

1 **Transient structural variations have strong effects**
2 **on quantitative traits and reproductive isolation in fission yeast.**

3
4 October 3, 2016

5
6 **Authors**

7 Daniel C. Jeffares^{1,2}, Clemency Jolly¹, Mimoza Hoti¹, Doug Speed², Liam Shaw^{1,2},
8 Charalampos Rallis^{1,2,3}, Francois Balloux^{1,2}, Christophe Dessimoz^{1,4,5,6*}, Jürg Bähler^{1,2*},
9 Fritz J. Sedlazeck^{7*}

10
11 **Affiliations:**

12 1. Department of Genetics, Evolution and Environment, University College London,
13 London WC1E 6BT, United Kingdom.

14 2. UCL Genetics Institute, University College London, London WC1E 6BT, United
15 Kingdom.

16 3. Current address: School of Health, Sport and Biosciences, University of East London,
17 London E15 4LZ, United Kingdom.

18 4. Department of Computer Science, University College London, London WC1E 6BT,
19 United Kingdom.

20 5. Department of Ecology and Evolution & Center for Integrative Genomics, University
21 of Lausanne, Biophore, 1015 Lausanne, Switzerland.

22 6. Swiss Institute of Bioinformatics, Biophore, 1015 Lausanne, Switzerland

23 7. Department of Computer Science, Johns Hopkins University, 21218 Baltimore, USA

24
25 * Correspondence should be addressed to C.D. (christophe.dessimoz@unil.ch), J.B.
26 (j.bahler@ucl.ac.uk) or F.J.S. (fritz.sedlazeck@jhu.edu)

27
28

29 **ABSTRACT**

30 **Large structural variations (SVs) in the genome are harder to identify than smaller**
31 **genetic variants but may substantially contribute to phenotypic diversity and**
32 **evolution. Here we analyze the effects of SVs on gene expression, quantitative traits,**
33 **and intrinsic reproductive isolation in the yeast *Schizosaccharomyces pombe*. We**
34 **establish a high-quality curated catalog of SVs in the genomes of a worldwide**
35 **library of *S. pombe* strains, including duplications, deletions, inversions and**
36 **translocations. We show that copy number variants (CNVs) frequently segregate**
37 **within closely related clonal populations, are weakly linked to single nucleotide**
38 **polymorphisms (SNPs), and show other genetic signals consistent with rapid**
39 **turnover. These transient CNVs produce stoichiometric effects on gene expression**
40 **both within and outside the duplicated regions. CNVs make substantial**
41 **contributions to quantitative traits such as cell shape, cell growth under diverse**
42 **conditions, sugar utilization in winemaking, whereas rearrangements are strongly**
43 **associated with reproductive isolation. Collectively, these findings have broad**
44 **implications for evolution and for our understanding of quantitative traits including**
45 **complex human diseases.**

46

47 **Keywords:** structural variants, yeast, copy number variants, reproductive isolation,

48 quantitative genetics, next generation sequencing

49

50

51 A variety of genetic changes can influence the biology of species, including single-
52 nucleotide polymorphisms (SNPs), small insertion-deletion events (indels), transposon
53 insertions and large structural variations. Structural variations (SVs), including deletions,
54 duplications, insertions, inversions and translocations, are the most difficult to type and
55 consequently the least well described.

56 Nevertheless, it is clear that SVs have strong effects on various biological
57 processes. Copy number variants (CNVs) in particular influence quantitative traits in
58 microbes, plants and animals, including agriculturally important traits and a variety of
59 human diseases¹⁻⁵. Inversions are known to influence reproductive isolation⁶⁻¹³ and other
60 evolutionary processes such as recombination⁸ and hybridization between species¹⁴, with
61 a variety of consequences¹⁵.

62 We and others have recently begun to develop the fission yeast
63 *Schizosaccharomyces pombe* as a model for population genomics and quantitative trait
64 analysis^{6,7,16-18}. This model organism combines the advantages of a small, well-annotated
65 haploid genome¹⁹, abundant tools for genetic manipulation and high-throughput
66 phenotyping²⁰, and considerable resources of genome-scale and gene-centric data²¹⁻²³.

67 Previous analyses of fission yeast have begun to describe both naturally occurring
68 and engineered inversions and reciprocal translocations^{6,7,18}. Given this evidence for SVs
69 and their effects in this model species, we recognized that a systematic survey of SVs
70 would progress our understanding of their biological influence. Here, we utilize the
71 recent availability of 161 fission yeast genomes and extensive data on quantitative traits
72 and reproductive isolation¹⁷ to describe the nature and effects of SVs in *S. pombe*.

73 We show that SVs have strong effects on various aspects of biology. They
74 contribute an average of 11% of trait variance (the much more abundant SNPs contribute
75 24% on average), with the largest effects coming from CNVs. We show that CNVs are
76 transient within clonal populations, and are frequently not well tagged by SNPs. We also
77 show that rearrangements (inversions and translocations) contribute to reproductive
78 isolation, whereas CNVs do not.

79

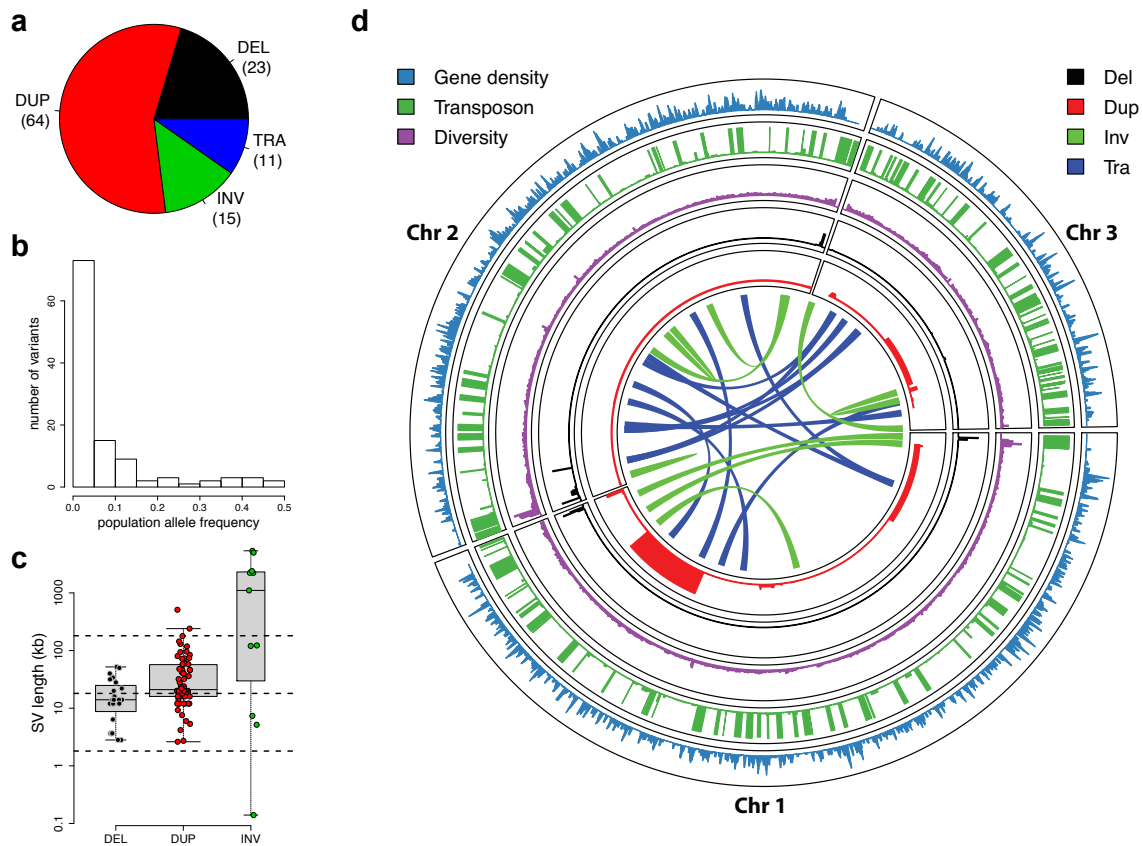
80 **RESULTS**

81 **Genome- and population-wide detection of structural variations**

82 To predict an initial set of SVs, we applied four inference software packages (Delly,
83 Lumpy, Pindel and cn.MOPs)²⁴⁻²⁷ to existing short-read data¹⁷, using parameters
84 optimized on simulated data (Methods). We then filtered these initial predictions,
85 accepting SVs detected by at least two callers, to obtain 315 variant calls (141 deletions,
86 112 duplications, 26 inversions, 36 translocations). We release this pipeline as an open-
87 source tool called SURVIVOR (Methods). To ensure a high specificity, we further
88 filtered the 315 variants by removing SV calls whose breakpoints overlapped with low
89 complexity regions or any that corresponded to previously annotated long terminal
90 repeats (LTRs)¹⁷. Finally, we manually vetted all the remaining SVs by visual inspection
91 of read alignments in multiple strains for all remaining candidates. This meticulous
92 approach aimed to ensure a high quality call set, to mitigate against the high uncertainty
93 associated with SV calling²⁵.

94 This curation produced a set of 113 SVs, comprising 23 deletions, 64
95 duplications, 11 inversions and 15 translocations (**Figure 1a**). Reassuringly, when

96 applying our variant calling methods to an engineered knockout strain, we correctly
97 identified the known deletions and called no false positives. Attempts to validate all
98 rearrangements by PCR and BLAST searches of *de novo* assemblies positively verified
99 76% of the rearrangements, leaving only a few PCR-intractable variants unverified (see
100 Methods for details).



101

102 **Figure 1. Characteristics of SVs in *S. pombe*.** (a) Relative proportions of SVs
103 identified. Duplications (DUP) were the most abundant SVs, followed by deletions
104 (DEL), inversions (INV), and translocations (TRA). (b) Population allele frequency
105 distribution of SVs, showing the frequencies of less abundant alleles in the population
106 (minor allele frequencies). (c) Length distributions of SVs, log₁₀ scale. Deletions were
107 smallest (2.8–52 kb), duplications larger (2.6–510 kb), and inversions often even larger,
108 spanning large portions of chromosomes (0.1 kb–5,374 kb, see (d)). Horizontal dotted
109 lines show the size of chromosome regions that contain an average of 1, 10 and 100 genes
110 in this yeast. (d) Locations of SVs on the three chromosomes compared to other genomic
111 features. From outside: density of essential genes, locations of *Tf*-type retrotransposons,
112 diversity (π , average pairwise diversity from SNPs), deletions (black), duplications (red),

113 and breakpoints of inversions and translocations as curved lines inside the concentric
114 circles (green and blue, respectively). Bar heights for retrotransposons, deletions and
115 duplications are proportional to minor allele frequencies. Diversity and retrotransposon
116 frequencies were calculated from 57 non-clonal strains as described by Jeffares, et al. ¹⁷.
117
118 Most SVs were present at low frequencies, with 28% discovered in only one of the strains
119 analyzed (**Figure 1b**). The deletions were generally slightly smaller (median length 14
120 kb, **Figure 1c**) than duplications (median length of 21 kb), with the largest duplication
121 extending to 510 kb and covering 200 genes (a singleton in strain JB1207/NBRC10570).
122 The majority of CNVs were present in copy numbers varying between zero and sixteen
123 (subsequently we refer to amplifications of two or more copies as ‘duplications’).
124 All SVs, particularly deletions and duplications, were biased towards the ends of
125 chromosomes (**Figure 1d, Supplementary Figures 1 and 2**), which are characterized by
126 high genetic diversity, frequent transposon insertions, and a paucity of essential genes¹⁷,
127 similar to *Saccharomyces cerevisiae* and *S. paradoxus*^{28,29}. All SVs preferentially
128 occurred in positions of low gene density and were strongly under-enriched in essential
129 genes (**Supplementary Figure 2**).
130
131 To describe SVs further, we conducted gene enrichment analysis with the AnGeLi tool
132 (**Supplementary Table 1**), which interrogates gene lists for functional enrichments using
133 multiple qualitative and quantitative information sources³⁰. The CNV-overlapping genes
134 were enriched for caffeine/rapamycin induced genes and genes induced during meiosis (P
135 = 4×10^{-7} and 1×10^{-5} , respectively); they also showed lower relative RNA polymerase II

136 occupancy and were less likely to contain genes known to produce abnormal cell
137 phenotypes ($P = 1.8 \times 10^{-5}$ and 3×10^{-5} , respectively). These analyses are all broadly
138 consistent with a paucity of CNVs in genes that encode essential mitotic functions.
139 Rearrangements disrupted only a few genes and showed no significant enrichments.

140

141 **Duplications are transient within clonal populations**

142 Our previous work identified 25 clusters of near-clonal strains, which differed by <150
143 SNPs within each cluster¹⁷. We expect that these clusters reflect either repeat depositions
144 of strains differing only at few sites (e.g. mating-type variants of reference strains h^{90} and
145 h differ by 14 SNPs) or natural populations of strains collected from the same location.
146 Such ‘clonal populations’ reflect products of mitotic propagation from a very recent
147 common ancestor, without any outbreeding. We therefore expected that SVs should be
148 largely shared within these clonal populations.

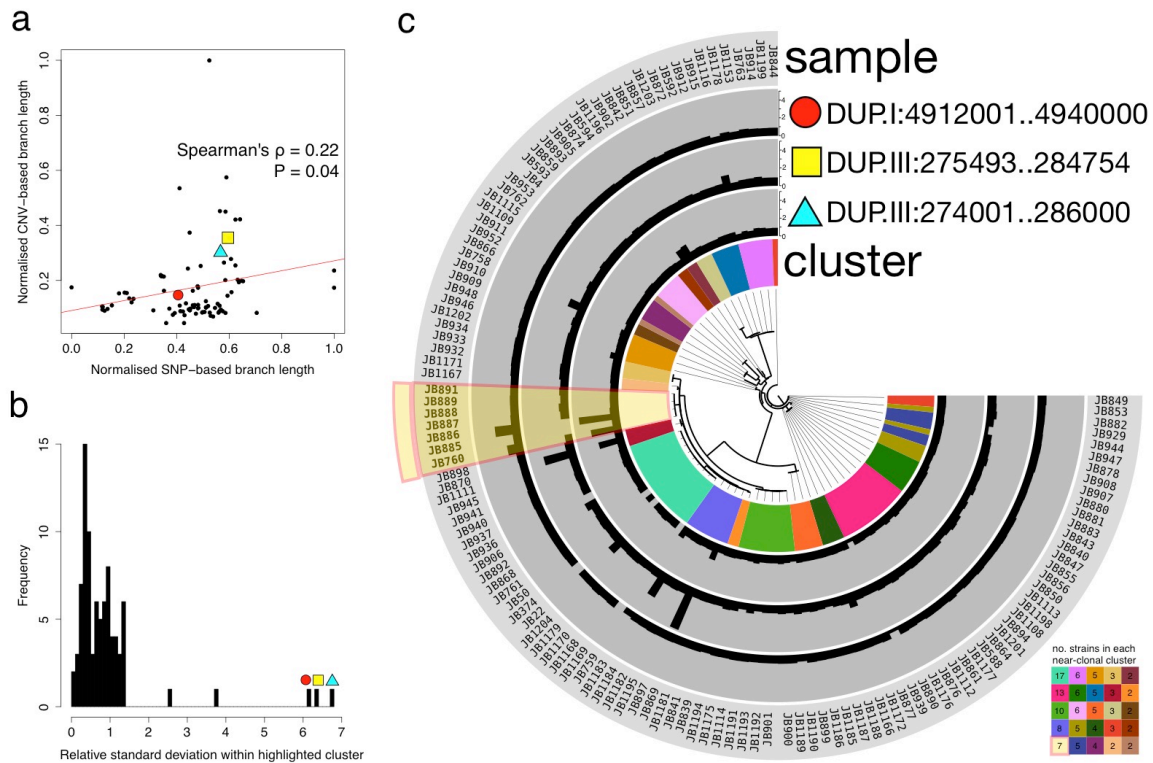
149 Surprisingly, our genotype predictions indicated that most SVs present in clonal
150 populations were segregating, i.e. were not fixed within the clonal population (68/95 SVs,
151 72%). Furthermore, we observed instances of the same SVs that were present in two or
152 more different clonal populations that were not fixed within any clonal population. These
153 SVs could be either incorrect allele calls in some strains, or alternatively, recent events
154 that have emerged during mitotic propagation. To distinguish between these two
155 scenarios, we re-examined the read coverage of all 49 CNVs present within at least one
156 clonal population. Since translocations and inversions were more challenging to
157 accurately genotype, we did not re-examine these variants. This analysis verified that 40
158 of these 49 CNVs (37 duplications, 3 deletions) were clearly segregating within at least

159 one clonal cluster (**Supplementary Figure 3**). For example, one clonal population of
160 seven closely related strains, collected together in 1966 from grape must in Sicily, have
161 an average pairwise difference of only 19 SNPs (diversity $\pi = 1.5 \times 10^{-6}$). Notably, this
162 collection showed four non-overlapping segregating duplications (**Fig. 3c**). This striking
163 finding suggests that CNVs can arise or disappear frequently during evolution.

