used a custom python script that counted the number of reads overlapping with at least one exonic position of an annotated gene feature. For each read only the best alignment was

24

1 considered for counting. An alignment was excluded from being counted
2 towards the expression of a gene if (i) at least one position in the alignment
3 overlapped to an annotated intron, (ii) the alignment fell entirely into a region
4 where two or more annotated genes overlap, and (iii) did not start at a position
5 that is part of an exon in all annotated isoforms. For each gene feature the
6 total number of reads passing these filters were used as the expression count.

7

8 **Effects of SVs on gene expression** We considered SVs with accurate
9 breakpoints (see **Validation by *denovo* contigs**). 119 TAIR10 genes
10 spanned SV breakpoints (i.e. were disrupted by SVs) and 6,909 were inside
11 them (transposed, inverted or duplicated). Genes were divided into three
12 categories: disrupted by breakpoints, within SV-regions and outside SVs and
13 compared with respect to mean and variance using t-tests. We also computed
14 the correlation of these genes with their local read anomaly values (for the six
15 read anomaly types), i.e. with the 10kb source region that contains the gene
16 and compared the mean and variance (by a t-test and an F-test, respectively)
17 of the Pearson correlation coefficients across categories.

18 **Acknowledgements**

19

20 MI, RM were supported by the Wellcome Trust Core Award Grant Number
21 090532/Z/09/Z. RMC was supported by NSF 0929262. EJO and RG were
22 supported by National Institutes of Health Genetics Training Grant T32
23 GM07464. We thank Fernando Rabanal for comments on the manuscript.

24

References

1. Simpson, J. T. & Pop, M. The Theory and Practice of Genome Sequence Assembly. *Annu. Rev. Genomics Hum. Genet.* (2015). doi:10.1146/annurev-genom-090314-050032
2. Chaisson, M. J. & Tesler, G. Mapping single molecule sequencing reads using basic local alignment with successive refinement (BLASR): application and theory. *BMC Bioinformatics* **13**, 238 (2012).
3. Jain, M. *et al.* Improved data analysis for the MinION nanopore sequencer. *Nat. Methods* (2015). doi:10.1038/nmeth.3290
4. Cao, J. *et al.* Whole-genome sequencing of multiple *Arabidopsis thaliana* populations. *Nat Genet* **43**, 956–963
5. Mills, R. E. *et al.* Mapping copy number variation by population-scale genome sequencing. *Nature* **470**, 59–65 (2011).
6. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
7. Rimmer, A. *et al.* Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications. *Nat. Genet.* **46**, 912–918 (2014).
8. Gan, X. *et al.* Multiple reference genomes and transcriptomes for *Arabidopsis thaliana*. *Nature* **477**, 419–423 (2011).
9. Ye, K., Schulz, M. H., Long, Q., Apweiler, R. & Ning, Z. Pindel: A pattern growth approach to detect break points of large deletions and medium

- 1 sized insertions from paired-end short reads. *Bioinformatics* **25**, 2865–
2 2871 (2009).
- 3 10. Chen, K. *et al.* BreakDancer: an algorithm for high-resolution mapping
4 of genomic structural variation. *Nat. Methods* **6**, 677–681 (2009).
- 5 11. Simpson, J. T., McIntyre, R. E., Adams, D. J. & Durbin, R. Copy number
6 variant detection in inbred strains from short read sequence data.
7 *Bioinformatics* **26**, 565–567
- 8 12. Manske, H. M. & Kwiatkowski, D. P. LookSeq: A browser-based viewer
9 for deep sequencing data. *Genome Res.* **19**, 2125–2132 (2009).
- 10 13. Sindi, S. S., Onal, S., Peng, L. C., Wu, H.-T. & Raphael, B. J. An
11 integrative probabilistic model for identification of structural variation in
12 sequencing data. *Genome Biol.* **13**, R22 (2012).
- 13 14. Rausch, T. *et al.* DELLY: structural variant discovery by integrated
14 paired-end and split-read analysis. *Bioinformatics* **28**, i333–i339 (2012).
- 15 15. Layer, R. M., Chiang, C., Quinlan, A. R. & Hall, I. M. LUMPY: a
16 probabilistic framework for structural variant discovery. *Genome Biol.*
17 **15**, R84 (2014).
- 18 16. Kronenberg, Z. N. *et al.* Wham: Identifying Structural Variants of
19 Biological Consequence. *PLoS Comput. Biol.* **11**, e1004572 (2015).
- 20 17. Durbin, R. M. *et al.* A map of human genome variation from population-
21 scale sequencing. *Nature* **467**, 1061–1073 (2010).
- 22 18. Cai, N. *et al.* Sparse whole-genome sequencing identifies two loci for
23 major depressive disorder. *Nature* (2015). doi:10.1038/nature14659

