

1 SUPPLEMENTAL INFORMATION

2 Supplemental Figures

3	Figure S1. Benchmark of structural variants callers using simulated events.	2
4	Figure S2. Breakpoint differences between members of each structural variation cluster.	3
5	Figure S3. Comparison of SVs detected by Wang et al 2018 with our SV dataset.	4
6	Figure S4. Principal components analysis for the DEL dataset using all (a, b) and 560 high-coverage	5
7	samples (c, d).	
8	Figure S5. Distribution of short deletions around transcription start sites (TSSs).	6
9	Figure S6. Average sequence complexity around the vicinity of transcription start sites.	8
10	Figure S7. SNP versus SV densities in 100kb sliding windows across the Nipponbare RefSeq.	9
11	Figure S8. SNP and SV distribution across the genome and colocalization of SV peaks and gene	10
12	classes.	
13	Figure S9. Known structural variants verified using the dataset.	11
14	Figure S10. Genome-wide association studies.	12
15	Figure S11. Sensitivity of callers on duplication prediction based on different limits for reciprocal	13
16	overlap.	
17	Figure S12. Ratio of caller sensitivity on 70% RO over 90% RO.	14
18	Figure S13. Impact on the rice genome of structural variants predicted by NGSEP.	15
19	Figure S14. Structure and enrichment analysis of CNVs.	16
20	Figure S15. Distribution of short deletions around TSS.	17
21	Figure S16. Distribution of the number of deletions in the vicinities of start and end of	18
22	transcription and translation.	
23	Figure S17. Dot plots for selected variants validated as true positives.	19
24	Figure S18. Assessment of copy number variation prediction.	20
25	Figure S19. Comparison of deletion and insertion sequences to known or potentially active TEs.	21

26

27 Supplemental Tables

28	Table S1. Transposable and repeat elements (size >50bp) among the structural variants.	22
29	Table S2. Strategies and SV types supported by each caller.	23
30	Table S3. Summary of structural variants (SV) identified by NGSEP on the complete dataset	24
31	combining read depth (RD), read pair (RP) and split read (SR) approaches.	
32	Table S4. Manually validated SVs using N22 (CX368) reference genome.	25

33

34 Supplemental Methods

26

35

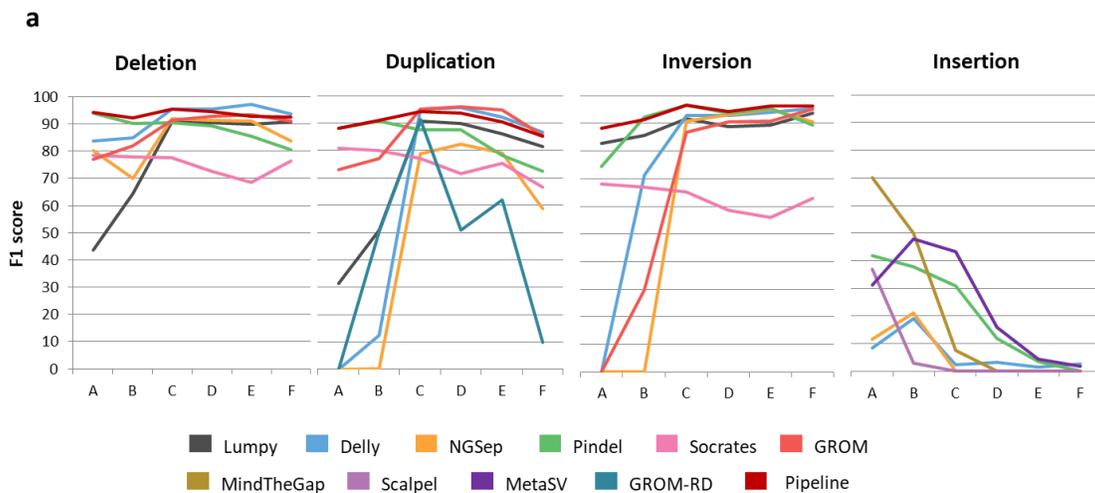
36

37

38

39 **Figure S1. Benchmark of structural variants callers using simulated events.** (a) F1-score (harmonic
 40 average of the precision and recall) of each caller is computed per variant type and size bin. Sizes
 41 are binned according to lengths: A (50-150 bp), B (151-500 bp), C (500-5000 bp), D (5-50 kb), E (50-
 42 250 kb) and F (0.25-1Mb). (b) Sensitivity of GATK-UG and the SV discovery pipeline according to the
 43 sizes of variants. GATK-UG is sensitive only for short insertions and deletions and must be
 44 complemented by other variant callers, such as our pipeline built for structural variants. (c)
 45 Percentage of GATK-UG within SV-pipeline deletions predicted in random samples, stratified by
 46 deletion size.

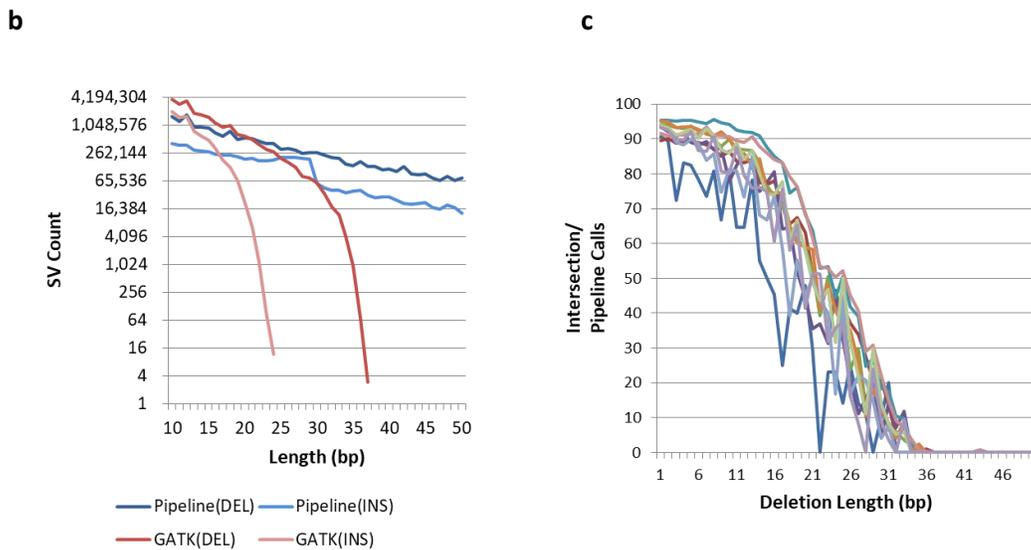
47



48

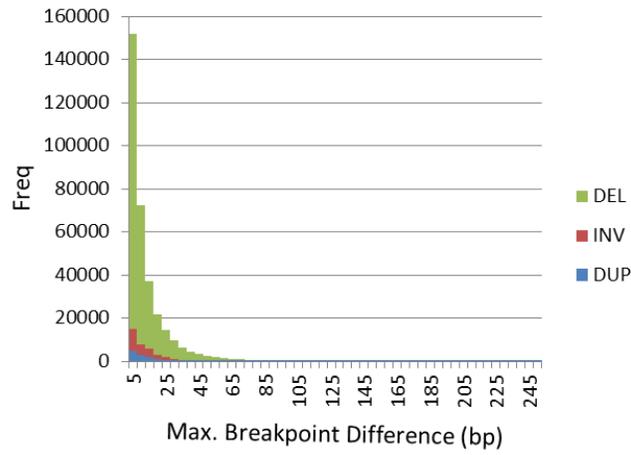
49

50



51 **Figure S2. Breakpoint differences between members of each structural variation cluster.** After
52 implementing hierarchical clustering, groupings were refined and the computed mean distance is
53 2.2%.