164 To examine whether this transience is a general feature of CNVs in this population,
165 we quantified the variation in copy number of each CNV relative to mutations in the
166 adjacent region of the genome. If a CNV was subject only to the same processes as these
167 adjacent regions, we would expect a strong correlation between the total mutation in
168 these regions and the total variation in copy number of the CNV. However, the variation
169 in copy number of CNVs across the dataset was only weakly correlated with SNP
170 variation in nearby regions of the genome (Spearman rank correlation $\rho = 0.22$, $P =$
171 0.041), indicating that CNVs are subject to additional or different evolutionary processes
172 (**Figure 2a**). Furthermore, some CNVs showed high rates of variation within closely-
173 related clusters relative to their variation in the rest of the dataset (**Figure 2b and 2c**,
174 **Supplementary Table 2, Supplementary Figure 4**). Finally, we found that many CNVs
175 represented the rare allele within the cluster, consistent with events that have short half-
176 lives (**Supplementary Figure 5**). Taken together, these results indicate that CNVs are
177 transient and variable features of the genome, even within extremely closely related
178 strains.

179

180



181

182 **Figure 2. Copy number variants are transient within fission yeast.** (a) For each of

183 the 87 CNVs we calculated the genetic distance between strains using SNPs in the region

184 around the CNV (20 kb up- and down-stream of the CNV, merged) as the total branch

185 length from an approximate maximum-likelihood tree (x -axis, SNP-based branch length

186 normalised to maximum value). We further calculated a CNV-based distance using the

187 total branch length from a neighbour-joining tree constructed from Euclidean distances

188 between strains based on their copy numbers (y -axis, CNV-based branch length

189 normalised to maximum value). The weak correlation indicates that CNVs are subject to

190 additional or different evolutionary processes. (b) Histogram of the standard deviation of

191 each CNV within a near-clonal cluster (see also Figure 2a), relative to its standard

192 deviation across strains not in the near-clonal cluster. Standard deviation is highly

193 correlated with CNV-based branch length (Spearman rank correlation $\rho = 0.90$, $P <$
194 0.001) (**Supplementary Figure 4b**). The highlighted CNVs have unusually high rates of
195 variation within this cluster compared to other clusters. **(c)** Copy number variation of
196 these highlighted CNVs plotted on a SNP-based phylogeny (20kb up- and down-stream
197 of the DUP.III:274001..286000 CNV) shows their relative transience within the cluster,
198 as well as their variation across other near-clonal clusters. SNP-based phylogenies for the
199 other two selected CNVs also do not separate the strains with different copy numbers
200 (individual plots for each CNV across clusters for its corresponding SNP-based
201 phylogeny are available as Supplementary data).

202

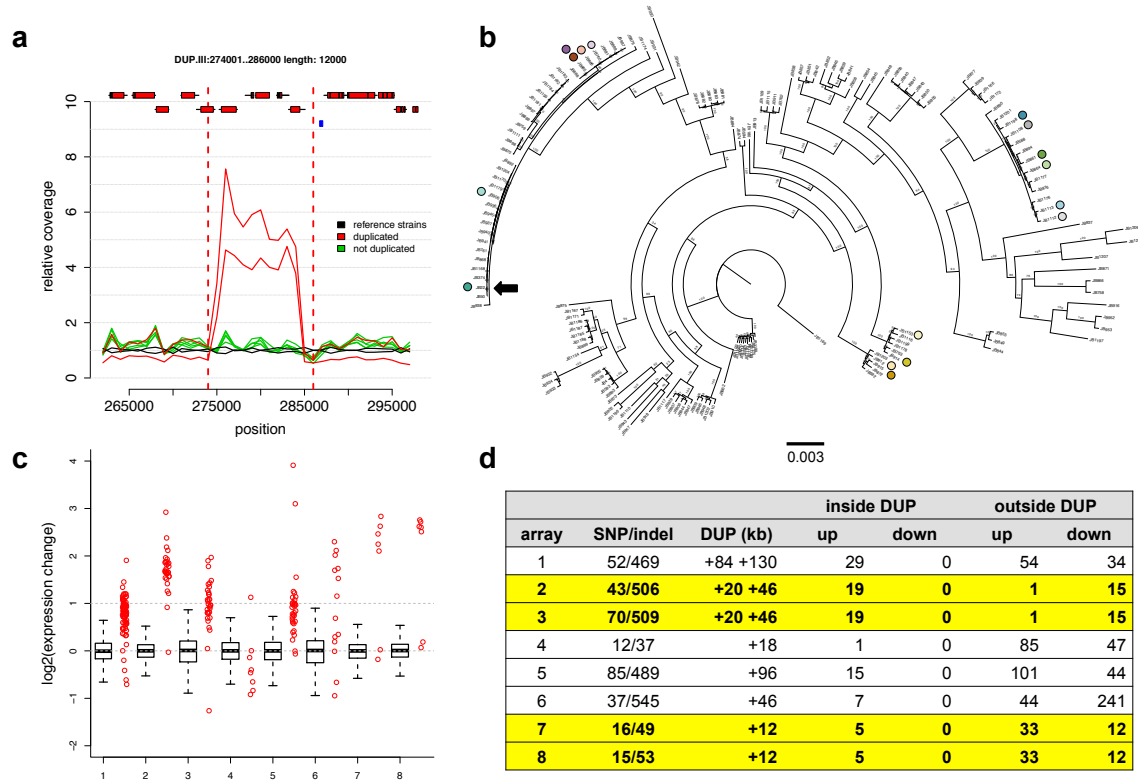
203

204 **Transient duplications affect gene expression**

205 Partial aneuploidies of 500-700 kb in the *S. pombe* reference strain are known to alter
206 gene expression levels within and, to some extent, outside of the duplicated region³¹. The
207 naturally occurring duplications described here are typically smaller (median length: 21
208 kb), including an average of 6.5 genes. To examine whether naturally occurring CNVs
209 have similar effects on gene expression, we examined eight pairs of closely related strains
210 (<150 SNPs among each pair) that contained at least one unshared duplication (**Figure 3,**
211 **Supplementary Table 3**). Several of these strain pairs have been isolated from the same
212 substrate at the same time, and all pairs are estimated to have diverged approximately 50
213 to 65 years ago (**Supplementary Table 3**). We assayed transcript expression from log
214 phase cultures using DNA microarrays, each time comparing a duplicated to a non-
215 duplicated strain from within the same clonal population. In seven of the eight strain
216 pairs, the expression levels of genes within duplications were significantly induced,

217 although the degree of expression changes between genes was variable (**Figure 3c**,
218 **Supplementary Figure 6**). The increased transcript levels correlated with the increased
219 genomic copy numbers, so that higher copy numbers produced correspondingly more
220 transcripts (Spearman rank correlation $\rho = 0.71$, $P = 0.014$, **Supplementary Figure 7**).
221 No changes in gene expression were evident immediately adjacent to the duplications
222 (**Supplementary Figure 7**), suggesting that the local chromatin state was not strongly
223 altered by the CNVs. This result not only confirms the previous observation that CNVs
224 alter the gene expression levels, but more importantly it reveals large copy number
225 differences between two genomes that are only 19 SNPs apart.

226 Interestingly, some genes outside the duplicated regions also showed altered
227 expression levels (**Figure 3d**, **Supplementary Table 4**). For example, two strain pairs
228 differ by a single 12 kb duplication. Here, five of seven genes within the duplication
229 showed induced expression, while 45 genes outside the duplicated region also showed
230 consistently altered expression levels (38 protein-coding genes, 7 non-coding RNAs)
231 (**Figure 3d**, arrays 7 and 8). As environmental growth conditions were tightly controlled,
232 these changes in gene expression could be due to either compensatory effects of the
233 initial perturbation caused by the 12kb duplication or changes that arise due to SNPs or
234 indels that segregate between the strains (**Supplementary Figure 6**). We conclude that
235 these evolutionary unstable duplications reproducibly affect the expression of distinct sets
236 of genes and thus have the potential to influence cellular function and phenotypes.



237

238 **Figure 3. Transient duplications affect gene expression.** (a) Duplications occur within
 239 near-clonal strains. Plot showing average read coverage in 1 kb windows for two clonal
 240 strains (JB760, JB886) with the duplication (red), five strains without duplication (green),
 241 and two reference strains (h^+ , and h^-) (black). Genes (with exons as red rectangles) and
 242 retrotransposon LTRs (blue rectangles) are shown on top (see **Supplementary Table 3**
 243 for details). (b) Eight pairs of closely related strains, differing by one or more large
 244 duplications, selected for expression analysis. The tree indicates the relatedness of these
 245 strain pairs (dots colored as in d). The position of the reference strain (Leupold's 972,
 246 JB22) is indicated with a black arrow. The scale bar shows the length of 0.003 insertions
 247 per site. (c) Gene expression increases for most genes within duplicated regions. For each
 248 tested strain pair, we show the relative gene expression (strains with duplication/strains
 249 without duplication) for all genes outside the duplication (as boxplot) and for all genes

250 within the duplication (red strip chart). In all but one case (array 4), the genes within the
251 duplication tend to be more highly expressed than the genes outside of the duplication (all
252 Wilcoxon rank sum test P-values $<1.5 \times 10^{-3}$). **(d)** Summary of expression arrays 1-8, with
253 strains indicated as colored dots (as in b), showing number of single-nucleotide
254 polymorphism differences between strains (SNP), sizes of duplications in kb (DUP,
255 where '+X +Y' indicates two duplications with lengths X and Y, respectively). We show
256 total numbers of induced (up) and repressed (down) genes, both inside and outside the
257 duplicated regions. Arrays 2,3 and 7,8 (in yellow shading) are replicates within the same
258 clonal population that contain the same duplications, so we list the number of up- and
259 down-regulated genes that are consistent between both arrays. See **Supplementary**
260 **Tables 3 and 4** for details.

261

262 **Copy number variants contribute to quantitative trait variance**

263 To test whether SVs affect phenotypes, we examined the contributions of SNPs, CNVs
264 and rearrangements to 227 quantitative traits (**Supplementary Table 5**), including 20 cell
265 shape parameters, colony size on solid media assaying 42 stress and nutrient conditions¹⁷,
266 126 growth parameters in liquid media conditions⁷ and three biochemical parameters
267 from wine fermentation³². For each phenotype, we used mixed model analysis to estimate
268 the total proportion of variance explained by the additive contribution of genomic
269 variants (the narrow-sense heritability).

270 When we determined heritability using only SNP data, estimates varied between
271 0% and 74% (median 30%). After adding CNVs and rearrangements to SNPs in a
272 composite model, the estimated overall heritability increased for nearly all traits,

273 explaining up to ~40% of trait variance (**Figure 4a**). This finding indicates that the CNVs
274 and rearrangements can explain a substantial proportion of the trait variance. Using this
275 composite model, we quantified the individual contributions of heritability best explained
276 by SNPs, CNVs and rearrangements (**Figure 4b**). On average, SNPs explained 24% of
277 trait variance, CNVs 7%, and rearrangements 4% (**Supplementary Table 5**). Analysis of
278 simulated data confirmed that the contribution of CNVs could not be explained by
279 linkage to causal SNPs alone (**Supplementary Figure 6**).

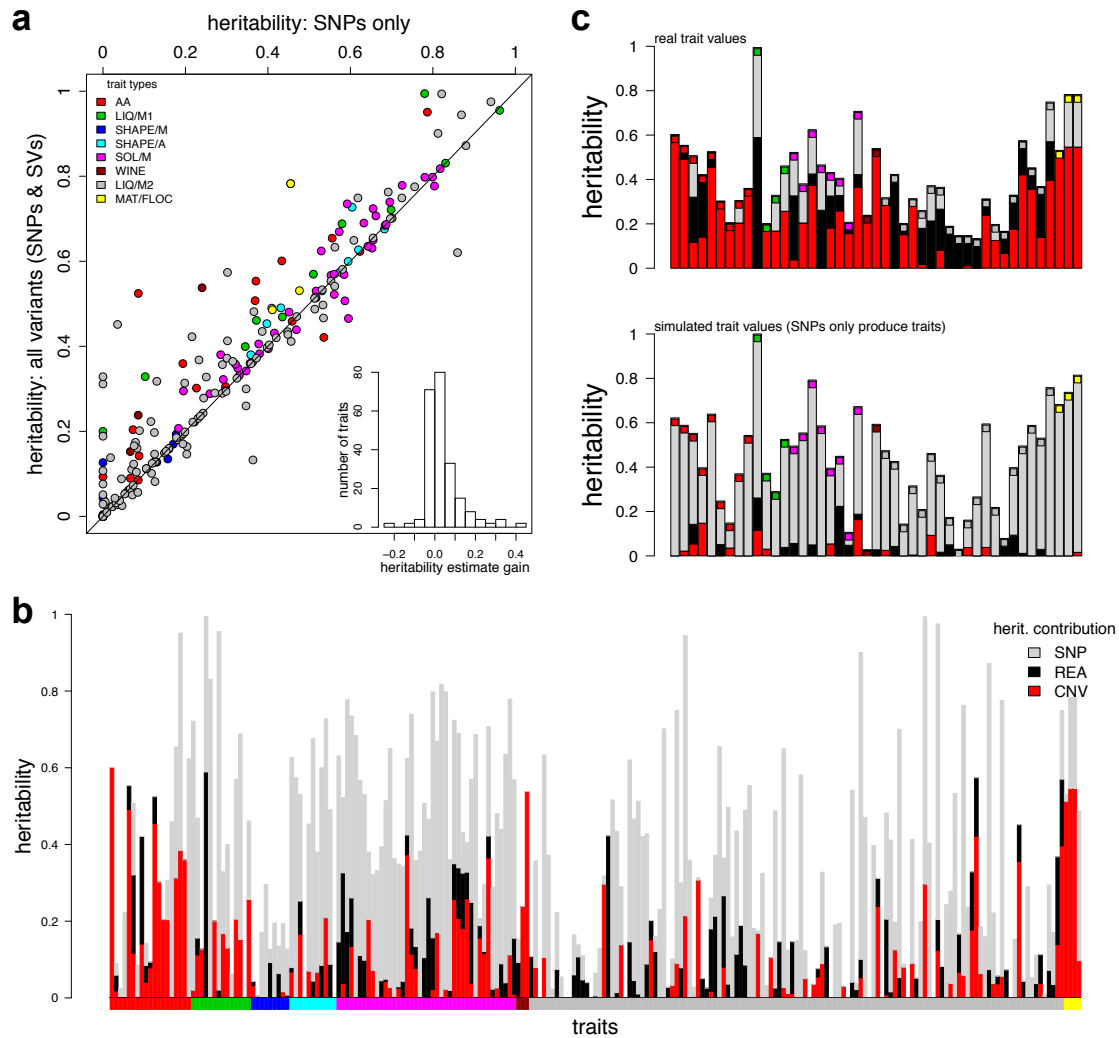
280 Many trait measures gathered using the same method (e.g., growth on solid
281 media, cell shape) are strongly correlated¹⁷. Thus, some groups of traits have consistently
282 larger contributions from SVs (**Figure 4b**) than from SNPs alone. These traits include
283 intracellular amino acid concentrations, growth under stress and several traits measured
284 during wine fermentation (**Figure 4c**). Since many of these strains have been collected
285 from fermentations (**Supplementary Table 6**), the substantial influence of CNVs may
286 represent recent strong selection and adaptation to fermentation conditions that has
287 occurred via recent CNV acquisition.