- 1 19. Nicod, J. *et al.* Genome-wide association of multiple complex traits in
2 outbred mice by ultra-low-coverage sequencing. *Nat. Genet.* (2016).
3 doi:10.1038/ng.3595
- 4 20. Davies, R. W., Flint, J., Myers, S. & Mott, R. Rapid genotype imputation
5 from sequence without reference panels. *Nat. Genet.* **48**, 965–9 (2016).
- 6 21. Yalcin, B. *et al.* Sequence-based characterization of structural variation
7 in the mouse genome. *Nature* **477**, 326–329 (2011).
- 8 22. Cao, J. *et al.* A whole-genome map of sequence variants in multiple
9 *Arabidopsis thaliana* populations. *Nat Genet In Press*, (2011).
- 10 23. Hu, T. T. *et al.* The *Arabidopsis lyrata* genome sequence and the basis
11 of rapid genome size change. *Nat. Genet.* **43**, 476–81 (2011).
- 12 24. Quadrana, L. *et al.* The *Arabidopsis thaliana* mobilome and its impact at
13 the species level. *Elife* **5**, e15716 (2016).
- 14 25. Kover, P. X. *et al.* A Multiparent Advanced Generation Inter-Cross to
15 fine-map quantitative traits in *Arabidopsis thaliana*. *Plos Genet.* **5**,
16 e1000551 (2009).
- 17 26. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-
18 Wheeler transform. *Bioinformatics* **26**, 589–595 (2010).
- 19 27. Lunter, G. & Goodson, M. Stampy: a statistical algorithm for sensitive
20 and fast mapping of Illumina sequence reads. *Genome Res* **21**, 936–
21 939 (2011).
- 22 28. Long, Q. *et al.* Massive genomic variation and strong selection in
23 *Arabidopsis thaliana* lines from Sweden. *Nat. Genet.* **45**, 884–90 (2013).

- 1 29. Fransz, P. F. *et al.* Integrated cytogenetic map of chromosome arm 4S
2 of *A. thaliana*: structural organization of heterochromatic knob and
3 centromere region. *Cell* **100**, 367–376 (2000).
- 4 30. Wijnker, E. *et al.* The genomic landscape of meiotic crossovers and
5 gene conversions in *Arabidopsis thaliana*. *Elife* **2**, e01426 (2013).
- 6 31. Lai, A. G., Denton-Giles, M., Mueller-Roeber, B., Schippers, J. H. &
7 Dijkwel, P. P. Positional information resolves structural variations and
8 uncovers an evolutionarily divergent genetic locus in accessions of
9 *Arabidopsis thaliana*. *Genome Biol Evol* (2011). doi:evr038
10 [pii]10.1093/gbe/evr038
- 11 32. Berlin, K. *et al.* Assembling large genomes with single-molecule
12 sequencing and locality-sensitive hashing. *Nat. Biotechnol.* **33**, 623–630
13 (2015).
- 14 33. Springate, D. A. & Kover, P. X. Plant responses to elevated
15 temperatures: a field study on phenological sensitivity and fitness
16 responses to simulated climate warming. *Glob. Chang. Biol.* **20**, 456–65
17 (2014).
- 18 34. Nishimura, N. *et al.* ABA hypersensitive germination2-1 causes the
19 activation of both abscisic acid and salicylic acid responses in
20 *Arabidopsis*. *Plant Cell Physiol.* **50**, 2112–22 (2009).
- 21 35. Thines, M. *et al.* A new species of *Albugo* parasitic to *Arabidopsis*
22 *thaliana* reveals new evolutionary patterns in white blister rusts
23 (*Albuginaceae*). *Persoonia Mol. Phylogeny Evol. Fungi* **22**, 123–128