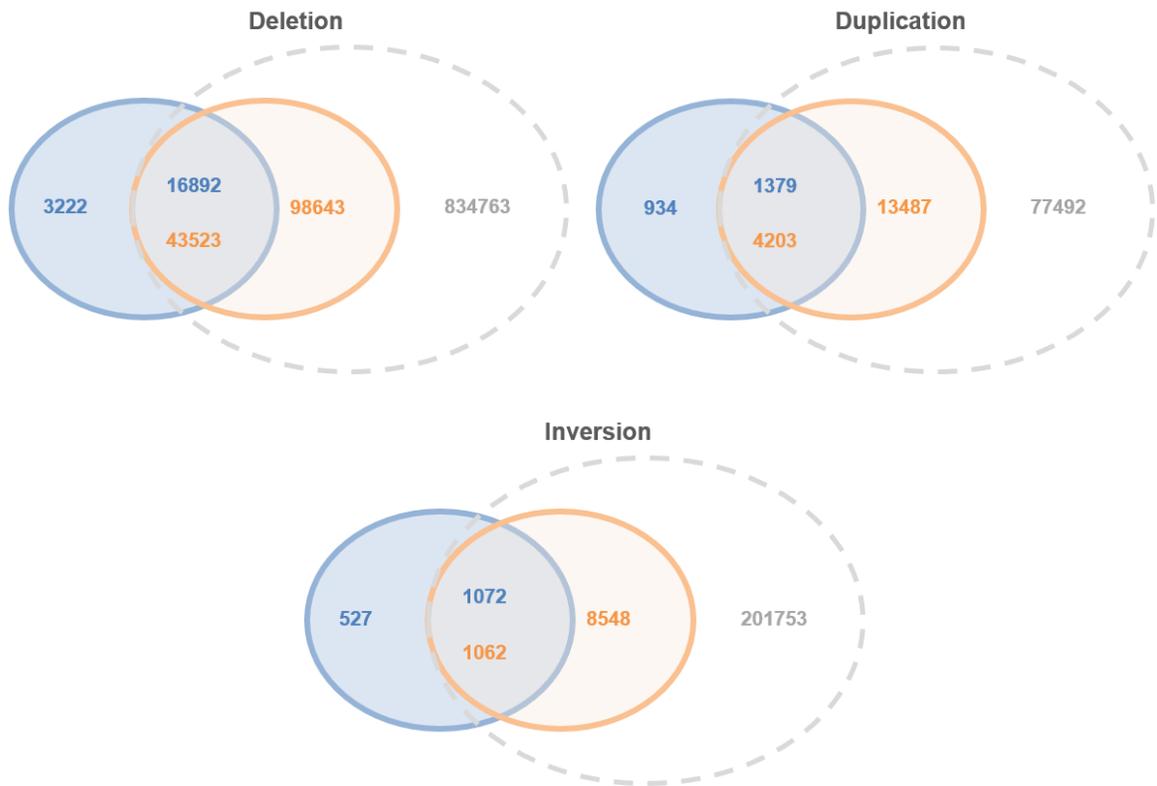
54



55

56

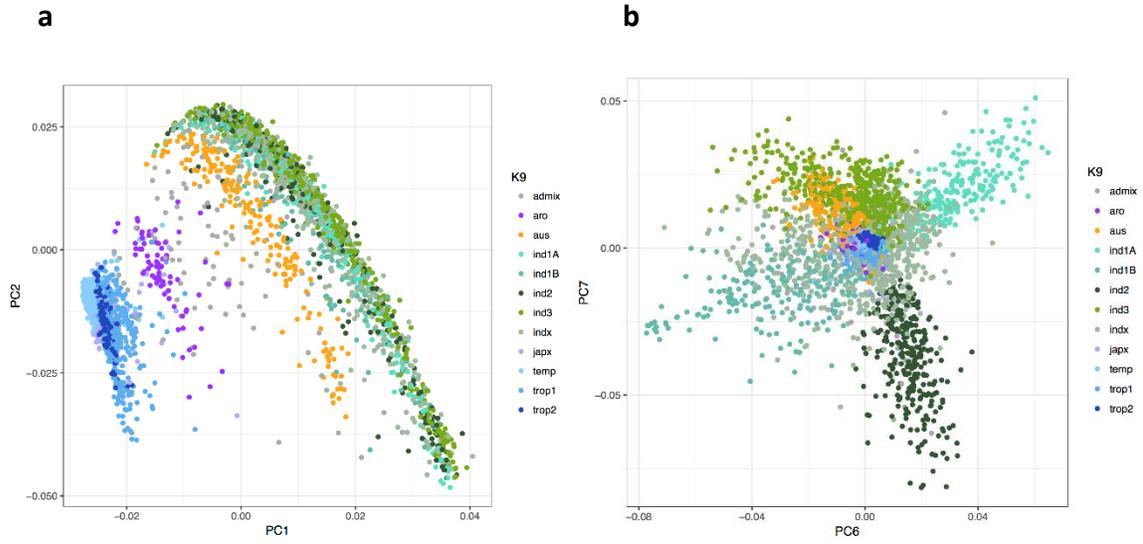
57 **Figure S3. Comparison of SVs detected by Wang et al 2018 with our SV dataset.** For a fair comparison
 58 we applied the same filtering as in Wang et al: size between 100 bp and 500 kb; frequencies
 59 between 6 samples and 80% of samples; only SVs from the 453 subset samples were used.
 60 Translocations were excluded from this comparison due to the poor performance of the callers in
 61 detecting this specific SV type. The comparison of our clusters with Wang et al clusters is not one-
 62 to-one due to differences in the clustering method; thus, the circle intersections list both the
 63 number of Wang et al clusters that have >50% reciprocal overlap with our clusters (blue), and the
 64 number of our clusters (orange). Large differences in cluster number are due to the different
 65 clustering approaches. Our full dataset, including all samples and variants of size below 100 bp, is
 66 represented by the broken circle.



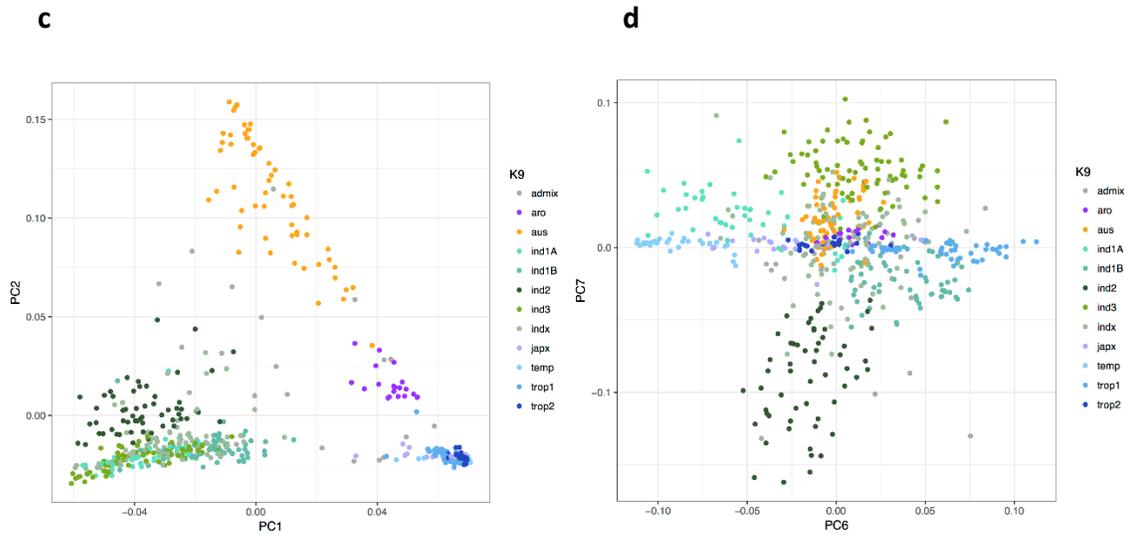
70

71 **Figure S4. Principal components analysis for the DEL dataset using all (a, b) and 560 high-coverage**
72 **samples (c, d). The color indicates clusters from Wang *et al.* defined by SNP data.**