288 Our analysis of heritability showed that SNPs are generally able to capture most,
289 but not all, of the genetic contribution of SVs (**Figure 4**). To examine whether trait-
290 influencing SVs would be effectively detected from SNPs alone in this population, we
291 examined the linkage of all 113 SVs with SNPs. We found that only 63 of these SVs
292 (55%) are in strong linkage to SNPs ($r^2 > 0.6$), leaving 45% of the SVs weakly linked.
293 This lack of linkage is consistent with SVs being transient, rather than persisting within
294 haplotypes. Such weakly linked SVs may be missed in SNP-only association studies.

295 To examine this possibility, and to locate specific SVs that affect these traits, we

296 performed mixed model genome-wide association studies, using all 68 SVs with minor
297 allele counts >5 (i.e. occurring in at least 5 strains) as well as 139,396 SNPs and 22,058
298 indels with minor allele counts >5. Trait-specific significance thresholds for 5%
299 familywise error rates were computed via permutation analysis, and were approximately
300 10^{-4} (SVs) and 10^{-6} (SNPs and indels). Nineteen SVs (28%) were significantly associated
301 with traits (15 duplications, 5 deletions, 1 translocation), as well as 228 SNPs (0.16%),
302 and 93 indels (0.42%) (**Supplementary Table 7**). SVs were associated with 20 different
303 traits, including amino acid concentrations, mating traits, and stress resistance in solid
304 and liquid media. Nine of these SVs were not strongly linked to SNPs ($r^2 < 0.6$). The
305 median effect size of these SVs was 14% (range 6-33%). While more detailed analyses of
306 these associations will be required to confirm any particular association, our findings are
307 consistent with the heritability analysis.

308 Collectively, these analyses indicate that even a small collection of SVs, most
309 notably CNVs, can contribute substantially to quantitative traits. Thus, GWAS analyses
310 conducted without genotyping SVs could fail to capture these important genetic factors.
311



312

313 **Figure 4. SVs contribute to quantitative traits.** (a) Heritability estimates are improved
314 by the addition of SVs. Heritability estimates for 227 traits (**Supplementary Table 5**),
315 using only SNP data (x axis) range from 0 to 96% (median 29%). Adding SV calls (y
316 axis) increases the estimates (median 34%), with estimates for some traits being
317 improved up to a gain of 43% (histogram inset). The diagonal line shows where estimates
318 after adding SVs are the same as those without ($x=y$). Inset: the distribution of the ‘gain’
319 in heritability after adding SV calls (median 0.4%, maximum 43%). Points are colored by
320 trait types, according to legend top left. (b) The contributions of SNPs (grey), CNVs (red)
321 and rearrangements (black) to heritability varied considerably between traits. Coloured

322 bars along the x axis indicate the trait types. heritability estimates are in **Supplementary**
323 **Table 5**. The panel below bars indicates trait types as in the legend for part (a). (c) Top
324 panel, for some traits, SVs explained more of the trait variation than SNPs. Boxes are
325 colored as legend in part (a). Lower panel, analysis of simulated data generated with
326 assumption that only SNPs cause traits indicates that the contribution of SVs to trait
327 variance is unlikely to be due to linkage.

328

329 **Structural variations contribute to intrinsic reproductive isolation**

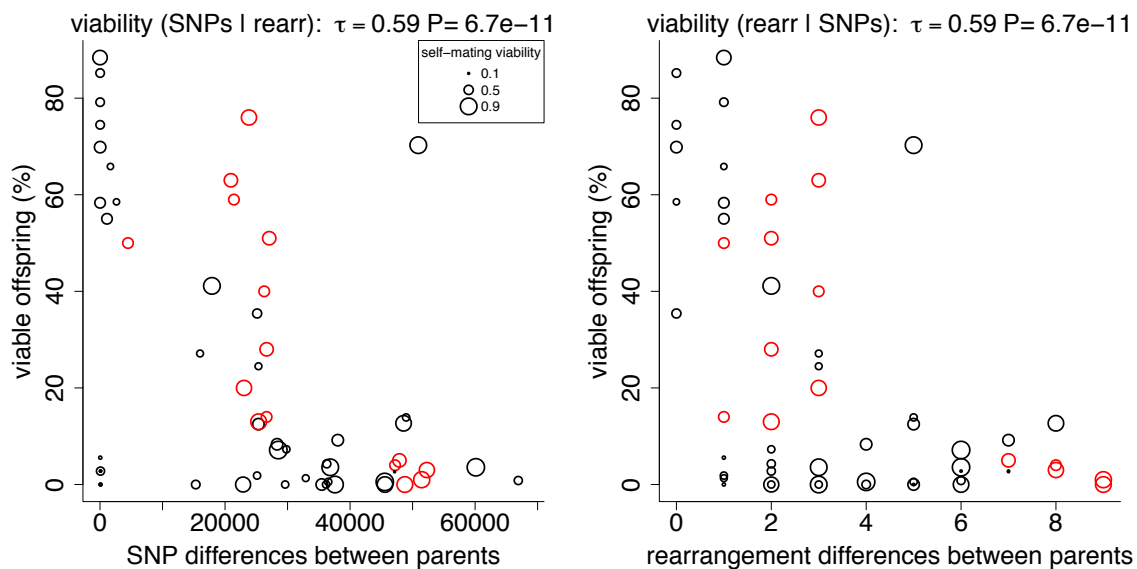
330 Crosses between *S. pombe* strains produce between <1% and 90% viable offspring^{6,18}.
331 We have previously shown that spore viability correlates inversely with the number of
332 SNPs between the parental strains¹⁷. This intrinsic reproductive isolation may be due to
333 the accumulation of Dobzhansky-Muller incompatibilities (variants that are neutral in one
334 population, but incompatible when combined)^{33,34}. However, genetically distant strains
335 also accumulate SVs, which are known to lower hybrid viability and drive reproductive
336 isolation⁹. In *S. pombe*, engineered inversions and translocations reduce spore viability by
337 ~40%⁶. At present the impact of naturally occurring rearrangements, sequence
338 divergence, and incompatible alleles in speciation within budding yeast is unclear¹²⁻
339 ^{14,35,36}.

340 To analyse intrinsic reproductive isolation in our population based on naturally
341 occurring SVs, we examined the relationship between viability, SNPs and SVs. Both SV-
342 distance (number of unshared SVs between parents) and SNP-distance inversely
343 correlated with hybrid viability (Kendall correlation coefficients, SVs: $\tau = -0.26$, $P = 5.6$
344 $\times 10^{-3}$, SNPs: $\tau = -0.35$, $P = 1.6 \times 10^{-4}$) (**Supplementary Figure 7**). While inversions and

345 translocations are known to lower hybrid viability as they affect chromosome pairing and
346 segregation during meiosis^{6,18,37}, CNVs are not expected to influence spore viability.
347 Consistent with this view, there was no significant correlation between CNVs and
348 viability (rearrangements, $\tau = -0.36$, $P = 2.0 \times 10^{-4}$; CNVs, $\tau = -0.10$, $P = 0.28$).

349 As the numbers of SNP and rearrangement differences between mating parents
350 are themselves correlated ($\tau = 0.53$, $P = 1.3 \times 10^{-8}$), we also estimated the influence of each
351 factor alone using partial correlations. When either SNPs or rearrangements were
352 controlled for, both remained significantly correlated with offspring viability ($P = 0.04$, P
353 $= 0.02$, respectively) (**Figure 5**). Taken together, these analyses indicate that both
354 rearrangements and SNPs contribute to reproductive isolation, but CNVs do not.

355



356

357

358 **Figure 5. Both SNPs and rearrangements contribute to intrinsic reproductive**

359 **isolation.** Spore viability was measured from 58 different crosses from Jeffares, et al.¹⁷

360 (black) or Avelar, et al.⁶ (red), with each circle in the plots representing one cross. An

361 additive linear model incorporating both SNP and rearrangement differences showed
362 highly significant correlations with viability ($P = 1.2 \times 10^{-6}$, $r^2 = 0.39$). Both genetic
363 distances measured using SNPs and rearrangements (inversions and translocations)
364 significantly correlated with viability when controlling for the other factor (Kendall
365 partial rank order correlations with viability SNPs|rearrangements $\tau = -0.19$, $P = 0.038$;
366 rearrangements|SNPs $\tau = -0.22$, $P = 0.016$). Some strains produce low viability spores
367 even when self-mated with their own genotype. The lowest self-mating viability of each
368 strain pair is indicated by circle size (see legend, smaller circles indicate lower self-
369 mating viability) to illustrate that low-viability outliers tend to include such cases (see
370 **Supplementary Table 8** for details).

371

372

373 **DISCUSSION**

374 Here we present the first genome- and population-wide catalog of SVs among *S. pombe*
375 strains. To account for the high discrepancy of available methods²⁵, we applied a
376 consensus approach to identify SVs (SURVIVOR), followed by rigorous filtering and
377 manual inspection of all calls. We focused on high specificity (the correctness of the
378 inferred SV) rather than high sensitivity (attempting to detect all SVs).

379 Our previous analyses of these strains, conducted without SV data¹⁷, attributed
380 both trait variations and reproductive isolation to SNPs and/or small indels. Here we
381 show that the small number of SVs we describe make substantial contributions to both of
382 these factors. We demonstrate that CNVs (duplications and deletions) contribute
383 significantly to our ability to describe quantitative traits, whereas variants that rearrange

384 the order of the genome (inversions and translocations) produce much weaker effects on
385 traits. In contrast, CNVs have no detectable influence on reproductive isolation, while
386 rearrangements contribute substantially to reproductive isolation, similar to other
387 species^{10,38}.

388 We show that CNVs and, to a lesser extent, rearrangements can produce
389 substantial contributions to trait variation. These CNVs subtly alter the expression of
390 genes within and beyond the duplications, and contribute considerably to quantitative
391 traits. Within small populations, CNVs may produce larger effects on traits in the short
392 term than SNPs, since their effect sizes can be substantial (SVs significant in GWAS
393 have a mean effect size of 16% in this study). Within budding yeast, clearly measured
394 effects of alterations to gene order in the DAL metabolic cluster³⁹ and the lethality of
395 some engineered rearrangements⁴⁰ indicates that rearrangements can also effect
396 phenotypic changes. Given the evidence for extensive ploidy and aneuploidy variation
397 with budding yeasts, including clinical and industrial budding yeasts^{29,41,42}, SVs can be
398 expected to have considerable impacts on phenotypic variation these fungi.

399 In this context, it is striking that CNVs appear to be transient within the clonal
400 populations that we studied. Our analysis is consistent with experimental studies with
401 budding yeast, indicating that both rearrangements and CNVs may be gained or lost at
402 rates in excess of point mutations. For example, frequent gain of duplications has been
403 observed in laboratory cultures of *S. pombe*, where spontaneous duplications suppress
404 *cdc2* mutants at least 100 times more frequently than point mutations. These suppressor
405 strains lose their duplications with equal frequency⁴³, indicating reversion of alleles.
406 Similarly, duplications frequently occur during experimental evolution with budding

407 yeast⁴⁴. This instability is likely facilitated by repeated elements, which are unstable
408 within both budding and fission yeast genomes⁴⁵⁻⁴⁸, which is also supported by the
409 enrichment of SVs in our population near retrotransposon LTRs (**Supplementary Figure**
410 **8**). Though we do not examine the stability of rearrangements, there is also evidence for
411 their instability. Transposon-mediated rearrangements are highly dynamic in laboratory
412 cultures during selection^{49,50}, and show elevated mutation rates at subtelomeric regions⁵¹.
413
414 This analysis also has relevance for human diseases, since *de novo* CNV formation in the
415 human genome occurs at a rate of approximately one CNV/10 generations⁵², and CNVs
416 are known to contribute to a wide variety of diseases⁴. Indeed, both the population
417 genetics and the effects of SVs within *S. pombe* seem similar to human, in that CNVs are
418 associated with stoichiometric changes on gene expression, and SVs are in weak linkage
419 with SNPs^{53,54}, and therefore may be badly tagged by SNPs in GWAS studies. We show
420 that CNVs and rearrangements in fission yeast not only rapidly emerge, but substantially
421 contribute to quantitative traits independent of weakly linked SNPs. These findings
422 highlight the need to identify SVs when describing traits using GWAS, and indicate that a
423 failure to call SVs can lead to an overestimation of the impact of SNPs to traits or
424 contribute to the problem that large proportions of the heritable component of trait
425 variation are not discovered in GWAS (the ‘missing heritability’). We observed a clear
426 example of this effect in two winemaking traits, where heritability was entirely due to
427 SVs.

428 In summary, we show that different types of SVs are transient within populations
429 of fission yeast, where they alter gene expression, impact phenotypes and can lead to

430 reproductive isolation.

431

432 **METHODS**

433 **Performance assessment of SV callers using simulated data**

434 To identify filtering parameters for DELLY, LUMPY and Pindel for the *S. pombe*
435 genome, we simulated seven datasets (s1-s7) of 40x coverage with a range of different
436 SV types and sizes (**Supplemental Table 7**). The simulated read sets contained
437 sequencing errors (0.4%), SNPs and indels (0.1%) within the range of actual data from *S.*
438 *pombe* strains and between 30 and 170 SVs. These data sets were produced by modifying
439 the reference genome using our in-house software (SURVIVOR, described below), and
440 simulating reads from this genome with Mason software⁵⁵.

441 After mapping the reads and calling SVs, we evaluated the calls. We defined a SV
442 correctly predicted if: i) the simulated and reported SV were of the same type (e.g.
443 duplication), ii) were predicted to be on same chromosome, and iii) their start and stop
444 locations were within 1 kb. We then defined caller-specific thresholds to optimize the
445 sensitivity and false discovery rate (FDR) for each caller. FDRs on the simulated data
446 were low: DELLY (average 0.13), LUMPY (average 0.06) and Pindel (average 0.04).