- 1 (2009).
- 2 36. Van Der Biezen, E. A., Freddie, C. T., Kahn, K., Parker, J. E. & Jones,
3 J. D. G. Arabidopsis RPP4 is a member of the RPP5 multigene family of
4 TIR-NB-LRR genes and confers downy mildew resistance through
5 multiple signalling components. *Plant J.* **29**, 439–451 (2002).
- 6 37. Yi, H. & Richards, E. J. Gene duplication and hypermutation of the
7 pathogen Resistance gene SNC1 in the arabidopsis bal variant.
8 *Genetics* **183**, 1227–1234 (2009).
- 9 38. Kang, H. M. *et al.* Efficient control of population structure in model
10 organism association mapping. *Genetics* **178**, 1709–1723 (2008).
- 11 39. Kang, H. M., Ye, C. & Eskin, E. Accurate discovery of expression
12 quantitative trait loci under confounding from spurious and genuine
13 regulatory hotspots. *Genetics* **180**, 1909–1925 (2008).
- 14 40. Durkin, K. *et al.* Serial translocation by means of circular intermediates
15 underlies colour sidedness in cattle. *Nature* **482**, 81–84 (2012).
- 16 41. Genovese, G. *et al.* Using population admixture to help complete maps
17 of the human genome. *Nat. Genet.* **45**, 406–14, 414–2 (2013).
- 18 42. Li, J.-Y., Wang, J. & Zeigler, R. S. The 3,000 rice genomes project: new
19 opportunities and challenges for future rice research. *Gigascience* **3**, 8
20 (2014).
- 21 43. Alonso-Blanco, C. *et al.* 1,135 Genomes Reveal the Global Pattern of
22 Polymorphism in Arabidopsis thaliana. *Cell* (2016).
23 doi:10.1016/j.cell.2016.05.063

- 1 44. Ito, H. & Kakutani, T. Control of transposable elements in *Arabidopsis*
2 *thaliana*. *Chromosom. Res.* **22**, 217–223 (2014).
- 3 45. Kent, W. J. BLAT--the BLAST-like alignment tool. *Genome Res.* **12**,
4 656–664 (2002).
- 5 46. Cabrera, C. P. *et al.* Uncovering networks from genome-wide
6 association studies via circular genomic permutation. *G3 (Bethesda)*. **2**,
7 1067–75 (2012).
- 8 47. Rozen, S. & Skaletsky, H. Primer3 on the WWW for general users and
9 for biologist programmers. *Methods Mol. Biol.* **132**, 365–386 (2000).
- 10 48. Kemen, E. *et al.* Gene gain and loss during evolution of obligate
11 parasitism in the white rust pathogen of *Arabidopsis thaliana*. *PLoS Biol.*
12 **9**, (2011).
- 13 49. Jean, G., Kahles, A., Sreedharan, V. T., De Bona, F. & Ratsch, G. RNA-
14 Seq read alignments with PALMapper. *Curr Protoc Bioinforma.* **Chapter**
15 **11**, Unit 11 6