73

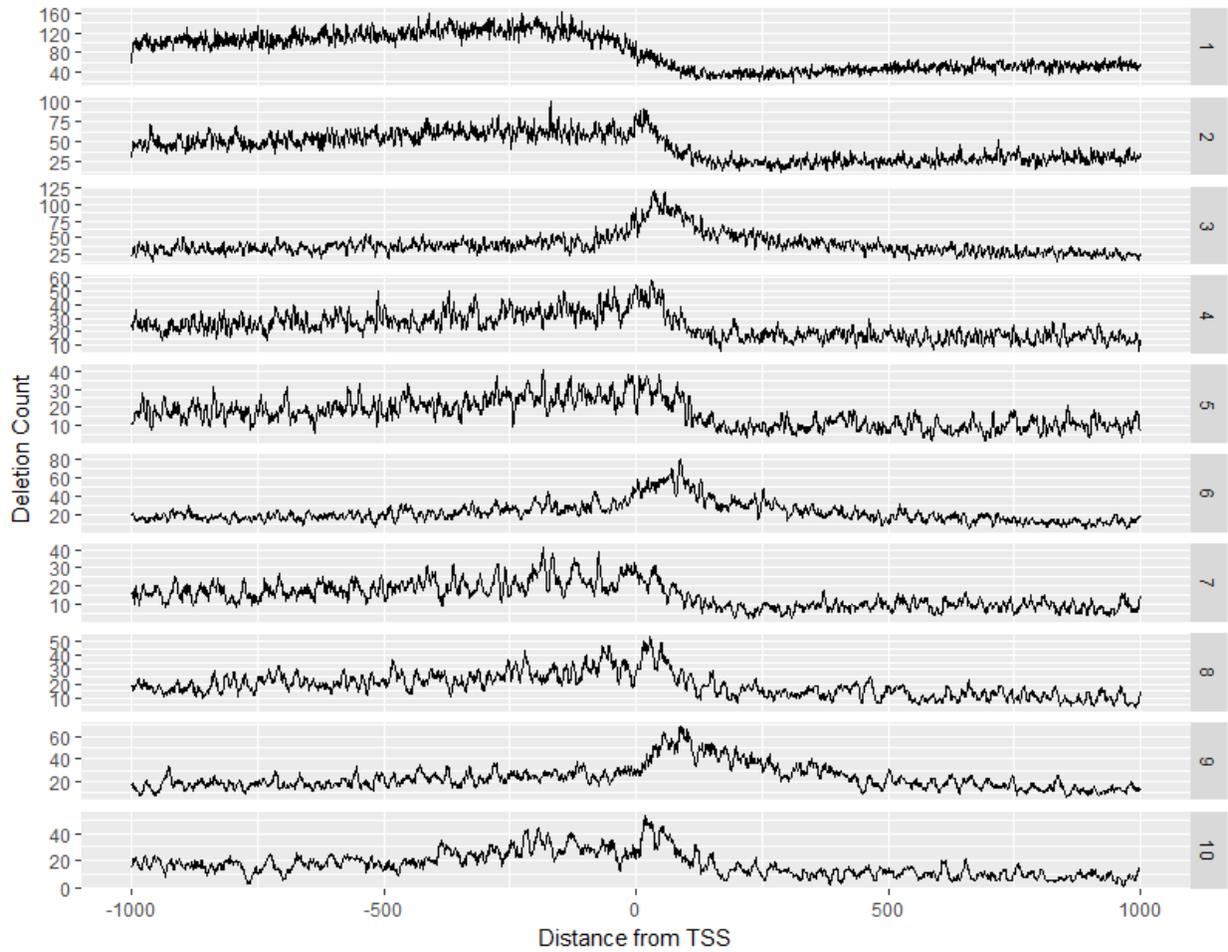


75

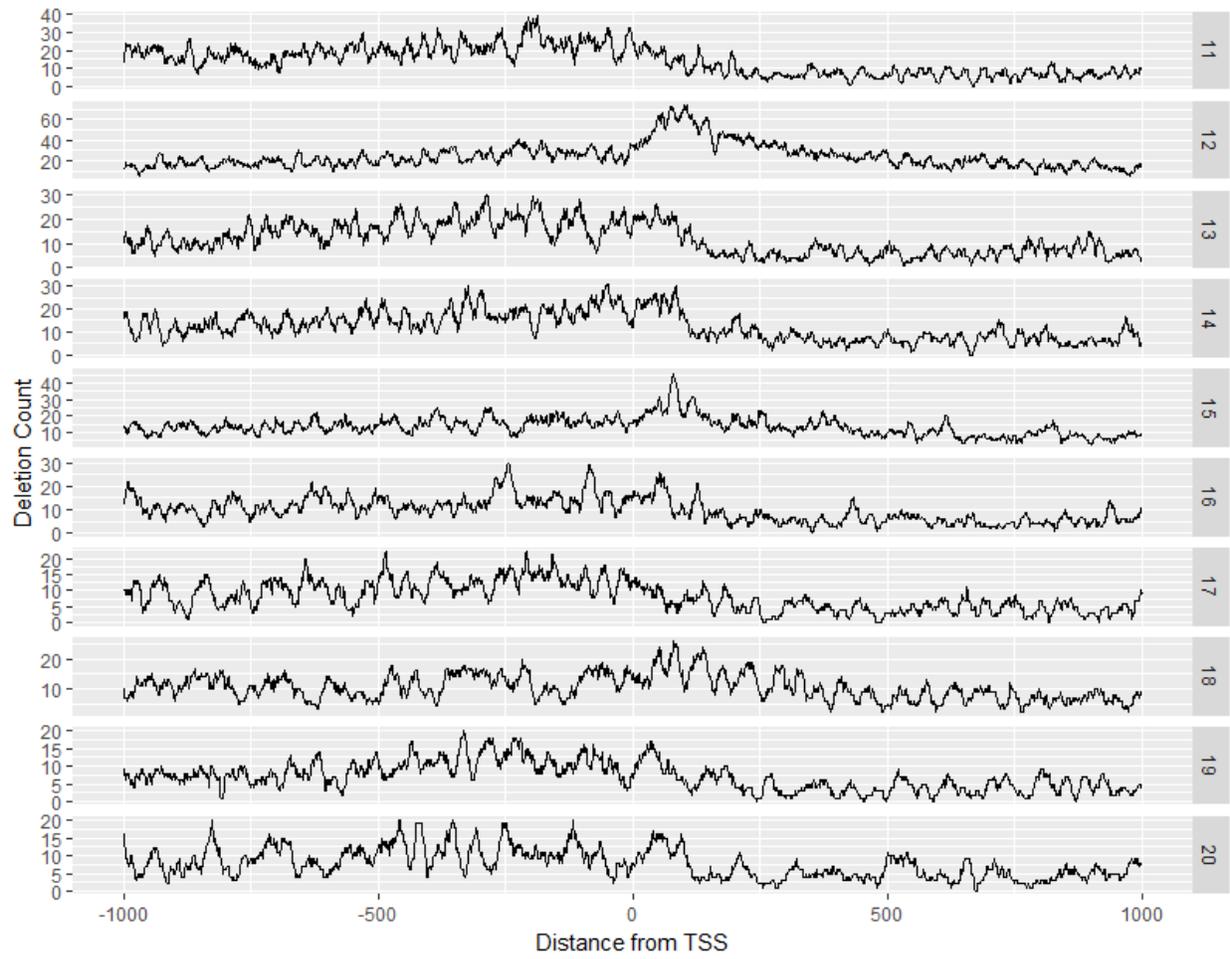


77

78 **Figure S5. Distribution of short deletions around transcription start sites (TSSs).** Deletions with sizes in
79 multiples of 3 nucleotides contribute largely to the peak in the 5' UTRs, and they have a 4-fold lower
80 distribution on both sides of these regions.
81



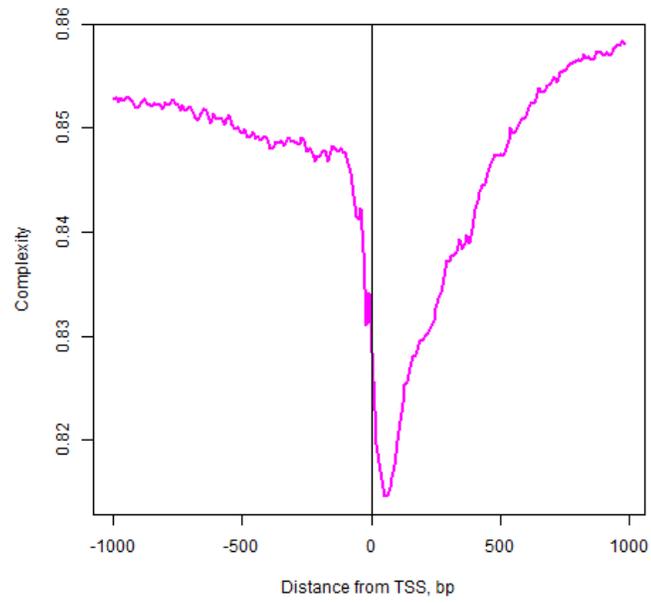
82



83

84

85 **Figure S6. Average sequence complexity around the vicinity of transcription start sites.**



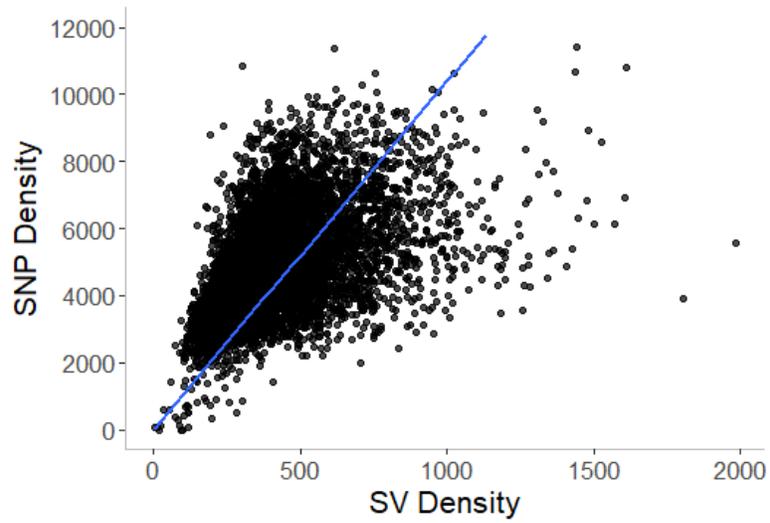
86

87

88

89 **Figure S7. SNP versus SV densities in 100kb sliding windows across the Nipponbare RefSeq.** SNP data
90 were taken from SNP-Seek database (Mansueto *et. al.* 2017).