447 Selecting calls that were present in at least two callers further reduced the FDR
448 (average of 0.01). DELLY had the highest sensitivity (average 0.75), followed by
449 SURVIVOR (average 0.70), LUMPY (average 0.62) and Pindel (0.55). We further used
450 simulated data to assess the sensitivity and FDR of our predictions. cn.mops was
451 evaluated with a 2 kb distance for start and stop coordinates. Our cn.mops parameters
452 were designed to identify large (above 12 kb) events and thus did not identify any SVs

453 simulated for s1-s6. Details of simulations and caller efficacy are provided in

454 **Supplementary Table 9.**

455

456 **SURVIVOR (StructURal Variant majorIty VOte) Software Tool**

457 We developed the SURVIVOR tool kit for assessing SVs for short read data that contains

458 several modules. The first module simulates SVs given a reference genome file (fasta)

459 and the number and size ranges for each SV (insertions, deletions, duplications,

460 inversions and translocations). After reading in the reference genome, SURVIVOR

461 randomly selects the locations and size of SV following the provided parameters.

462 Subsequently, SURVIVOR alters the reference genome accordingly and prints the so

463 altered genome. In addition, SURVIVOR provides an extended bed file to report the

464 locations of the simulated SVs.

465 The second module evaluates SV calls based on a variant call format (VCF) file ⁵⁶

466 and any known list of SVs. A SV was identified as correct if i) they were of same type

467 (e.g. deletion); ii) they were reported on same chromosome, and iii) the start and stop

468 coordinates of the simulated and identified SV were within 1 kb (user definable).

469 The third module of SURVIVOR was used to filter and combine the calls from

470 three VCF files. In our case, these files were the results of DELLY, LUMPY and Pindel.

471 This module includes methods to convert the method-specific output formats to a VCF

472 format. SVs were filtered out if they were unique to one of the three VCF files. Two SVs

473 were defined as overlapping if they occur on the same chromosome, their start and stop

474 coordinates were within 1 kb, and they were of the same type. In the end, SURVIVOR

475 produced one VCF file containing the so filtered calls. SURVIVOR is available at
476 github.com/fritzsedlazeck/SURVIVOR.

477

478 **Read mapping and detection of structural variants**

479 Illumina paired-end sequencing data for 161 *S. pombe* strains were collected as described
480 in Jeffares, et al.¹⁷, with the addition of Leupold's reference 975 h^+ (JB32) and excluding
481 JB374 (known to be a gene-knockout version of the reference strain, see below).
482 Leupold's 968 h^{90} and Leupold's 972 h^- were included as JB50 and JB22, respectively
483 (**Supplementary Table 6**). For all strains, reads were mapped using NextGenMap
484 (version 0.4.12)⁵⁷ with the following parameter (-X 1000000) to the *S. pombe* reference
485 genome (version ASM294v2.22). Reads with 20 base pairs or more clipped were
486 extracted using the script *split_unmapped_to_fasta.pl* included in the LUMPY package
487 (version 0.2.9)²⁵ and were then mapped using YAHA (version 0.1.83)⁵⁸ to generate split-
488 read alignments. The two mapped files were merged using Picard-tools (version 1.105)
489 (<http://broadinstitute.github.io/picard>), and all strains were then down-sampled to 40x
490 coverage using Samtools (version 0.1.18)⁵⁹.

491 Subsequently, DELLY (version 0.5.9, parameters: “-q 20 -r”)²⁶, LUMPY
492 (version 0.2.9, recommended parameter settings)²⁵ and Pindel (version 0.2.5a8, default
493 parameter)²⁷ were used to independently identify SVs in the 161 strains using our
494 SURVIVOR software. This included merging any variants of the same type (duplication,
495 deletion *etc*) whose start and end coordinates were within 1 kb. Merging was justified by
496 the finding that most allele calls were close to the defined call (only 5% of start or end
497 positions were >300nt from the defined consensus boundary). We then retained all

498 variants predicted by at least two methods. These SVs calls were genotyped using
499 DELLY.

500 To identify further CNVs, we ran cn.MOPS²⁴ with parameters tuned to collect
501 large duplications/deletions as follows: read counts were collected from bam alignment
502 files (as above) with *getReadCountsFromBAM* and WL=2000, and CNVs predicted using
503 *haplocn.mops* with minWidth= 6, all other parameters as default. Hence, the minimum
504 variant size detected was 12 kb. CNV were predicted for each strain independently by
505 comparing the alternative strain to the two reference strains (JB22, JB32) and four
506 reference-like strains that differed from the reference by less than 200 SNPs (JB1179,
507 JB1168, JB937, JB936).

508 After CNV calling, allele calling was achieved by comparing counts of coverage
509 in 100bp windows for the two reference strains (JB22, JB32) to each alternate strain
510 using custom R scripts. Alleles were called as non-reference duplications if the one-sided
511 Wilcoxon rank sum test p-values for both JB22 and JB32 vs alternate strain were less
512 than 1×10^{-10} (showing a difference in coverage) and the ratio of alternate/reference
513 coverage (for both JB22 and JB32) was >1.8 (duplications), or <0.2 (deletions). Manual
514 inspection of coverage plots showed that the vast majority of the allele calls were in
515 accordance with what we discerned by eye. These R scripts were also used to examine
516 CNVs predicted to be segregating within clusters (clonal populations). All such CNVs
517 were examined in all clusters that contained at least one non-reference allele call
518 **(Supplementary Table 10).**

519 Finally, we manually mapped two large duplications that did not satisfy these
520 criteria (DUP.I:2950001..3190000, 240kb and DUP.I:5050001..5560000, 510kb – both

521 singletons in JB1207), but were clearly visible in chromosome-scale read coverage plots
522 **(Supplementary Figure 9).**

523

524 **Reduction of false discovery rate**

525 This filtering produced 315 variant calls. However, because 31 of these 315 (~10%) were
526 called within the two reference strains (JB22, JB32), we expected that this set still
527 contained false positives. To further reduce the false positive rate, we looked for
528 parameters that would reduce calls made in reference strains (JB22 and JB32) but not
529 reduce calls in strains more distantly related to the reference (JB1177, JB916 and JB894
530 that have 68223, 60087 and 67860 SNP differences to reference¹⁷). The reasoning was
531 that we expected to locate few variants in the reference, and more variants in the more
532 distantly related strains. This analysis showed that paired end support, repeats and
533 mapping quality were of primary value.

534 We therefore discarded all SVs that had a paired end support of 10 or less. In
535 addition, we ignored SVs that appeared in low mapping quality regions (i.e. regions
536 where reads with MQ=0 map) or those where both start and end coordinates overlapped
537 with previously identified retrotransposon LTRs¹⁷.

538 Finally, to ensure a high specificity call set, these filtered SVs were manually
539 curated using IGV⁶⁰ **(Supplementary Tables 11,12)**. We assigned each SVs a score (0:
540 not reliable, 1: unclear, 2: reliable based on inspection of alignments through IGV). We
541 utilized different visualizations from IGV to identify regions were pairs of the reads
542 mapped to different loci, for example, which we interpreted as possible artefacts. Overall,
543 we investigated whether the alignments of the breakpoints and reads in close proximity

544 had a reliable mapping in terms of mapping quality and clearness of the distortions of the
545 pairs. Only calls passing this manual curation as reliable (score 2) were included in the
546 final data set of 113 variants utilized for all further analyses. These filtering and manual
547 curation steps reduced our variant calls substantially, from 315 to 113. At this stage only
548 1/113 (~1%) of these variants was called within the two standard reference strains
549 (Leupolds's *h+* and *h-*, JB22 and JB32 in our collection).

550

551 **PCR validation**

552 PCR analysis was performed to confirm 10 of the 11 inversions and all 15 translocations
553 from the curated data set. One inversion was too small to examine by PCR
554 (INV.AB325691:6644..6784, 140 nt). Primers were designed using Primer3⁶¹ to amplify
555 both the reference and alternate alleles. PCR was carried out with each primer set using a
556 selection of strains that our genotype calls predict to include at least one alternate allele
557 and at least one reference allele (usually 6 strains). Products were scored according to
558 product size and presence/absence (**Supplementary Tables 13,14**).

559 Inversions: 9/10 variants were at least partially verified by either reference or
560 alternate allele PCR (3 variants were verified by both reference and alternate PCRs), and
561 7/10 inversions also received support from BLAST (see below). Translocations: 10/15
562 were at least partially verified by either reference or alternate allele PCR (5/15 variants
563 were verified by both reference and alternate PCRs). One additional translocation
564 received support from BLAST (see below), meaning that 11/15 translocations were
565 supported by PCR and/or BLAST. Three of the four translocations that could not be
566 verified were probably nuclear copies of mitochondrial genes (NUMTs)⁶², because one

567 breakpoint was mapped to the mitochondrial genome. Details of the 113 curated variants
568 are presented in Supplementary Table 15.

569

570 **Validation by BLAST of *de novo* assemblies**

571 We further assessed the quality of the predicted breakpoints for the inversions and
572 translocations by comparing them to the previously created *de novo* assemblies for each
573 of the 161 strains¹⁷. To this end, we created blast databases for the scaffolds of each
574 strain that were >1kb. We then created the predicted sequence for 1 kb around each
575 junction of the validated 10 inversions and 15 translocations. These sequences were used
576 to search the blast databases using BLAST+ with --gapopen 1 --gapextend 1 parameters.
577 We accepted any blast hsp with a length >800 bp as supporting the junction (because
578 these must contain at least 300 bp at each side of the break point). Four inversions and
579 three translocations gained support from these searches (Supplementary File Tables2-
580 PCR.xlsx).

581

582 **Knockout strain control**

583 Our sample of sequenced strains included one strain (JB374) that is known to contain
584 deletions of the *his3* and *ura4* genes. Our variant calling and validation methods
585 identified only two variants in this strain, both deletions that corresponded to the
586 positions of these genes, as below:
587 *his3* gene location is chromosome II, 1489773-1488036, deletion detected at II:1488228-
588 1489646.

589 *ura4* gene location is chromosome III, 115589-116726, deletion detected at III:115342-
590 117145.

591 This strain was not included in the further analyses of the SVs.

592

593 **Microarray expression analysis**

594 Cells were grown in YES (Formedium, UK) and harvested at $OD_{600} = 0,5$. RNA was
595 isolated followed by cDNA labeling⁶³. Agilent 8 x 15K custom-made *S. pombe*
596 expression microarrays were used. Hybridization, normalization and subsequent washes
597 were performed according to the manufacturer's protocols. The obtained data were
598 scanned and extracted using GenePix and processed for quality control and normalization
599 using in-house developed R scripts. Subsequent analysis of normalized data was
600 performed using R. Microarray data have been submitted to ArrayExpress (accession
601 number E-MTAB-4019). Genes were considered as induced if their expression signal
602 after normalization was >1.9 , and repressed if <0.51 .

603

604 **Time to most recent common ancestor (TMRCA) estimates**

605 Previously, based on the genetic distances between these strains and the 'dated tip' dating
606 method implemented in BEAST⁶⁴, we have estimated the divergence times between all
607 161 *S. pombe* strains sequenced¹⁷. To determine the TMRCA for pairs of strains, we re-
608 examined the BEAST outputs using FigTree to obtain the medium and 95% confidence
609 intervals.

610

611 **SNP and indel calling**

612 SNPs were called as described¹⁷. Insertions and deletions (indels) were called in 160
613 strains using stampy-mapped, indel-realigned bams as described previously¹⁷. We
614 accepted indels that were called by both the Genome Analysis Toolkit HaplotypeCaller⁶⁵
615 and Freebayes⁶⁶, and then genotyped all these calls with Freebayes.

616 Briefly, indels were called on each strains bam with HaplotypeCaller, and filtered
617 for call quality >30 and mapping quality >30 (bcftools filter --include 'QUAL>30 &&
618 MQ>30'). Separately, indels were called on each strains bam with Freebayes, and filtered
619 for call quality >30. All Freebayes vcf files were merged, accepting only positions called
620 by both Freebayes and HaplotypeCaller. These indels were then genotyped with
621 Freebayes using a merged bam (containing reads from all strains), using the --variant-
622 input flag for Freebayes to genotyped only the union calls. Finally indels were filtered for
623 by score, mean reference mapping quality and mean alternate mapping quality >30
624 (bcftools filter --include 'QUAL>30 && MQM>30 & MQMR>30'). These methods
625 identified 32,268 indels. Only 50 of these segregated between Leupold's h⁻ reference
626 (JB22) and Leupold's h⁹⁰ reference (JB50), whereas 12109 indels segregated between the
627 JB22 reference and the divergent strain JB916.

628

629 **Heredity and GWAS**

630 We selected 53 traits that contained at least values from 100 strains¹⁷, and so included
631 multiple individuals from within clonal populations (growth rates on 42 different solid
632 media and 11 cell shape characters measured with automated image analysis). Trait
633 values were normalized using a rank-based transformation in R, for each trait vector y ,
634 $\text{normal.y} = \text{qnorm}(\text{rank}(y)/(1+\text{length}(y)))$. Total heritability, and the contribution of SNPs,

635 CNVs and rearrangements were estimated using LDAK (version 5.94)⁶⁷, with kinship
636 matrices derived from all SNPs, 146 CNVs, and 15 rearrangements. All genotypes,
637 including CNVs were encoded as binary values (1 or 0) for heritability and GWAS. To
638 assess whether the contribution of CNVs could be primarily due to linkage with causal
639 SNPs, we simulated trait data using the --make-phenos function of LDAK with the
640 relatedness matrix from all SNPs, assuming that all variants contributed to the trait (--
641 num-causals -1). We made one simulated trait data set per trait, for each of the 53 traits,
642 with total heritability defined as predicted from the real data. We then estimated the
643 heritability using LDAK, including the joint matrix of SNPs, CNVs and rearrangements.
644 To assess the extent to which the contribution of SNPs to heritability was overestimated,
645 we performed another simulation using the relatedness matrix from the 87 segregating
646 CNVs alone, and then estimated the contribution of SNPs, CNVs and rearrangements in
647 this simulated data as above.

648 Genome-wide associations were performed with LDAK (version 5) using default
649 parameters. To account for the unequal relatedness of strains, we used a kinship matrix
650 derived from all 172,368 SNPs called previously Jeffares, et al.¹⁷. Association analysis
651 was used to find associations between traits, testing SVs, SNPs and indels with a minor
652 allele count ≥ 5 . Analysis was run separately for 68 SVs, 139,396 SNPs and 22,058 indels
653 (each used the kinship derived from all SNPs). We examined the same 53 traits as for the
654 heritability analysis (above). For each trait, we carried out 1000 permutations of trait
655 data, and define the 5th percentile of these permutations as the trait-specific P-value
656 threshold.
657

658 **Model details for Heritability and GWAS Analysis**

659 To estimate the heritability contribution of SNPs, we computed a kinship matrix (K_{SNP})
660 using all 172,368 SNPs that we had discovered in our previous published analysis¹⁷
661 (elements of this matrix represent pairwise allelic correlations across all SNPs)⁶⁷, onto
662 which we regressed the phenotypic values assuming the following model:

$$Y \sim N(0, K_{\text{SNP}} \sigma_{\text{SNP}}^2 + \sigma_e^2 I)$$

663

664 We estimated the two variance components, σ_{SNP}^2 and σ_e^2 , using REML (restricted
665 maximum likelihood), based on which our estimates of the heritability of SNPs is

$$\frac{\sigma_{\text{SNP}}^2}{\sigma_{\text{SNP}}^2 + \sigma_e^2}$$

666

667 To estimate the heritability of CNVs and rearrangements, we repeated this analysis using
668 instead K_{CNV} then K_{REA} , computed using only 146 segregating CNVs and 15 segregating
669 rearrangements, respectively.