1 **Supplementary Figures legends**

2 **Figure S1** Manhattan plots and founder contributions for high read coverage
3 in a *cis* (**a,c**) and a *trans* (**b,d**) SV-QTL. In the manhattan plots the red line
4 shows the source and the association peak the sink of the SV-QTL. In the
5 founder contributions tables rows and columns correspond to founder
6 haplotypes at the sink and source, respectively. The colour hue at each cell is
7 the trait value for each combination of founder haplotypes, darker colour
8 means higher value. In **c** trait values range from 0 to 1000. The figure shows a
9 duplication (confirmed by *denovo* contigs) present in 3 founders, namely Bur-
10 0, Col-0 and Edi-0. In **d** trait values (high read coverage) range from 0 to 600
11 and the figure is showing a *trans* QTL in chromosome 5, present in Bur-0, Oy-
12 0, Po-0, Rsch-4 and Wu-0.

13 **Figure S2** Distribution of SV allele frequencies, defined as the fraction of
14 founders carrying the SV allele at SV-QTLs.

15 **Figure S3** Alignment of a 175 kb manually assembled contig from Ler-0,
16 chr5:22.05 – 25.23Mb³¹ to the reference. See legend of Figure 3 for
17 explanation.

18 **Figure S4** Validation of SV QTL predictions using paired-end data from high
19 coverage sequence in the 19 founder genomes. The figure shows results for
20 2,391 SV-QTLs for which predictions of founder haplotypes carrying SVs were
21 reliable. Each bar corresponds to a subsample of the SV QTLs with maximum
22 genome-wide logP exceeding a given threshold, with bar height showing the

1 fraction that were supported by reads at $P < 0.05$. The right-most bar shows the
 2 results in the entire set of 2,391 QTLs. The test compares the number of
 3 reads linking the 10kb source region to a variably sized window around the
 4 sink (see **Methods**). The colours are coding the different window lengths used
 5 – where different window sizes gave significant results for the same SV QTL
 6 we report the smallest.

7

1 **Supplementary Table legends**

2 **Table S1** The 19 founder accessions of the MAGIC population of
3 Recombinant Inbred Lines. Shown are the stock centre numbers, the
4 accessions' names and the place of origin.

5 **Table S2** The mosaic reconstructions of 488 MAGIC lines. Each row
6 represents one segment of a MAGIC line. **magic**: name of the MAGIC line;
7 **chr**: the chromosome of the segment; **acc**: the founder accession present at
8 this locus; **from.bp**, **to.bp**: the start and end bp coordinates of the segment
9 (TAIR10); **from.site**, **to.site**: the start and end position in terms of sites;
10 **len.bp**: the length of the segment in bp; **sites**: the length of the segment in
11 sites; **errors**: the number of errors (sites whose genotype does not agree with
12 the founder accession genotype); **error.site**: the number of errors divided by
13 the number of sites; **error.bp**: the number of errors divided by the length of
14 the segment

15

16 **Table S3**: Catalogue of SV-QTLs detected by genetic mapping of read pair
17 anomalies. Each row represents a distinct SV. SV-QTLs with coincident
18 sources and sinks for different read pair anomalies are merged into a single
19 row. **src.chr**, **src.pos**: genomic location of the source, defined as the start of
20 the 10kb region in which read pair anomalies were measured; **sink.chr**,
21 **sink.pos**: genomic location of the sink, defined as the peak of association;
22 **read.anomaly.traits**: read anomaly traits with the sink SV-QTL;
23 **max.gw.logp**: maximum genome-wide logP estimated by the genome scan

1 (λ_A , see **Methods**), fitted.p – extreme value distribution fitted p-value (γ_A - see
2 **Methods**), - **qtl.distance**: *cis* or *trans*; **SV.type**: prediction of the type of
3 structural variant based on the read anomaly traits that had a SV-QTL ;
4 **sink.founders** – founder haplotypes predicted to carry the SV, NA means
5 unknown (**Methods**); **known.indel.dist**: distance between the midpoints of
6 the source loci to the nearest large (>100bp) SV from ¹⁷;
7 **read.support.founder**: P-value of FET comparing read pairs connecting the
8 source and the sink in founders carrying the SV allele from the high-coverage
9 sequencing of the 19 founders (NA if founders are the sink were unknown) ¹⁷,
10 **read.support.founder.window**: size of window (distance from association
11 peak) containing the significant association in read.support.founder,
12 **read.support.line**, **read.support.line.window**: same validation test using
13 read pairs from the 488 MAGIC lines; **denovo.contigs**: Boolean variable
14 showing confirmation by at least one *denovo* contig.