91

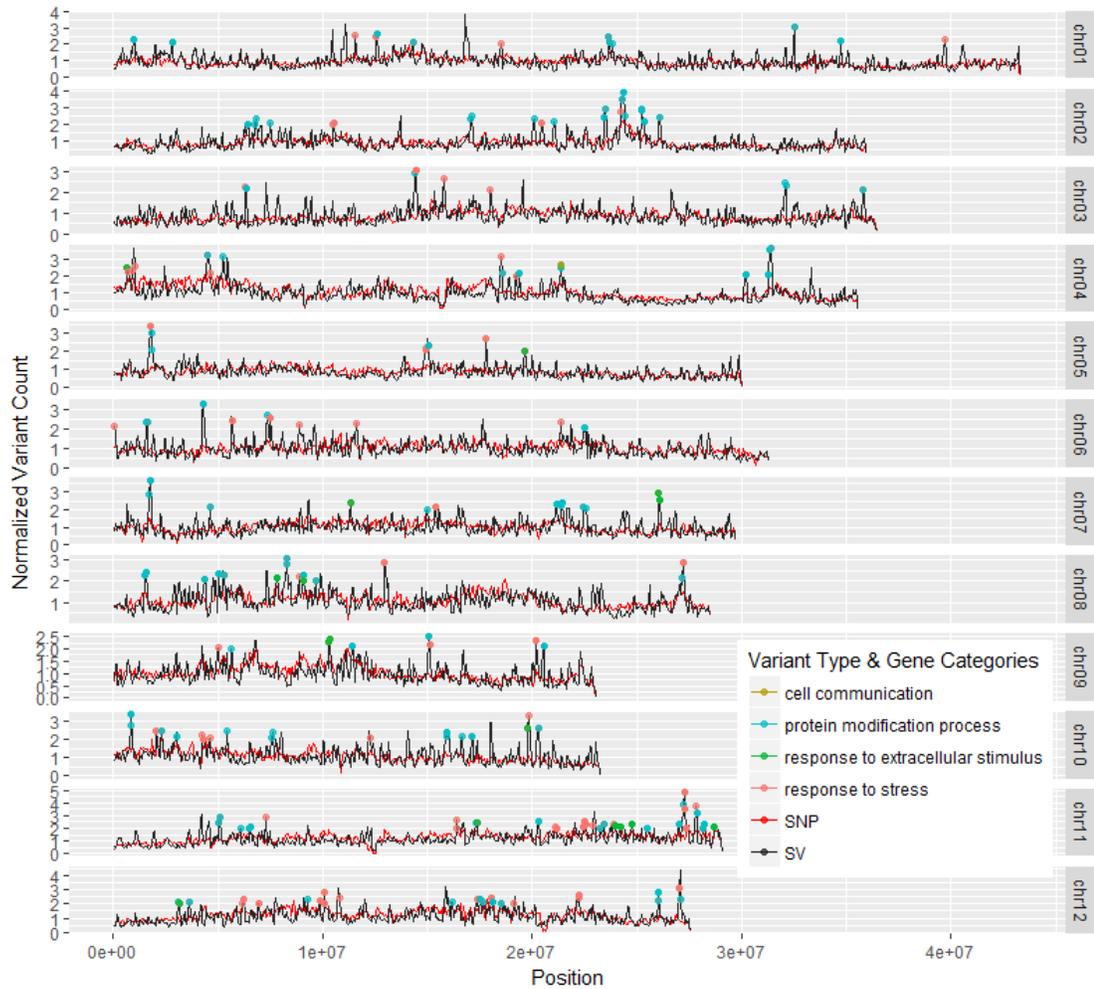


92

93

94 **Figure S8. SNP and SV distribution across the genome and colocalization of SV peaks and gene classes.**

95 The red and black lines shows the correlated distribution of SNPs and SVs, respectively, across the
96 chromosomes. Dots indicate regions with higher numbers of structural variants that contain genes
97 associated with stress and other responses.



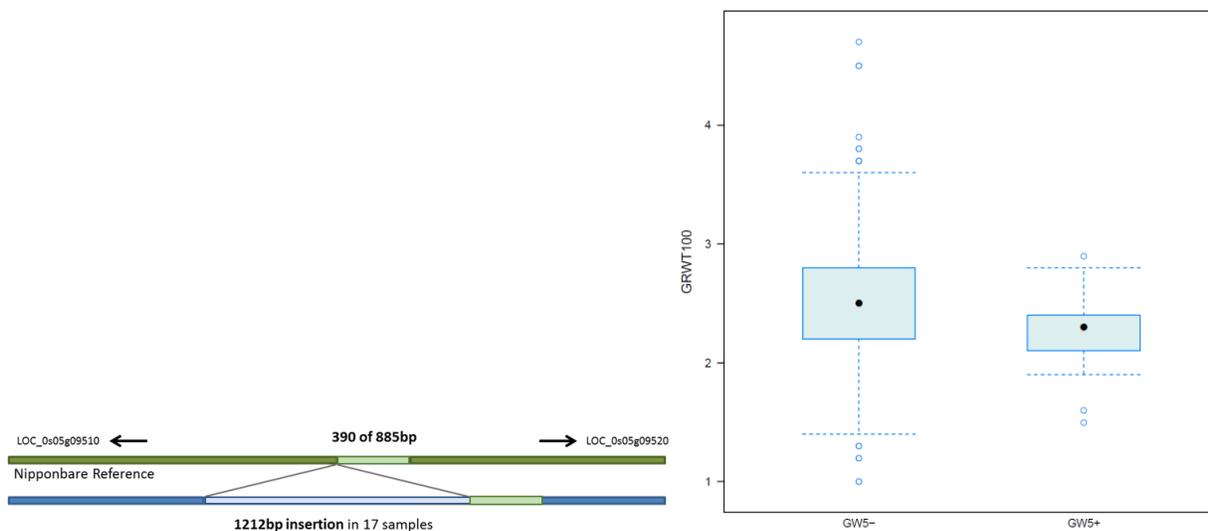
98

99

100

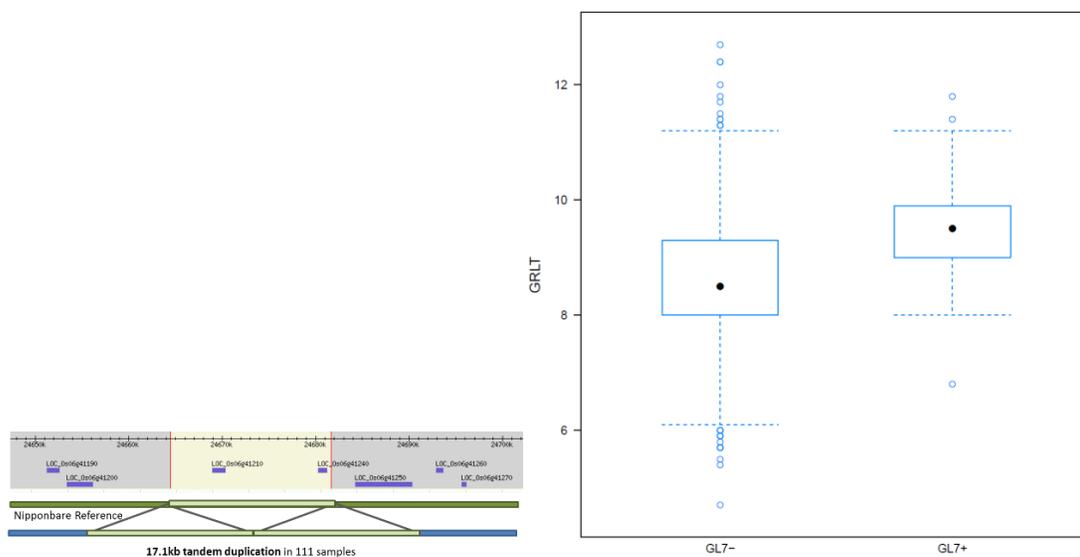
101 **Figure S9. Known structural variants verified using the dataset. a)** *GW5* gene is partially deleted from the
 102 Nipponbare RefSeq (green) and effect of the observed 1212 bp insertion on grain weight, p-value=0.04.
 103 **b)** tandem duplication near *GL7* (*LOC_Os06g41200*) gene and effect of the event on grain length, p-value
 104 = 1.0e-10.