670

671 We additionally considered the model

$$672 Y \sim N(0, K_{\text{SNP}} \sigma_{\text{SNP}}^2 + K_{\text{CNV}} \sigma_{\text{CNV}}^2 + K_{\text{REA}} \sigma_{\text{REA}}^2 + \sigma_e^2 I),$$

673

674 Having estimated the four variance components, again using REML, the relative
675 contributions of SNPs, CNVs and rearrangements are, respectively,

$$676 \frac{\sigma_{\text{SNP}}^2}{S}, \frac{\sigma_{\text{CNV}}^2}{S} \text{ and } \frac{\sigma_{\text{REA}}^2}{S}$$

677 where $S = \sigma_{\text{SNP}}^2 + \sigma_{\text{CNV}}^2 + \sigma_{\text{REA}}^2$

678

679 To test the specificity of this analysis, we generated phenotypes for which only one
680 predictor type contributed (e.g., only SNPs), then analyzed using the individual and joint
681 models above, which allowed us to assess how accurately we can distinguish between
682 contributions of different predictor types.

683 For the mixed model association analysis, we used the same the SNP kinship
684 matrix. As the predictors (variants that we examined for effects on a trait), we chose to
685 analyse SNPs, indels and SVs with a minor allele count ≥ 5 (68 SVs, 139396 SNPs and
686 22,058 indels).

687 Then for each predictor X_j we considered the model

$$688 Y \sim N(\beta_j X_j K_{\text{SNP}} \sigma_{\text{SNP}}^2 + \sigma_e^2 I),$$

689 where β_j is the effect size of predictor X_j

690

691 Having solved using REML, we used a likelihood ratio test (comparing to the null model
692 ($\beta_j = 0$) to assess whether β_j is significantly non-zero. Each of these analyses used the
693 kinship derived from all SNPs.

694

695

696 **Offspring viability and genetic distance**

697 Cross spore viability data and self-mating viability were collected from previous analyses
698 ^{6,17}. The number of differences between each pair was calculated using vcftools vcf-
699 subset ⁵⁶, and correlations were estimated using R, with the ppcor package. When
700 calculating the number of CNVs differences between strains, we altered our criteria for

701 ‘different’ variants (to merge variants whose starts and ends were within 1 kb), and
702 merged CNVs if their overlap was >50% and their allele calls were the same.

703

704 **Transience analysis**

705 For each CNV, we extracted all SNPs from 20 kb upstream and 20 kb downstream. 86/87
706 CNVs showed variation in these regions (DUP.MT:1..19382 was the only CNV with no
707 corresponding SNPs). We then used these concatenated SNPs to build a local SNP-based
708 tree with FastTree (version 2.1.9)⁶⁸. To build a CNV-based tree from the copy number
709 variation in each CNV region, we used a neighbour-joining tree estimation based on the
710 Euclidean distances between strains.

711 The total branch length of the CNV-based tree was strongly correlated (Spearman
712 rank correlation $\rho=0.90$, $P < 0.001$) with the standard deviation of copy number variation
713 (Supplementary Figure 4). We therefore used this standard deviation to define a relative
714 rate of transience for each cluster, $\sigma_{rc} = \sigma_{ic}/\sigma_{oc}$ where σ_{ic} and σ_{oc} are the within cluster and
715 without cluster standard deviations respectively, meaning that CNVs which were highly
716 relatively transient within a given cluster would have high values of σ_{rc} . This was used to
717 select the three CNVs visualised in Figure 2c. See Supplementary Table 2 for all values
718 of σ_{rc} , Supplementary Figure 4 for visualization as heatmap). Visualisations of all 86/87
719 CNVs with their SNP-based phylogenies are available at:

720 https://figshare.com/projects/fission_yeast_structural_variation/15798.

721 Circle plots were used to visualize the variation in copy number over the SNP-based
722 phylogeny for each CNV using Anvi'o (version 2.0.3)⁶⁹.

723

724

725

726 References

- 727 1 Chen, C. *et al.* A comprehensive survey of copy number variation in 18 diverse pig populations and
728 identification of candidate copy number variable genes associated with complex traits. *BMC Genomics*
729 **13**, 733, doi:10.1186/1471-2164-13-733 (2012).
- 730 2 Stefansson, H. *et al.* CNVs conferring risk of autism or schizophrenia affect cognition in controls. *Nature*
731 **505**, 361-366, doi:10.1038/nature12818 (2014).
- 732 3 Wang, Y. *et al.* Copy number variation at the GL7 locus contributes to grain size diversity in rice. *Nat*
733 *Genet* **47**, 944-948, doi:10.1038/ng.3346 (2015).
- 734 4 Zhang, F., Gu, W., Hurler, M. E. & Lupski, J. R. Copy number variation in human health, disease, and
735 evolution. *Annu Rev Genomics Hum Genet* **10**, 451-481, doi:10.1146/annurev.genom.9.081307.164217
736 (2009).
- 737 5 Zhang, H. *et al.* Gene copy-number variation in haploid and diploid strains of the yeast *Saccharomyces*
738 *cerevisiae*. *Genetics* **193**, 785-801, doi:10.1534/genetics.112.146522 (2013).
- 739 6 Avelar, A. T., Perfeito, L., Gordo, I. & Ferreira, M. G. Genome architecture is a selectable trait that can be
740 maintained by antagonistic pleiotropy. *Nat Commun* **4**, 2235, doi:10.1038/ncomms3235 (2013).
- 741 7 Brown, W. R. *et al.* A Geographically Diverse Collection of *Schizosaccharomyces pombe* Isolates Shows
742 Limited Phenotypic Variation but Extensive Karyotypic Diversity. *G3 (Bethesda)* **1**, 615-626,
743 doi:10.1534/g3.111.001123 (2011).
- 744 8 McGaugh, S. E. & Noor, M. A. Genomic impacts of chromosomal inversions in parapatric *Drosophila*
745 species. *Philos Trans R Soc Lond B Biol Sci* **367**, 422-429, doi:10.1098/rstb.2011.0250 (2012).
- 746 9 Rieseberg, L. H. Chromosomal rearrangements and speciation. *Trends Ecol Evol* **16**, 351-358 (2001).
- 747 10 Ayala, D., Guerrero, R. F. & Kirkpatrick, M. Reproductive isolation and local adaptation quantified for a
748 chromosome inversion in a malaria mosquito. *Evolution* **67**, 946-958, doi:10.1111/j.1558-
749 5646.2012.01836.x (2013).
- 750 11 Stevison, L. S., Hoehn, K. B. & Noor, M. A. Effects of inversions on within- and between-species
751 recombination and divergence. *Genome Biol Evol* **3**, 830-841, doi:10.1093/gbe/evr081 (2011).
- 752 12 Liti, G., Barton, D. B. & Louis, E. J. Sequence diversity, reproductive isolation and species concepts in
753 *Saccharomyces*. *Genetics* **174**, 839-850, doi:10.1534/genetics.106.062166 (2006).
- 754 13 Hou, J., Friedrich, A., de Montigny, J. & Schacherer, J. Chromosomal rearrangements as a major
755 mechanism in the onset of reproductive isolation in *Saccharomyces cerevisiae*. *Curr Biol* **24**, 1153-
756 1159, doi:10.1016/j.cub.2014.03.063 (2014).
- 757 14 Leducq, J. B. *et al.* Speciation driven by hybridization and chromosomal plasticity in a wild yeast. *Nat*
758 *Microbiol* **1**, 15003, doi:10.1038/nmicrobiol.2015.3 (2016).
- 759 15 Ortiz-Barrientos, D., Engelstadter, J. & Rieseberg, L. H. Recombination Rate Evolution and the Origin of
760 Species. *Trends Ecol Evol* **31**, 226-236, doi:10.1016/j.tree.2015.12.016 (2016).
- 761 16 Fawcett, J. A. *et al.* Population genomics of the fission yeast *Schizosaccharomyces pombe*. *PLoS One* **9**,
762 e104241, doi:10.1371/journal.pone.0104241 (2014).
- 763 17 Jeffares, D. C. *et al.* The genomic and phenotypic diversity of *Schizosaccharomyces pombe*. *Nat Genet*
764 **47**, 235-241, doi:10.1038/ng.3215 (2015).
- 765 18 Zanders, S. E. *et al.* Genome rearrangements and pervasive meiotic drive cause hybrid infertility in
766 fission yeast. *Elife* **3**, e02630, doi:10.7554/eLife.02630 (2014).
- 767 19 Wood, V. *et al.* The genome sequence of *Schizosaccharomyces pombe*. *Nature* **415**, 871-880,
768 doi:10.1038/nature724 (2002).
- 769 20 Sabatinos, S. A. & Forsburg, S. L. Molecular genetics of *Schizosaccharomyces pombe*. *Methods Enzymol*
770 **470**, 759-795, doi:10.1016/S0076-6879(10)70032-X (2010).
- 771 21 Kim, D. U. *et al.* Analysis of a genome-wide set of gene deletions in the fission yeast
772 *Schizosaccharomyces pombe*. *Nat Biotechnol* **28**, 617-623, doi:10.1038/nbt.1628 (2010).
- 773 22 Marguerat, S. *et al.* Quantitative analysis of fission yeast transcriptomes and proteomes in proliferating
774 and quiescent cells. *Cell* **151**, 671-683, doi:10.1016/j.cell.2012.09.019 (2012).
- 775 23 Ryan, C. J. *et al.* Hierarchical modularity and the evolution of genetic interactomes across species. *Mol*
776 *Cell* **46**, 691-704, doi:10.1016/j.molcel.2012.05.028 (2012).
- 777 24 Klambauer, G. *et al.* cn.MOPS: mixture of Poissons for discovering copy number variations in next-
778 generation sequencing data with a low false discovery rate. *Nucleic Acids Res* **40**, e69,
779 doi:10.1093/nar/gks003 (2012).
- 780 25 Layer, R. M., Chiang, C., Quinlan, A. R. & Hall, I. M. LUMPY: a probabilistic framework for structural
781 variant discovery. *Genome Biol* **15**, R84, doi:10.1186/gb-2014-15-6-r84 (2014).
- 782 26 Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end and split-read analysis.
783 *Bioinformatics* **28**, i333-i339, doi:10.1093/bioinformatics/bts378 (2012).

- 784 27 Ye, K., Schulz, M. H., Long, Q., Apweiler, R. & Ning, Z. Pindel: a pattern growth approach to detect break
785 points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics* **25**,
786 2865-2871, doi:10.1093/bioinformatics/btp394 (2009).
- 787 28 Bergstrom, A. *et al.* A high-definition view of functional genetic variation from natural yeast genomes.
788 *Mol Biol Evol* **31**, 872-888, doi:10.1093/molbev/msu037 (2014).
- 789 29 Dunn, B., Richter, C., Kvitek, D. J., Pugh, T. & Sherlock, G. Analysis of the *Saccharomyces cerevisiae* pan-
790 genome reveals a pool of copy number variants distributed in diverse yeast strains from differing
791 industrial environments. *Genome Res* **22**, 908-924, doi:10.1101/gr.130310.111 (2012).
- 792 30 Bitton, D. A. *et al.* AnGeLi: A Tool for the Analysis of Gene Lists from Fission Yeast. *Front Genet* **6**, 330,
793 doi:10.3389/fgene.2015.00330 (2015).
- 794 31 Chikashige, Y. *et al.* Gene expression and distribution of *Swi6* in partial aneuploids of the fission yeast
795 *Schizosaccharomyces pombe*. *Cell Struct Funct* **32**, 149-161 (2007).
- 796 32 Benito, A. *et al.* Selected *Schizosaccharomyces pombe* Strains Have Characteristics That Are Beneficial
797 for Winemaking. *PLoS One* **11**, e0151102, doi:10.1371/journal.pone.0151102 (2016).
- 798 33 Dobzhansky, T. On the Sterility of the Interracial Hybrids in *Drosophila Pseudoobscura*. *Proc Natl Acad
799 Sci U S A* **19**, 397-403 (1933).
- 800 34 Muller, H. J. Reversibility in Evolution Considered from the Standpoint of Genetics. *Biological Reviews*
801 **14**, 261-280, doi:10.1111/j.1469-185X.1939.tb00934.x (1939).
- 802 35 Fischer, G., Rocha, E. P., Brunet, F., Vergassola, M. & Dujon, B. Highly variable rates of genome
803 rearrangements between hemiascomycetous yeast lineages. *PLoS Genet* **2**, e32,
804 doi:10.1371/journal.pgen.0020032 (2006).
- 805 36 Gordon, J. L., Byrne, K. P. & Wolfe, K. H. Additions, losses, and rearrangements on the evolutionary route
806 from a reconstructed ancestor to the modern *Saccharomyces cerevisiae* genome. *PLoS Genet* **5**,
807 e1000485, doi:10.1371/journal.pgen.1000485 (2009).
- 808 37 Delneri, D. *et al.* Engineering evolution to study speciation in yeasts. *Nature* **422**, 68-72,
809 doi:10.1038/nature01418 (2003).
- 810 38 Noor, M. A., Grams, K. L., Bertucci, L. A. & Reiland, J. Chromosomal inversions and the reproductive
811 isolation of species. *Proc Natl Acad Sci U S A* **98**, 12084-12088, doi:10.1073/pnas.221274498 (2001).
- 812 39 Naseeb, S. & Delneri, D. Impact of chromosomal inversions on the yeast DAL cluster. *PLoS One* **7**,
813 e42022, doi:10.1371/journal.pone.0042022 (2012).
- 814 40 Naseeb, S. *et al.* Widespread Impact of Chromosomal Inversions on Gene Expression Uncovers
815 Robustness via Phenotypic Buffering. *Mol Biol Evol* **33**, 1679-1696, doi:10.1093/molbev/msw045
816 (2016).
- 817 41 Zhu, Y. O., Sherlock, G. & Petrov, D. A. Whole Genome Analysis of 132 Clinical *Saccharomyces cerevisiae*
818 Strains Reveals Extensive Ploidy Variation. *G3 (Bethesda)* **6**, 2421-2434, doi:10.1534/g3.116.029397
819 (2016).
- 820 42 Gallone, B. *et al.* Domestication and Divergence of *Saccharomyces cerevisiae* Beer Yeasts. *Cell* **166**,
821 1397-1410 e1316, doi:10.1016/j.cell.2016.08.020 (2016).
- 822 43 Carr, A. M., MacNeill, S. A., Hayles, J. & Nurse, P. Molecular cloning and sequence analysis of mutant
823 alleles of the fission yeast *cdc2* protein kinase gene: implications for *cdc2+* protein structure and
824 function. *Mol Gen Genet* **218**, 41-49 (1989).
- 825 44 Dunham, M. J. *et al.* Characteristic genome rearrangements in experimental evolution of *Saccharomyces
826 cerevisiae*. *Proc Natl Acad Sci U S A* **99**, 16144-16149, doi:10.1073/pnas.242624799 (2002).
- 827 45 Chan, J. E. & Kolodner, R. D. A genetic and structural study of genome rearrangements mediated by high
828 copy repeat Ty1 elements. *PLoS Genet* **7**, e1002089, doi:10.1371/journal.pgen.1002089 (2011).
- 829 46 Coulon, S. *et al.* Slx1-Slx4 are subunits of a structure-specific endonuclease that maintains ribosomal
830 DNA in fission yeast. *Mol Biol Cell* **15**, 71-80, doi:10.1091/mbc.E03-08-0586 (2004).
- 831 47 Gadaleta, M. C. *et al.* Swi1Timeless Prevents Repeat Instability at Fission Yeast Telomeres. *PLoS Genet*
832 **12**, e1005943, doi:10.1371/journal.pgen.1005943 (2016).
- 833 48 Vincens, M. D., Legendre, M., Caldara, M., Hagihara, M. & Verstrepen, K. J. Unstable tandem repeats in
834 promoters confer transcriptional evolvability. *Science* **324**, 1213-1216, doi:10.1126/science.1170097
835 (2009).
- 836 49 Chang, S. L., Lai, H. Y., Tung, S. Y. & Leu, J. Y. Dynamic large-scale chromosomal rearrangements fuel
837 rapid adaptation in yeast populations. *PLoS Genet* **9**, e1003232, doi:10.1371/journal.pgen.1003232
838 (2013).
- 839 50 Gresham, D. *et al.* The repertoire and dynamics of evolutionary adaptations to controlled nutrient-
840 limited environments in yeast. *PLoS Genet* **4**, e1000303, doi:10.1371/journal.pgen.1000303 (2008).
- 841 51 Nishant, K. T. *et al.* The baker's yeast diploid genome is remarkably stable in vegetative growth and
842 meiosis. *PLoS Genet* **6**, e1001109, doi:10.1371/journal.pgen.1001109 (2010).
- 843 52 Itsara, A. *et al.* De novo rates and selection of large copy number variation. *Genome Res* **20**, 1469-1481,
844 doi:10.1101/gr.107680.110 (2010).