15

16 **Table S4.** Breakpoints of 420 SVs detected using *denovo* contigs. **founder**:
17 founder genome in which the breakpoint was detected; **source.chr**,
18 **source.pos**: source position in which the read anomalies were measured;
19 **sink.chr**, **sink.pos**: SV-QTL position; **source.break.from**, **source.break.to**:
20 breakpoints detected by *denovo* contigs corresponding to the source area;
21 **sink.break.from**, **sink.break.to**: breakpoints corresponding to sink;
22 **source.length**, **sink.length**: length of the structurally variant region in the
23 source and sink regions.

24

Table S5. PCR validation results **(a)** Results per SV-QTL. **QTL_ID**: id of the SV-QTL; **QTL**: coordinates of the tested SV-QTL, position of the source followed by the position of the sink; **confirm**: Y – yes, A- ambiguous, N – no; **dist**: Boolean indicator of whether the SV-QTL is *cis* or *trans*; **read pred**: type of the SV (e.g. transposition, indel etc) predicted by read anomalies, PCR pred: prediction of the type of SV based on the contigs and PCR results. **(b)** Results per experiment. Each row corresponds to a single experiment (unique combination of primers) performed on all 19 founders. **QTL_ID**: ID of the SV-QTL predicted; **type**: type of experiment 1 or 2 (**see Methods**); **Forward**, **Reverse**: unique identifier of the primers used, in the form chr_pos_orientation, (e.g. id 4_1853989F means the sequence starting from chr 4, 1853989bp and with forward orientation). The identifier INV means that the sequence has the opposite orientation than expected. The remaining columns correspond to each of the founder genomes, 1 meaning that the experiment amplified, 0 that it did not, and >1 that it produced multiple bands. **(c)** Primer sequences.

Table S6 Nine physiological traits measured in MAGIC lines that were used in this study.

Table S7. Effects of SVs on physiological phenotypes. **Phenotype**: physiological phenotype; **SV trait type**: type of read pair anomaly, **SV-QTL**: number of SV traits that have SV-QTLs, **trans**: number of trans SV-QTLs; **max.source**: position of the source of the maximum-contributing SV trait; **max.sink**: position of the sink QTL of the maximum-contributing SV, or NA if it

1 is not mapped, **physio.QTL**: position of the QTL of the physiological
2 phenotype; **overlap.source**, **overlap.sink**: Boolean indicator of whether the
3 source or sink of the SV trait overlaps (within 200kb) with the QTL of the
4 physiological phenotype, NA if there is no QTL.

5 **Table S8**. Association of SVs with gene expression. **gene** – gene id, **chrom**,
6 **start**, **end** – gene coordinates, Same QTL: position of the gene relative to
7 SVs (break: the gene is spanning an SV breakpoint, within: the gene is
8 between SV breakpoints so it may be transposed, inverted or duplicated,
9 outside: the gene lies outside SVs), **mean**: mean gene expression, **var**:
10 variance of gene expression, **zeroes**: proportion of silenced (zero expression)
11 transcripts, **r_unmapped_largeisize**, **r_impop_paired_LT**, **r_same_strand**,
12 **r_largeisize**, **r_excess_reads**, **r_unmapped**: Pearson correlation
13 coefficients between expression levels and values of the local read anomaly
14 trait, for each of the six anomaly types.

15

16