105 **a**



106

107 **b**



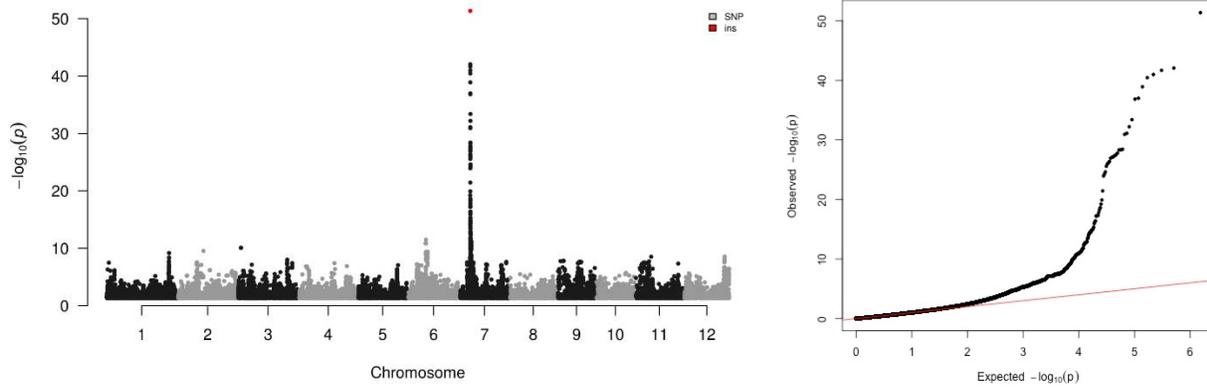
108

109

110 **Figure S10. Genome-wide association studies.** Manhattan and QQ plots for the **a) *Rc*** (red pericarp) trait
111 and **b) grain length** trait.

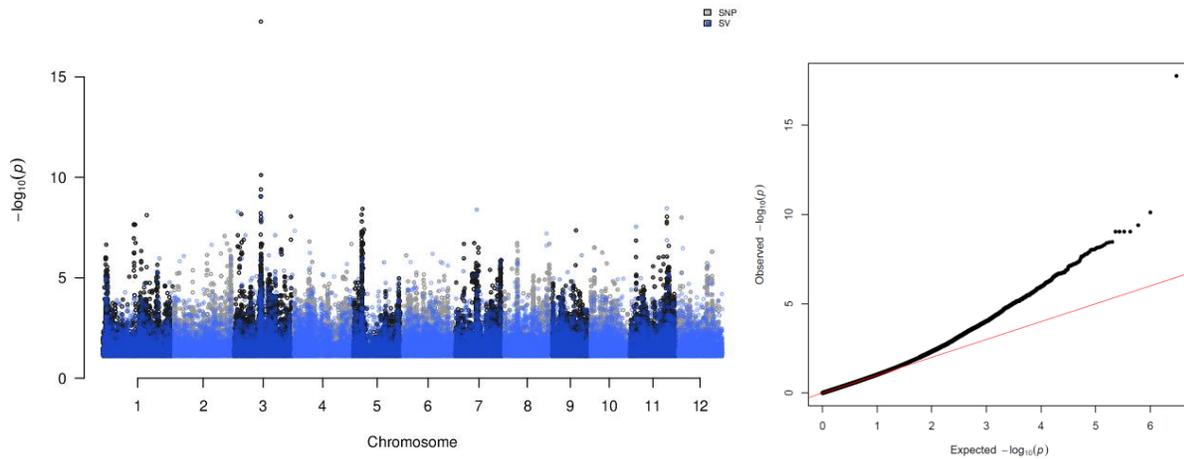
112

113 **a**



114

115 **b**



116

117

118 **Figure S11. Sensitivity of callers on duplication prediction based on different limits for reciprocal**
 119 **overlap.** The profiles of sensitivity show robustness of callers to the reciprocal overlap parameters.

120

121

122

123

124

125

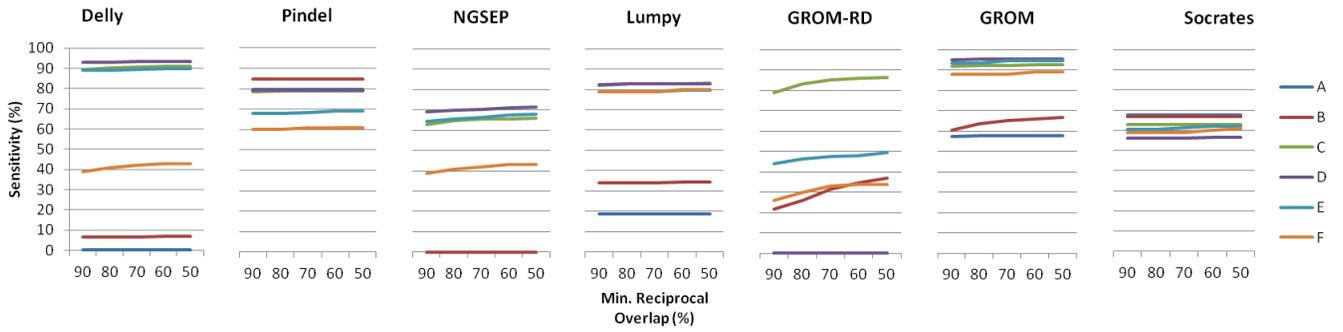
126

127

128

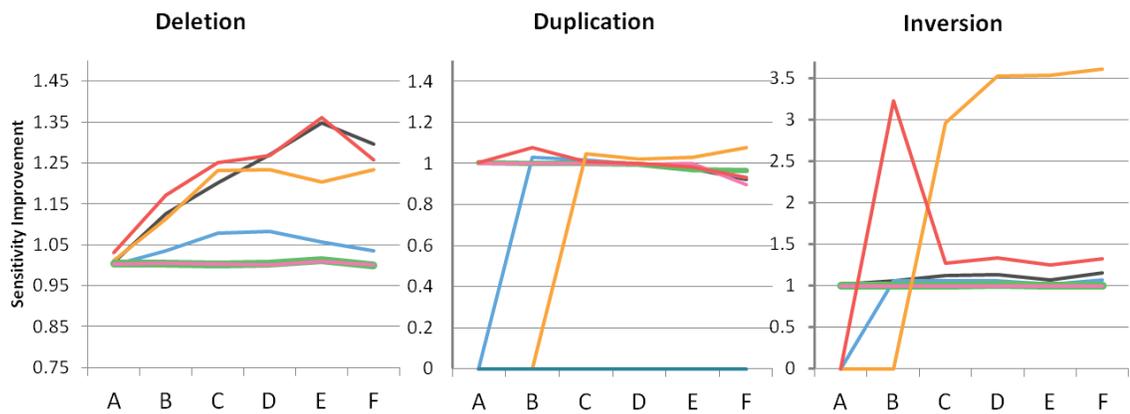
129

130



131 **Figure S12. Ratio of caller sensitivity on 70% RO over 90% RO.** A ratio of 1.0 means that the sensitivity
 132 did not change after relaxing the reciprocal overlap (RO) threshold. A near horizontal line shows
 133 reliability of breakpoint resolution across different size bins. Both Socrates (Schröder et al. 2014)
 134 and Pindel (Ye et al. 2009) consistently predicted precise breakpoints, thus reducing RO to 70% will
 135 barely improve their sensitivity.