- 845 53 Stranger, B. E. *et al.* Relative impact of nucleotide and copy number variation on gene expression
846 phenotypes. *Science* **315**, 848-853, doi:10.1126/science.1136678 (2007).
847 54 Sudmant, P. H. *et al.* An integrated map of structural variation in 2,504 human genomes. *Nature* **526**,
848 75-81, doi:10.1038/nature15394 (2015).
849 55 Holtgrewe, M. Mason-A Read Simulator for Second Generation Sequencing Data. (Institut für
850 Mathematik und Informatik, Freie Universität Berlin, 2010).
851 56 Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156-2158,
852 doi:10.1093/bioinformatics/btr330 (2011).
853 57 Sedlazeck, F. J., Rescheneder, P. & von Haeseler, A. NextGenMap: fast and accurate read mapping in
854 highly polymorphic genomes. *Bioinformatics* **29**, 2790-2791, doi:10.1093/bioinformatics/btt468
855 (2013).
856 58 Faust, G. G. & Hall, I. M. YAHA: fast and flexible long-read alignment with optimal breakpoint detection.
857 *Bioinformatics* **28**, 2417-2424, doi:10.1093/bioinformatics/bts456 (2012).
858 59 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079,
859 doi:10.1093/bioinformatics/btp352 (2009).
860 60 Thorvaldsdottir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV): high-
861 performance genomics data visualization and exploration. *Brief Bioinform* **14**, 178-192,
862 doi:10.1093/bib/bbs017 (2013).
863 61 Untergasser, A. *et al.* Primer3--new capabilities and interfaces. *Nucleic Acids Res* **40**, e115,
864 doi:10.1093/nar/gks596 (2012).
865 62 Lenglez, S., Hermand, D. & Decottignies, A. Genome-wide mapping of nuclear mitochondrial DNA
866 sequences links DNA replication origins to chromosomal double-strand break formation in
867 *Schizosaccharomyces pombe*. *Genome Res* **20**, 1250-1261, doi:10.1101/gr.104513.109 (2010).
868 63 Lyne, R. *et al.* Whole-genome microarrays of fission yeast: characteristics, accuracy, reproducibility,
869 and processing of array data. *BMC Genomics* **4**, 27, doi:10.1186/1471-2164-4-27 (2003).
870 64 Drummond, A. J., Suchard, M. A., Xie, D. & Rambaut, A. Bayesian phylogenetics with BEAUti and the
871 BEAST 1.7. *Mol Biol Evol* **29**, 1969-1973, doi:10.1093/molbev/mss075 (2012).
872 65 DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA
873 sequencing data. *Nat Genet* **43**, 491-498, doi:10.1038/ng.806 (2011).
874 66 Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. *arXiv preprint*
875 *arXiv 1207.3907* (2012).
876 67 Speed, D., Hemani, G., Johnson, M. R. & Balding, D. J. Improved heritability estimation from genome-
877 wide SNPs. *Am J Hum Genet* **91**, 1011-1021, doi:10.1016/j.ajhg.2012.10.010 (2012).
878 68 Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood trees for large
879 alignments. *PLoS One* **5**, e9490, doi:10.1371/journal.pone.0009490 (2010).
880 69 Eren, A. M. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* **3**, e1319,
881 doi:10.7717/peerj.1319 (2015).
882
883

884

885 **Supplementary data**

886 SNP, indel and SVs calls, genotypes and copy numbers are available on figshare at:

887 https://figshare.com/projects/fission_yeast_structural_variation/15798

888 Array data is available at ArrayExpress, accession number: E-MTAB-4019.

889

890 **Acknowledgments**

891 We thank Günter Klambauer for advice on cn.MOPS and Michael C. Schatz for helpful
892 discussions and comments on the manuscript. F.S. was supported through National
893 Science Foundation awards (DBI-1350041) and National Institutes of Health award
894 (R01-HG006677). D.J., M.H., C.R. were supported by a Wellcome Trust Senior
895 Investigator Award to J.B. (grant 095598/Z/11/Z). J.B. was supported by a Royal Society
896 Wolfson Research Merit Award.

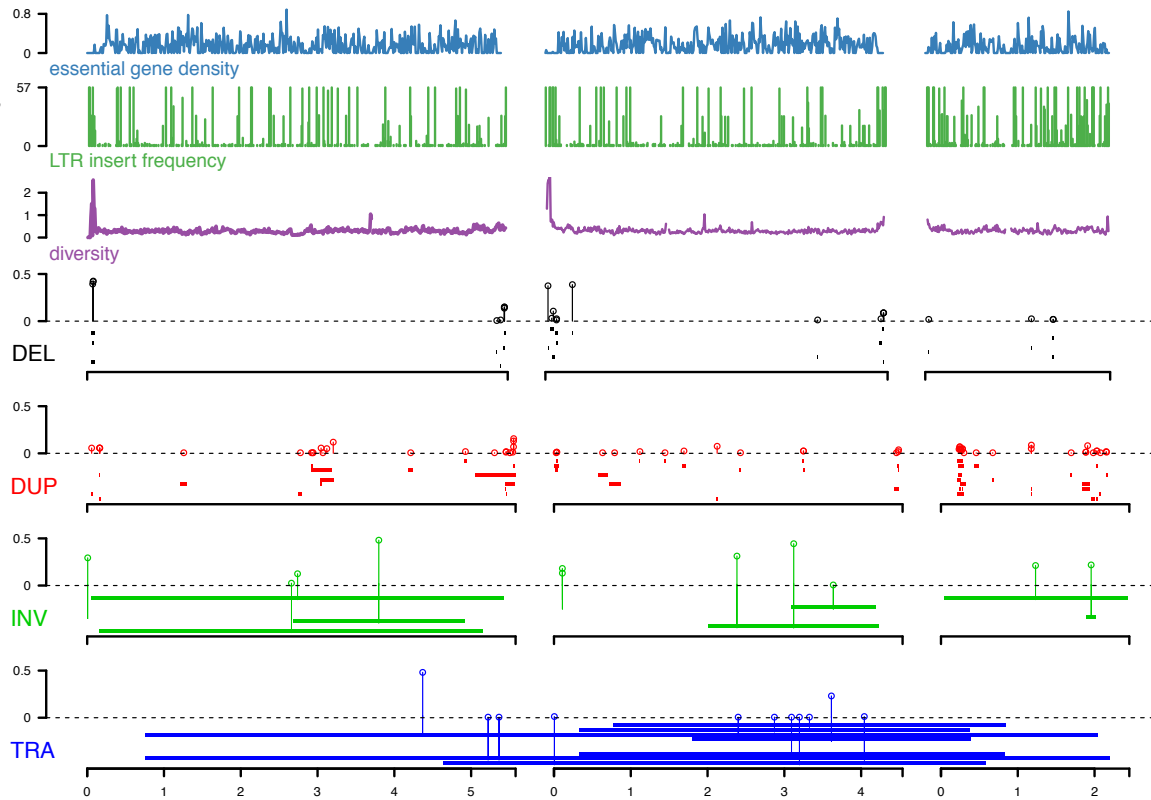
897

898 **Author contributions**

899 DJ, FS, CD and JB conceived and developed the study. DJ, LS, CJ and FS conducted the
900 bioinformatics analysis. DJ designed the laboratory work. FS designed and implemented
901 SURVIVOR. DS contributed to analysis of heritability and GWAS. CR and MH
902 produced the expression array analysis. MH conducted PCR validation of variants. JB
903 provided the bulk of funding for personnel and research costs. DJ, FS, CJ, LS, FB, CD
904 and JB wrote the manuscript.

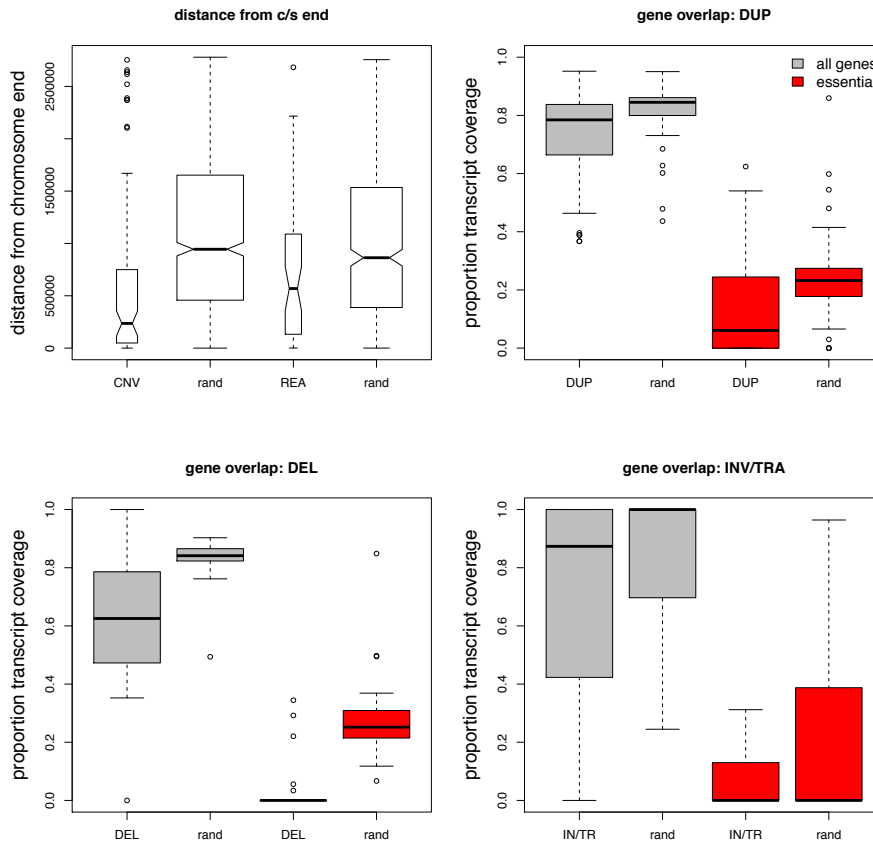
905

906 **Supplementary Figures**
907



908
909

910 **Supplementary Figure 1. Locations and minor allele frequencies of all structural**
911 **variants in curated data set.** Each of the three chromosomes is indicated by black bar,
912 with scale (in megabases) at bottom. From top (same data as Fig 1): density of essential
913 genes (blue), locations of *Tf*-type retrotransposons (green), and diversity (π , average
914 pairwise diversity from SNPs, purple). Bar heights for deletions and duplications are
915 proportional to minor allele frequency, the scale for retrotransposons is the frequency of
916 the insertion in the 57 non-clonal strains. Diversity and retrotransposon were calculated
917 from 57 non-clonal strains as described in Jeffares, et al. ¹⁷. Below, we show different
918 types of SVs: deletions (black), duplications (red), inversions (green) and translocations
919 (blue). The vertical lines terminating with open circles above dotted lines emit from the
920 mid-point of each SV and indicate the minor allele frequencies in the population of 161
921 strains.
922



923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

Supplementary Figure 2. Structural variations are biased towards chromosome

ends and to low gene density regions. Top left panel, both CNVs and rearrangements

are biased towards the ends of chromosomes. CNVs; median distance to chromosome

ends 236 kb compared to chromosome- and size-matched random sites 944 kb, Wilcoxon

rank sum test $P = 1.3 \times 10^{-11}$, rearrangements median distance 569 kb vs, matched random

863 kb, Wilcoxon test $P = 0.03$). All other panels, calculated proportion of each

duplication and deletion that contained all protein-coding or essential genes. Box plots

show the distributions of these proportions for all genes (grey), and proportion of

coverage by essential genes (red), compared to the null distribution (rand). All

comparisons were significantly less than the null distributions (Wilcoxon rank sum test,

P -values $< 1.6 \times 10^{-4}$). The same analysis was performed with the junctions of inversions

and translocations, by calculating the transcript coverage in the region 500 bp up- and

down-stream of the predicted start and end junctions. These rearrangements are slightly

biased away from genes ($P = 1.9 \times 10^{-3}$), but not significantly biased away from essential

genes ($P > 0.05$). The null distributions were determined by selecting 10 regions for each

actual variant/junction that were the same size, and were placed in random positions on

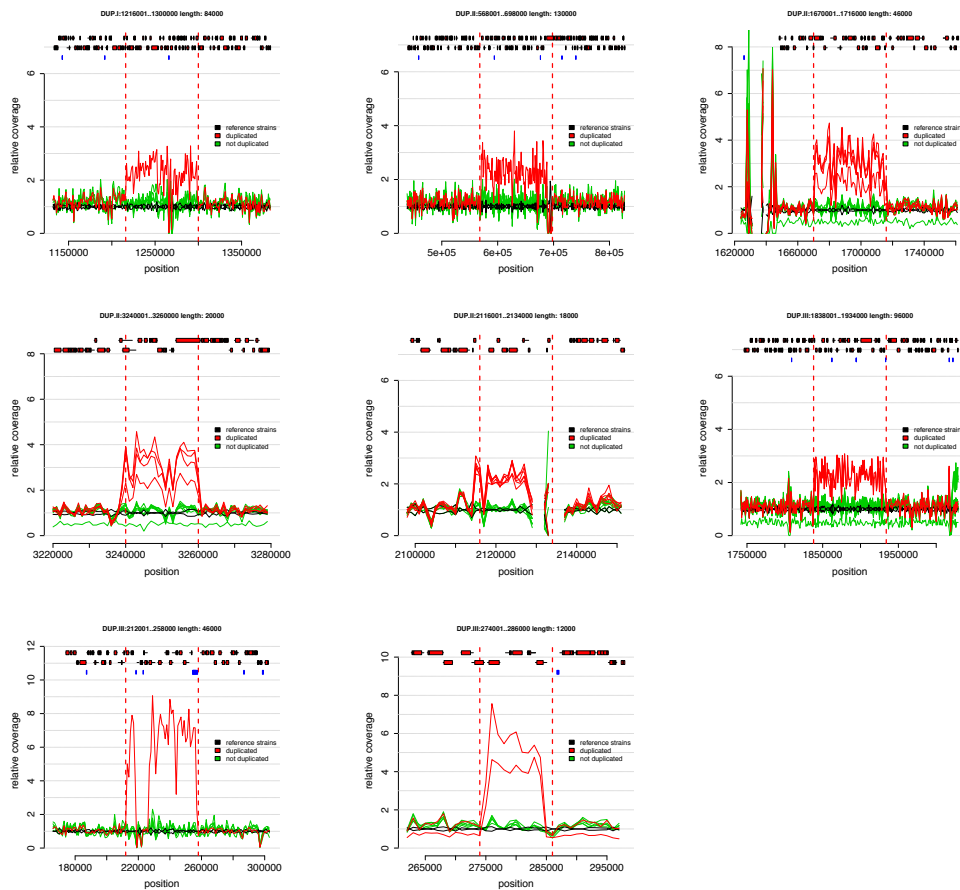
the same chromosome and calculating the gene coverage of these regions. Essential genes

were those with the Fission Yeast Phenotype Ontology term defined as FYPO:0002061

("inviable") in PomBase.