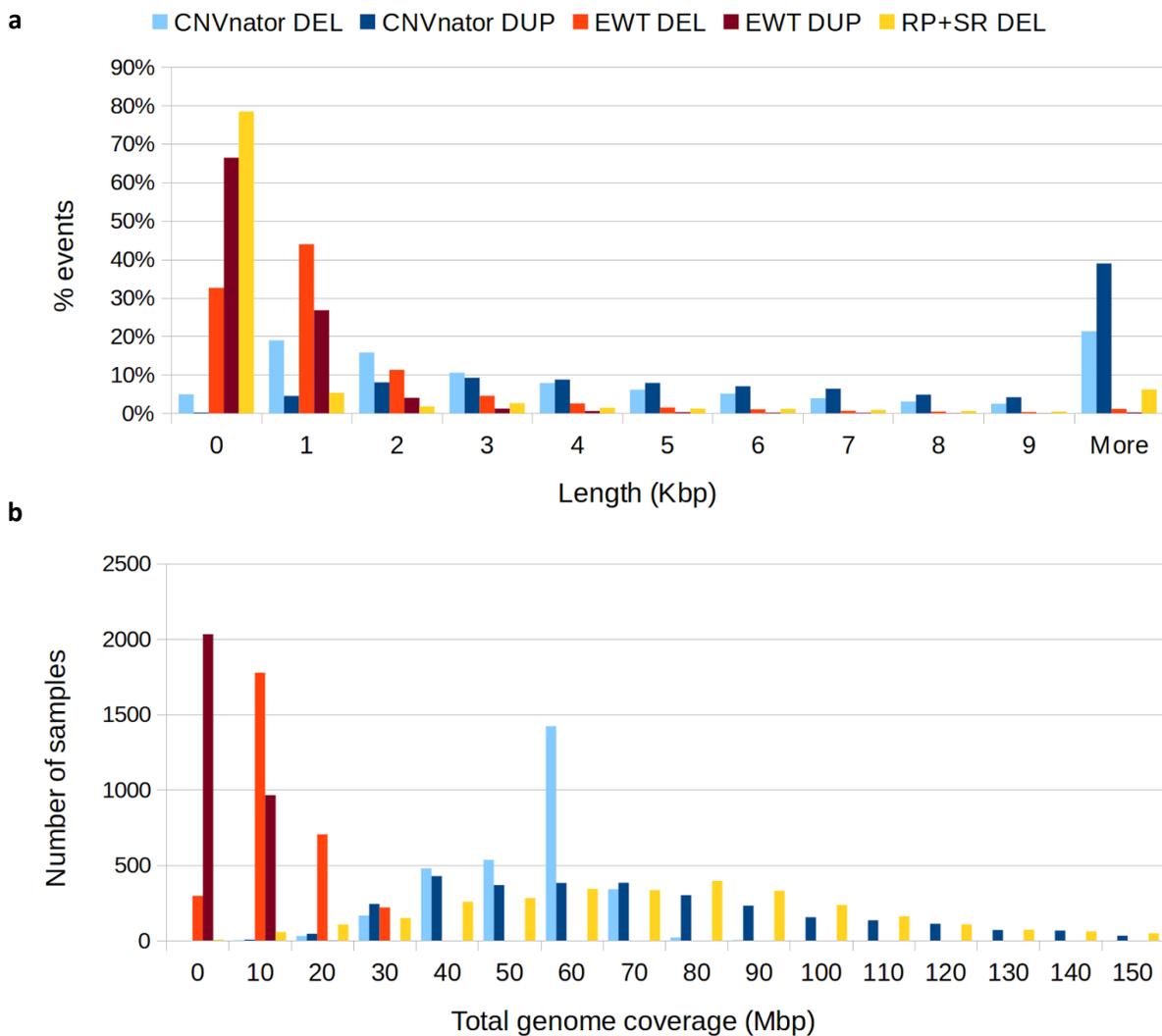
136



137

138

139 **Figure S13. Impact on the rice genome of structural variants predicted by NGSEP.** a) Length distribution
 140 of raw structural variants called by NGSEP. b) Distribution of the amount of reference genome
 141 covered by CNV structural variants across the 3K RG samples.
 142

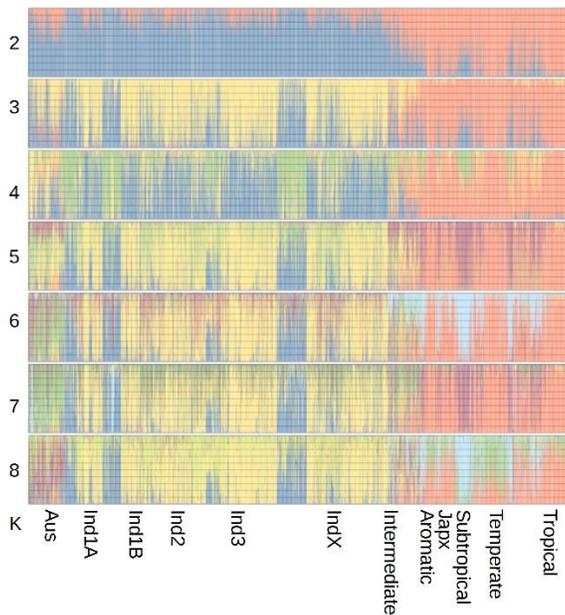


143

144

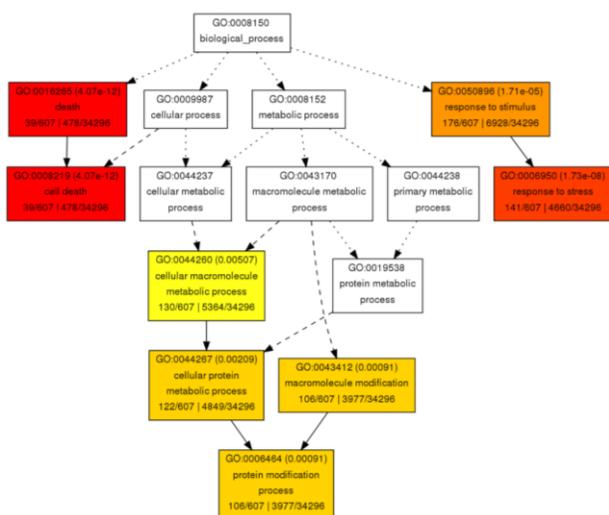
145 **Figure S14. Structure and enrichment analysis of CNVs.** a) Structure analysis using 2,839 CNV genotyped
 146 in at least 2,000 of the 3,023 samples, in non-repetitive regions of the genome and having the
 147 major allele in at most 80% of the samples. Enriched biological process (b) molecular function (c)
 148 gene ontology terms in genes covered by CNVs.

149 **a**

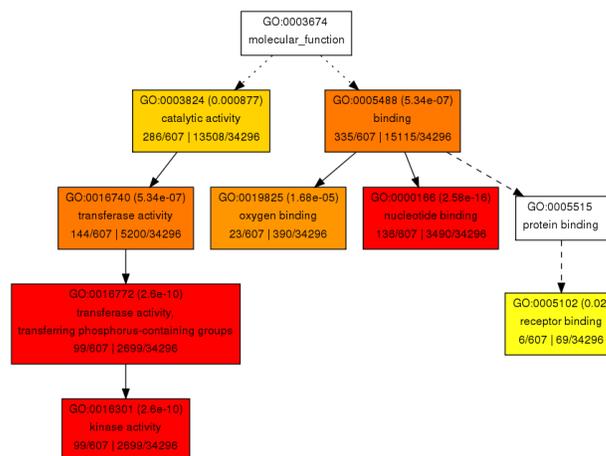


150

151 **b**



c

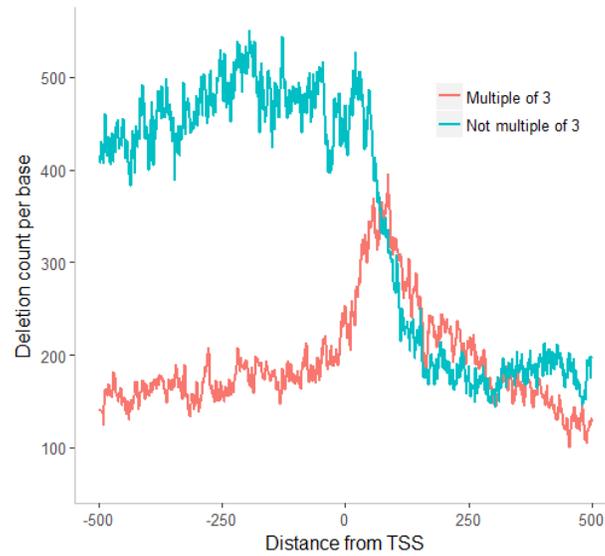


153 **Figure S15. Distribution of short deletions around TSS.** The plot shows a peak of deletions with size in
154 multiples of 3 bp downstream of TSS.

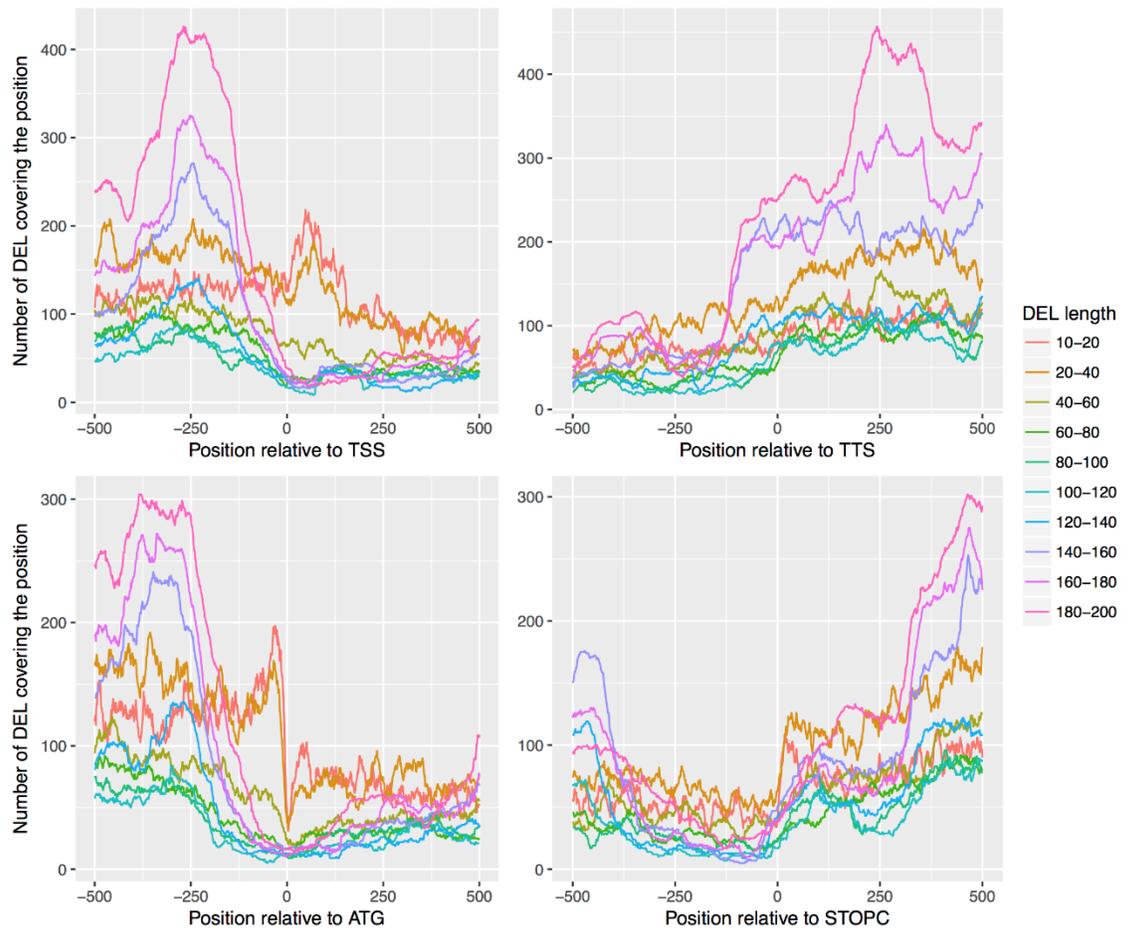
155

156

157

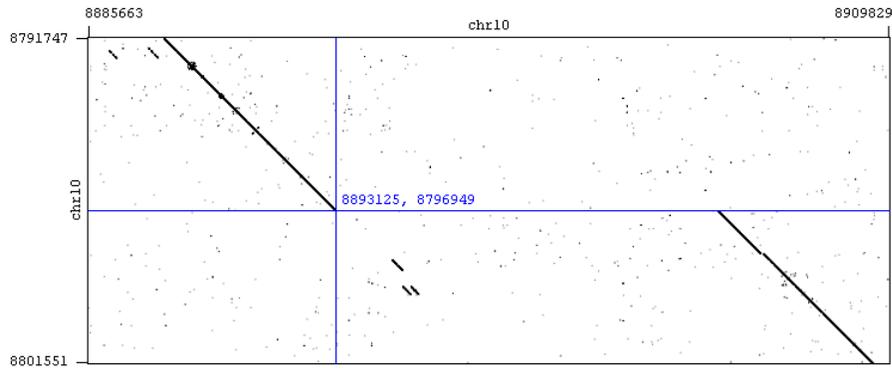


158 **Figure S16.** Distribution of the number of deletions in the vicinities of start and end of transcription and
159 translation.
160



161 **Figure S17. Dot plots for selected variants validated as true positives. a) Chromosome 10 11-kb deletion.**
162 **b) Chromosome 10 218-bp tandem duplication and c) Chromosome 12 640-bp inversion.** The blue
163 cross lines mark the start position of the structural variant.