943

944



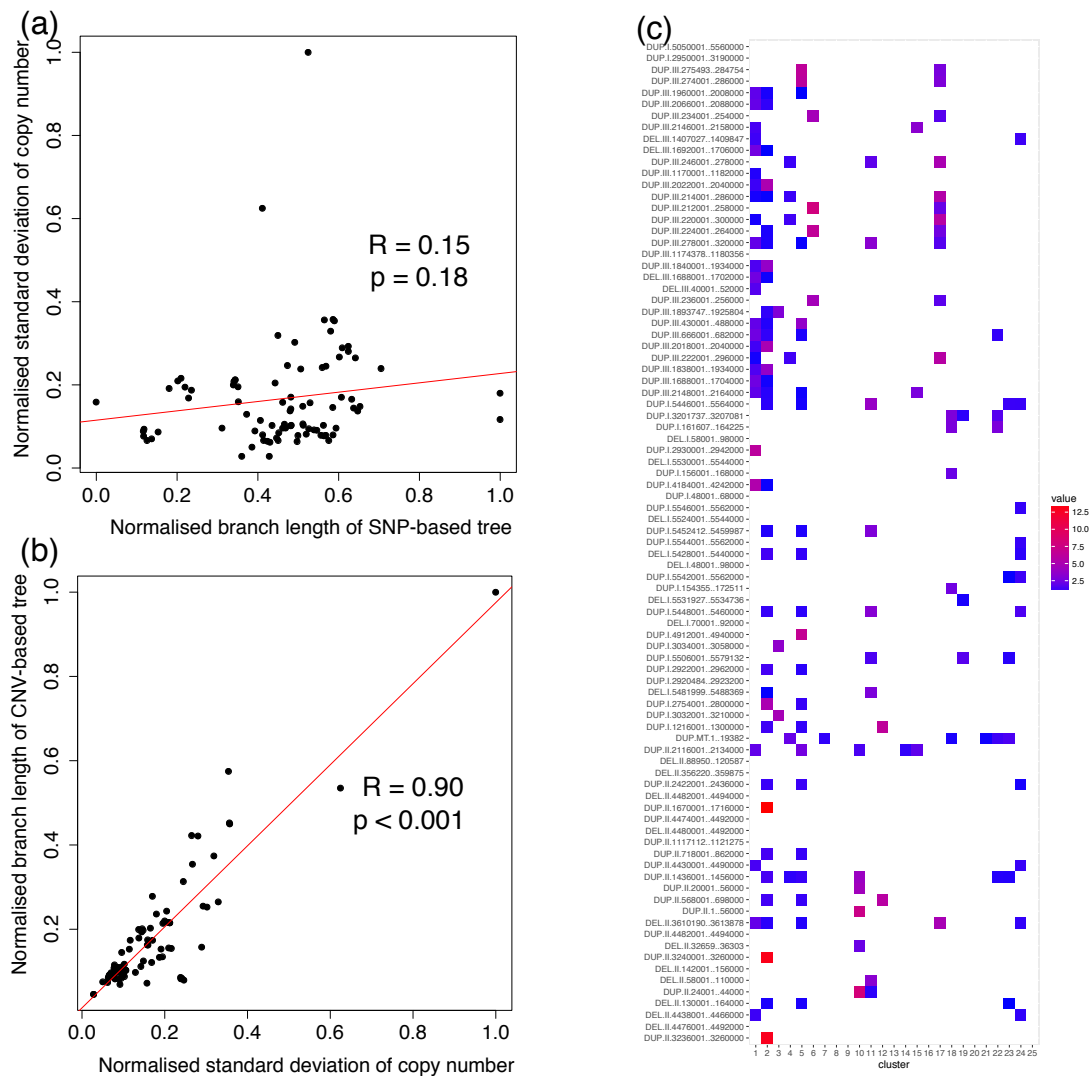
945

946 **Supplementary Figure 3. Duplications that segregate within closely related strains.**

947 Plots show the average coverage in 1 kb non-overlapping windows for strains with a
948 duplication (red) and all closely related strains without duplication (green); all these
949 strains differ by <150 SNPs. The coverage of the two standard reference strains (h^+ and h^-)
950) is shown in black. Top row, from left: variant DUP.I:1216001..1300000 (cluster 12,
951 from Japan in 57), DUP.II:568001..698000 (cluster 12), DUP.II:1670001..1716000
952 (cluster 2, unknown origin), second row DUP.II:3240001..3260000 (cluster 2),
953 DUP.II:2116001..2134000 (cluster 1, includes reference strain from French grapes in
954 1947), DUP.III:1838001..1934000 (cluster 2, various locations 1921-22). Bottom row:
955 DUP.III:212001..258000 (cluster 6, Jamaica/USA), and DUP.III:274001..286000 (cluster
956 5, Sicily 1966). Genes are shown on top of plots with exons as red rectangles and
957 retrotransposon LTRs as blue rectangles.

958

959
960

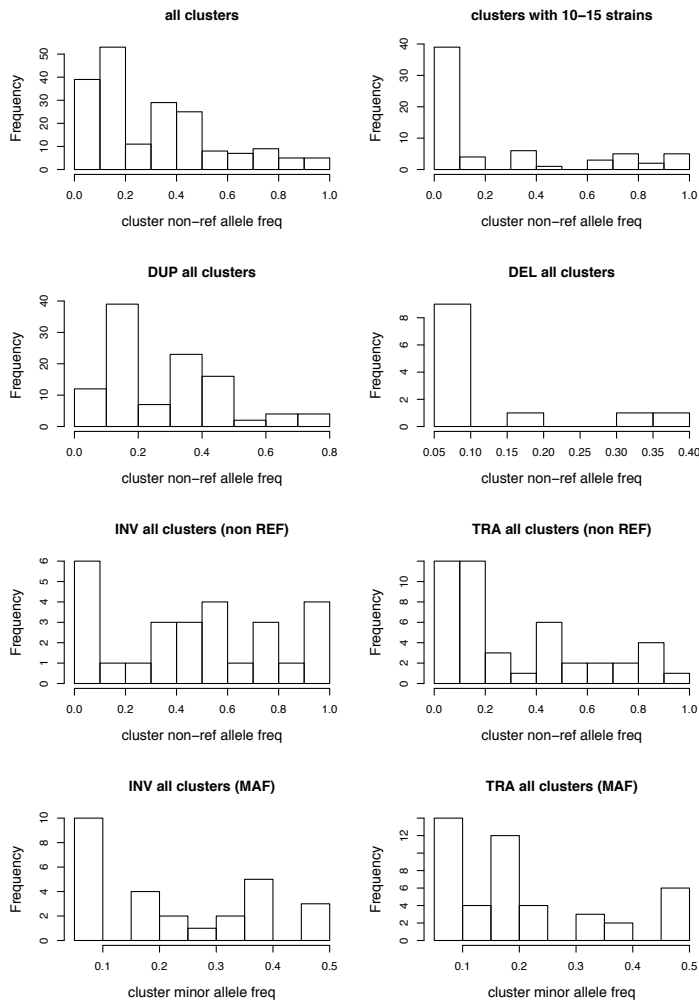


961
962
963
964
965
966
967
968
969

Supplementary Figure 4. Relative standard deviation of copy number variation within clusters. (a) Standard deviation of copy number for a CNV across the dataset is only weakly correlated with the total branch length of SNP-based phylogeny from the 20kb up- and down-stream phylogeny. (b) Standard deviation is highly correlated with the branch length of a CNV-based neighbour-joining tree. (c) The relative standard deviation of each CNV within each identified cluster of strains (<150 SNPs apart) relative to its change in rest of the dataset. For clarity, all relative standard deviations <1 are shown as white.

970

971



972

973

974

975

976

977

978

979

980

981

982

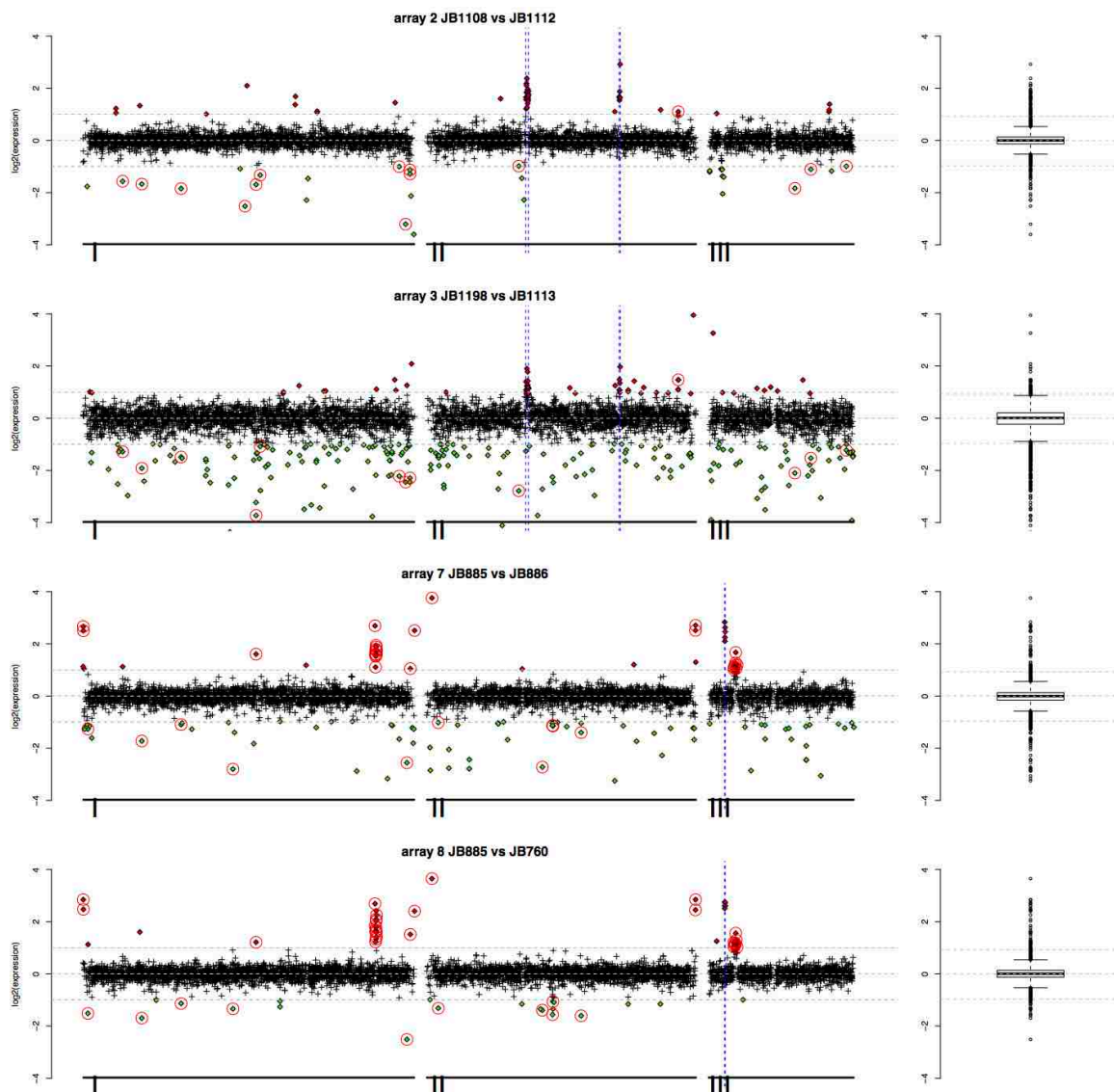
983

984

Supplementary Figure 5. Copy number variants are usually rare alleles within clonal populations. Clonal clusters, or clonal populations all differ by < 150 SNPs. In rows, from top left; we show the within-cluster frequency of the non-reference allele for all SVs, which is skewed to rare alleles. Limiting this analysis to cluster with 10 to 15 strains highlights the low frequency of non-reference alleles. Second row; CNVs (duplications and deletions) are skewed to rare alleles, because the non-reference allele is usually the derived allele. Third row; inversions and translocations are not skewed to the non-reference allele, but here non-reference alleles are not necessarily the derived allele. Bottom row; the *minor allele* of inversions and translocations, however, is skewed to rare alleles.

985

986



987

988

989

990

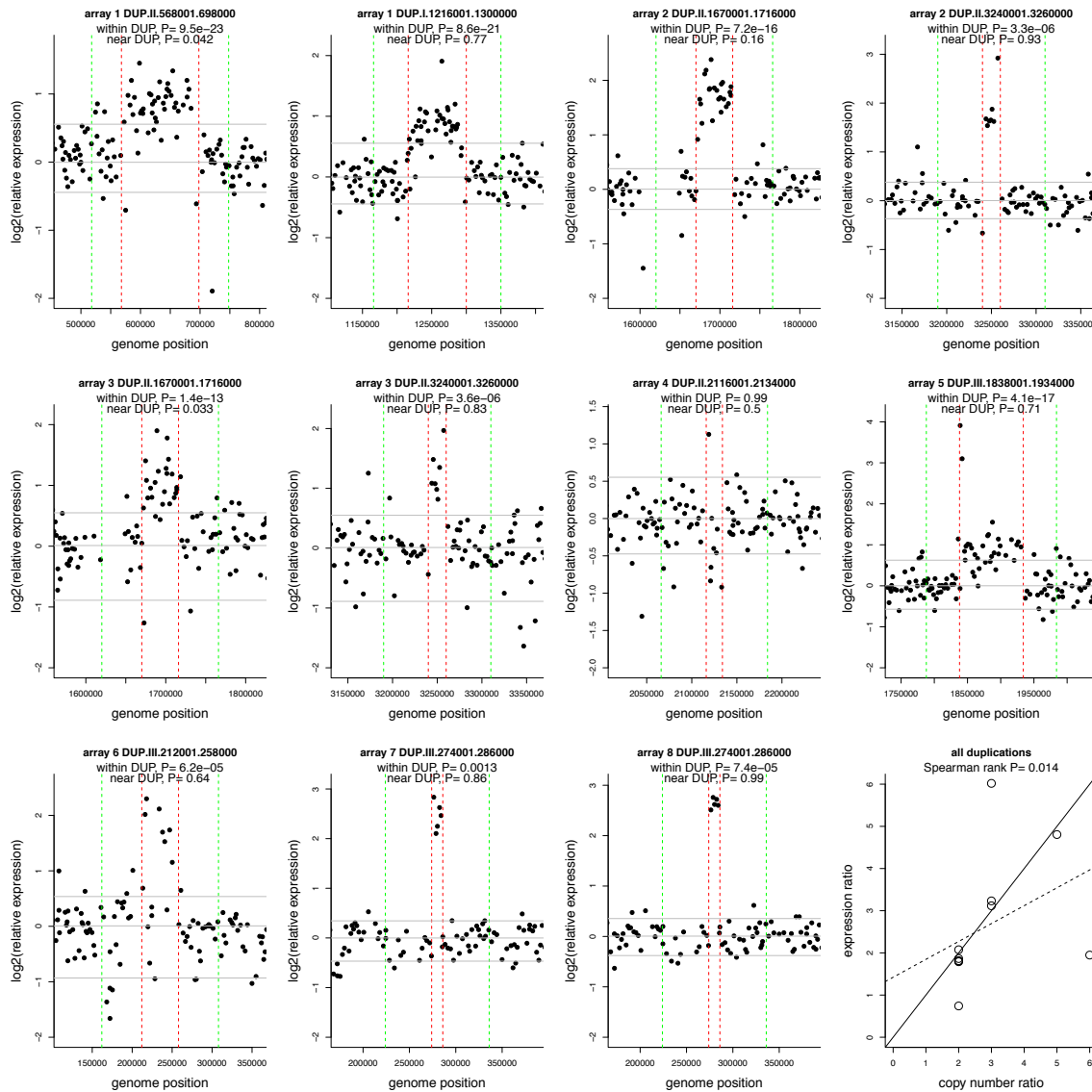
991

992

993

994

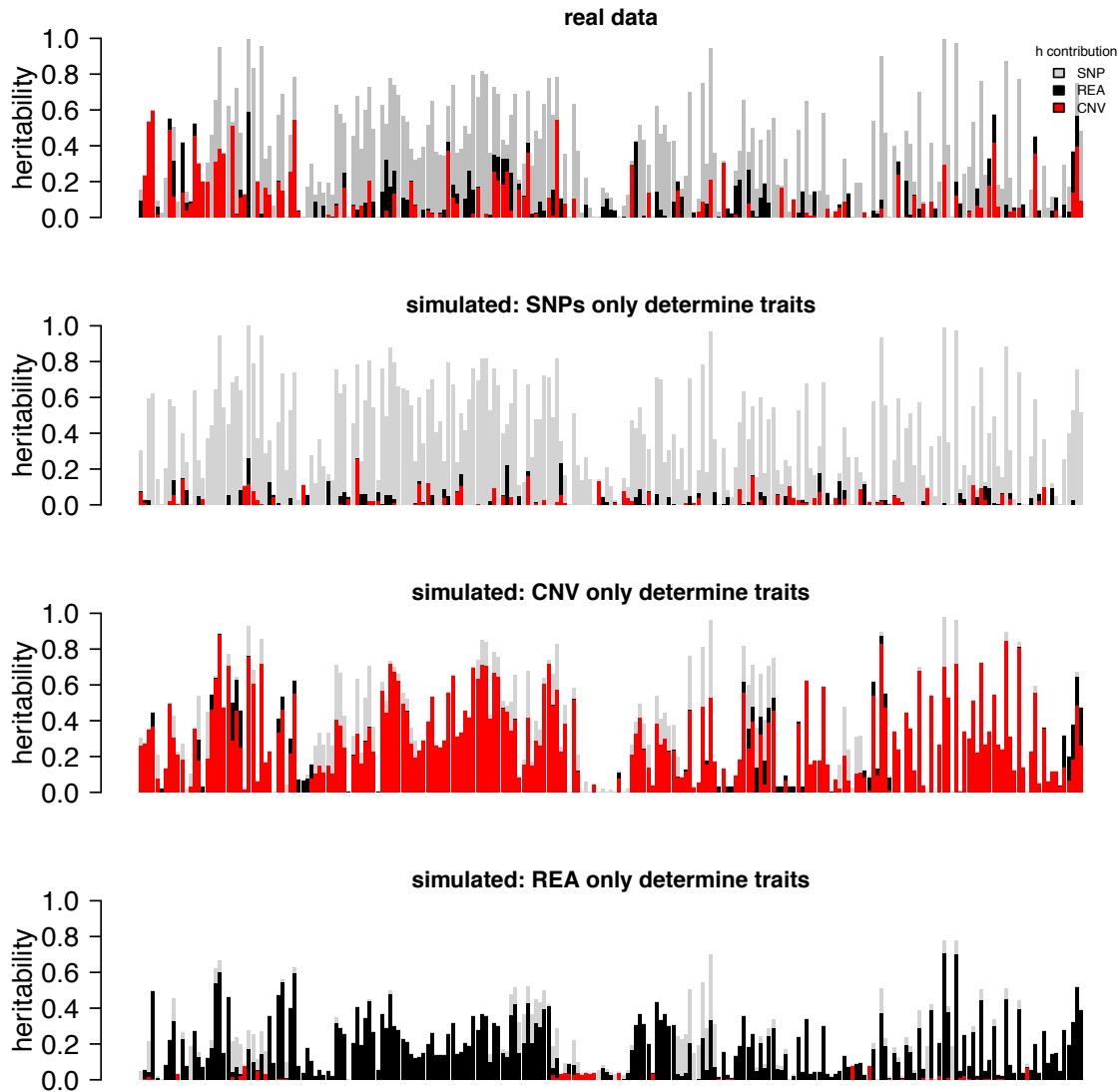
Supplementary Figure 6. Chromosome-scale view of gene expression changes. The relative gene expression levels (strain1/strain2) for arrays 2 and 3, and arrays 7 and 8 are shown with their positions on the three chromosomes. Filled circles indicates genes that we consider to be up-regulated (red) or repressed (green). Those highlighted with open red circles are consistently altered in both arrays (either 2+3, or 7+8). The blue lines show where the segregating duplications are. Box plots at right show the spread of data.



995

996 **Supplementary Figure 7. No significant increase in gene expression immediately**
 997 **adjacent to duplications.** For each duplication examined with DNA arrays, we show the
 998 relative expression (strain 1 vs strain 2) near the duplication. P-values show the support
 999 for the DUP genes within the duplication (red vertical lines), or the 50 kb adjacent to the
 1000 duplication (green vertical lines) being more highly expressed than all other genes in the
 1001 chromosome (one-sided Wilcoxon rank sum tests). The grey horizontal lines show the 5th,
 1002 50th and 95th percentiles for gene expression data on the chromosome. The bottom right
 1003 panel shows that the median increase in expression level within a duplication correlates
 1004 with the increase in genomic copy number. The solid back line shows the expected
 1005 increase for the 1:1 correspondence between genomic copy number and relative
 1006 expression (the line $y=x$), and the dashed line shows the linear model for the data. Copy
 1007 number and relative expression change are correlated (Spearman rank correlation $\rho =$
 1008 0.71 and $P = 0.014$).

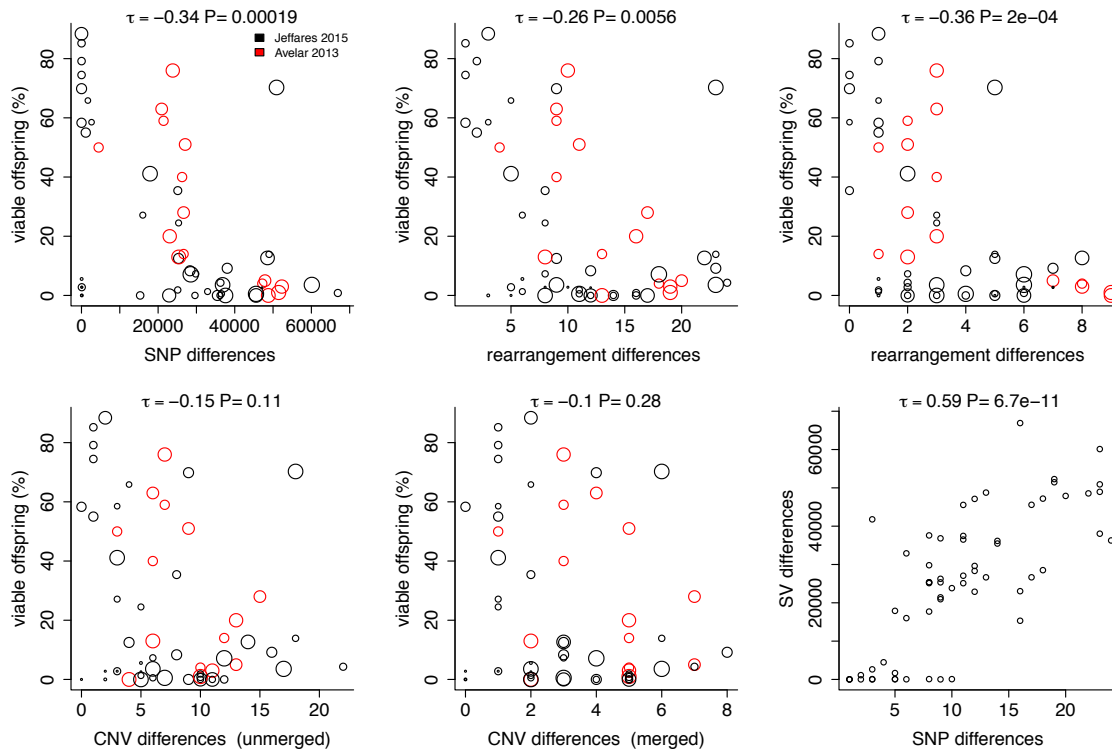
1009
1010
1011



1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022

Supplementary Figure 6. Contributions of SNPs, CNVs and rearrangements to traits. Top panel: for 227 traits, we show the total heritability estimated by the combination of 243,289 SNPs (green), 87 CNVs (red), and 26 rearrangements (grey). We then simulated data that was entirely due to the effects of SNPs (second panel), entirely due to the effects of CNVs (next panel) or entirely due to the effects of rearrangements (lower). In the second panel (entirely due to the effects of SNPs), the contribution of CNVs and rearrangements are artefacts, but these are relatively minor. This analysis indicates that the estimates are not strongly affected by linkage.

1023



1024

1025

1026 **Supplementary Figure 7. Correlations between spore viability, parental SNP-genetic**

1027 **distance and parental SV-genetic distance.** Spore viability was measured for 58

1028 crosses in total, including data from both Jeffares, et al. ¹⁷ (black) and Avelar, et al. ⁶

1029 (red), with each circle representing one cross. Unmerged CNV differences count any

1030 CNV as being different between parents when either start or end coordinates are more

1031 than 1 kb apart. Because this definition can cause us to count largely overlapping events

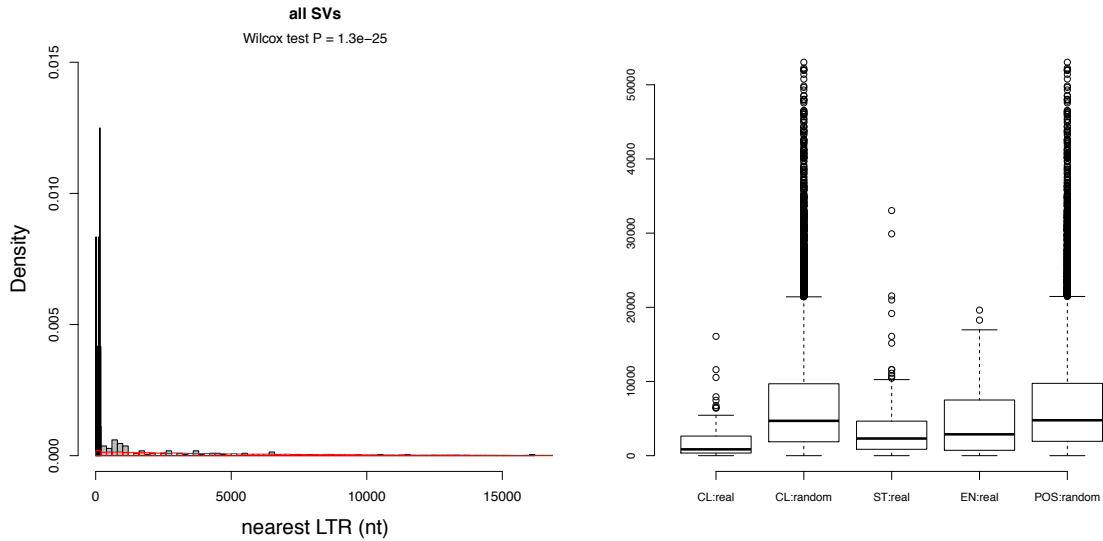
1032 as ‘different’, we also counted ‘merged’ differences where two CNVs were considered

1033 different only if their overlap was >50% of the total of both variants. This approach will

1034 exclude nested CNVs. CNV-genetic distance is not significantly correlated with viability

1035 in either case.

1036

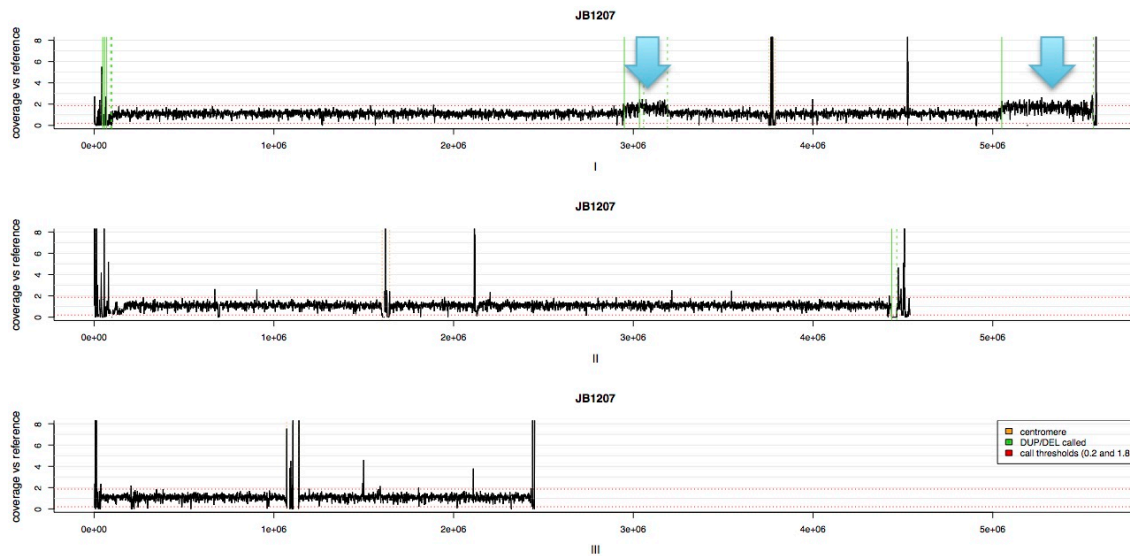


1037
1038
1039
1040
1041
1042
1043
1044
1045
1046

Supplementary Figure 8. SVs are enriched close to retrotransposon LTRs. For all SVs, we computed the closest distance of start or end coordinates to any LTR discovered previously¹⁷. As a control, we compute the closest distance of 10 random coordinates on the same chromosome. Left: the distributions of distances for real SVs (grey), those that are within 200nt (black) or random coordinates (red). Right: using the same analysis, we show the closest distance of real SVs (CL: real), and random coordinates (C: random). We also show that both start and end coordinates of SVs (ST:real, EN:real) are closer than random positions (POS:random).

1047

1048



1049

1050

1051

1052

1053

1054

1055

1056

1057

Supplementary Figure 9. Chromosome-scale read coverage plots for three chromosomes of strain JB1207. Coverage is calculated relative to the reference strain (JB22 in our collection). Two large duplications that did not satisfy the criteria used to detect CNVs with cn.MOPs are indicated with blue arrows (DUP.I:2950001..3190000, 240kb and DUP.I:5050001..5560000, 510kb).