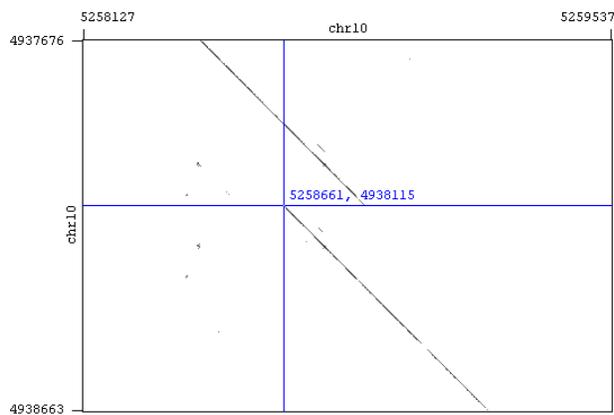
164 **a**



165

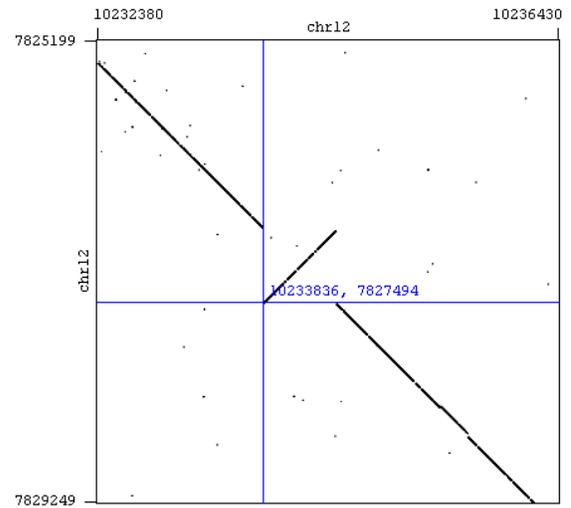
166

167 **b**



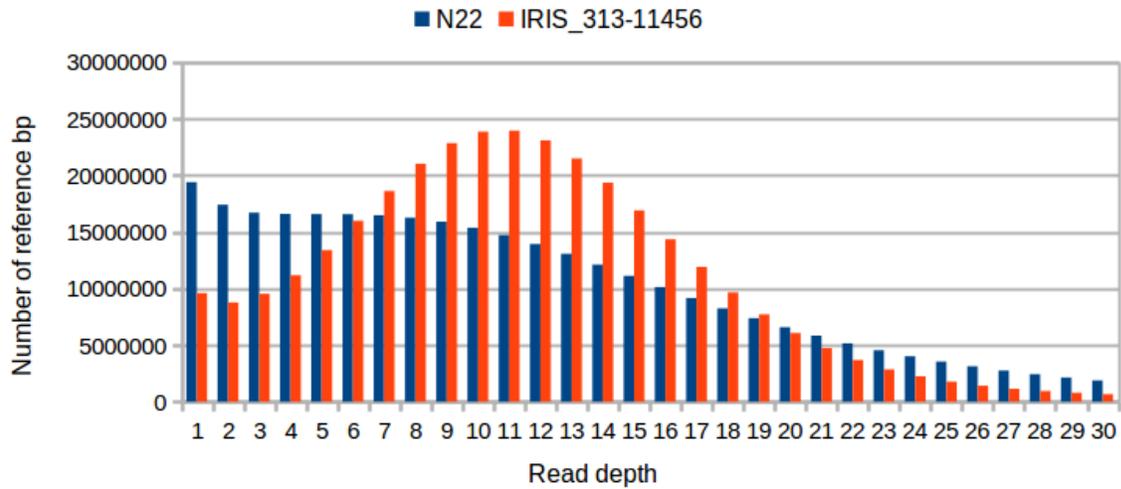
168

c

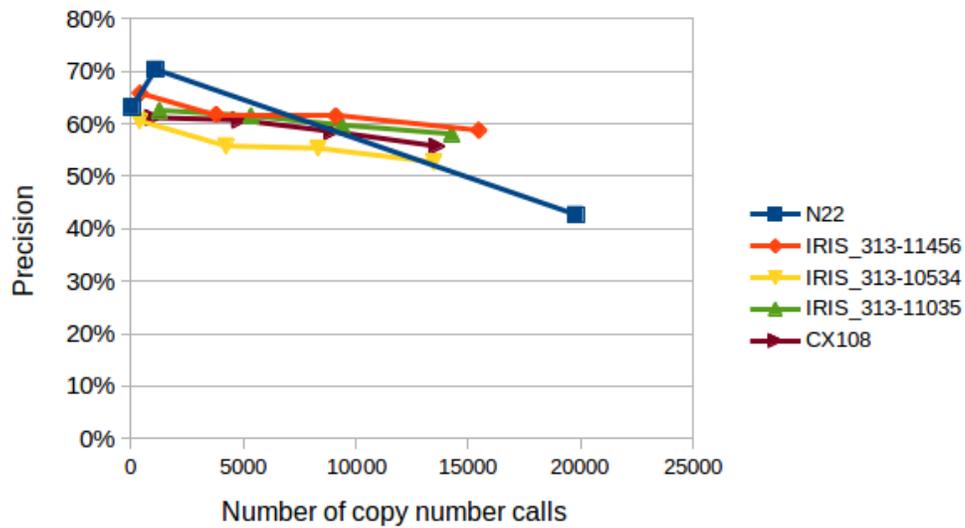


169 **Figure S18. Assessment of copy number variation prediction.** a) Read depth distribution for the Aus
 170 samples N22 and IRIS_313_11456 sequenced at about 12× average read depth. b) Number of copy
 171 number predictions and precision estimated by BLAST searches to the N22 assembly for five Aus
 172 samples sequenced at different average read depths.

173 **a**



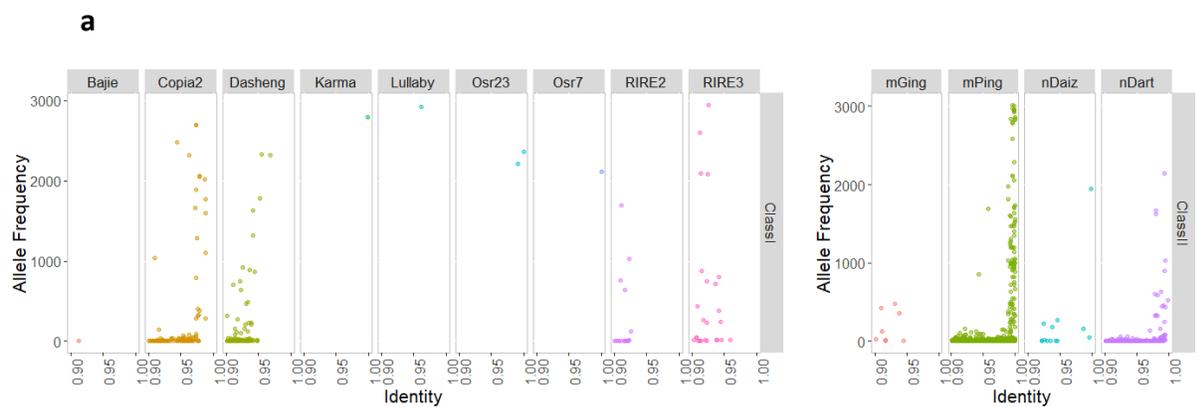
183 **b**



186 **Figure S19. Comparison of deletion and insertion sequences to known or potentially active TEs.** Each
 187 dot represents a cluster or an SV event. **a)** For deletions, extremely high or low- frequency TEs may
 188 indicate recent activity depending on whether they are retrotransposons (Class I; copy and paste)
 189 or DNA transposons (Class II; cut and paste). **b)** For insertions, cases with low allele frequency and
 190 high identity may indicate recent events. SV breakpoints do not always have precise breakpoints,
 191 hence lower sequence identity is expected especially for longer events. **c)** Insertional preference of
 192 INS events that matched active TE families.

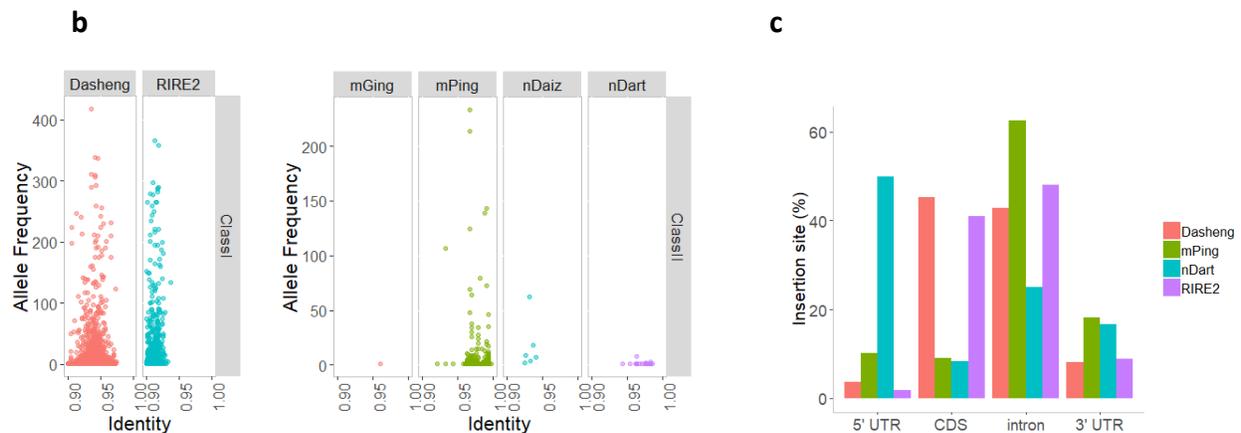
193

194



195

196



197

198

199 **Supplemental Tables**

200 **Table S1.** Transposable and repeat elements (size >50bp) among the structural variants.

Order	Superfamily	SV type (<i>no. observations / no. events</i>)			
		Deletion	Insertion*	Duplication	Inversion
Class I (Retrotransposons)					
LTR	Copia (RLC)	1,015,423 / 2,775	11,611 / 342	125,049 / 791	5,892 / 516
	Gypsy (RLG)	2,814,736 / 7,749	59,524 / 2,772	287,159 / 2,091	11,485 / 995
	ERV (RLE)	12,442 / 158	0	10 / 4	1,521 / 15
LINE	L1 (RIL)	187,686 / 2,105	58 / 5	1,738 / 18	221 / 42
SINE	- (RSU)	57,580 / 1,780	9,546 / 133	432 / 24	43 / 21
Class II (DNA transposons)					
TIR	Tc1-Mariner (DTT)	638,513 / 19,293	92,620 / 2,248	3,765 / 74	7,269 / 2,271
	hAT (DTA)	601,502 / 8,224	28,791 / 661	2,052 / 32	1,326 / 378
	Mutator (DTM)	1,445,375 / 16,853	106,758 / 2,699	1,880 / 67	16,362 / 1,455
	Harbinger (DTH)	1,864,512 / 36,689	305,963 / 6,899	2,358 / 91	5,033 / 1,482
	CACTA (DTC)	620,698 / 2,181	5,232 / 334	1032 / 41	1,133 / 129
	Micropon (DTN)	28,803 / 1,716	81 / 14	3 / 1	3 / 3
	MITE (DTI)	186,273 / 3,065	0	295 / 44	2,917 / 415
Helitron	Helitron (DHH)	316,409 / 2,155	1,093 / 48	2,227 / 73	890 / 98
Structural Repeat					
	Tandem repeats (SRR)	0	2 / 1	7 / 2	0
	Centromeric repeats (SRC)	1,287 / 81	3 / 3	1,334 / 71	266 / 48
	Telomeric repeats (SRT)	123 / 29	0	143 / 32	13 / 5
	Satellite (SRS)	0	0	0	0
	Simple repeats (SSS)	85,653 / 1,979	0	503 / 21	194 / 55
	Low Complexity (SSL)	1,776 / 151	0	90 / 8	50 / 18
	Total	9,878,791 / 106,983	621,282 / 16,159	430,077 / 3,485	54,618 / 7,946

201 * MetaSV insertions detected in 562 high-coverage samples

202

203

204 **Table S2.** Strategies and SV types supported by each caller.

205

Caller	Paired End	Split read	Read depth	<i>de novo</i> Assembly	Variant Types*
Delly	✓	✓			DEL,INS,DUP,INV,TRA
GROM	✓	✓	✓		DEL,INS,DUP,INV,TRA
GROM-RD			✓		DEL,DUP
Lumpy	✓	✓			DEL,INS,DUP,INV
MetaSV	✓	✓	✓	✓	DEL,INS
MindTheGap				✓	INS
NGSep	✓	✓	✓		DEL,DUP,INS,INV
Pindel	✓	✓			DEL,INS,DUP,INV
Scalpel				✓	DEL,INS,DUP,INV
Socrates		✓			DEL,INS,DUP,INV,TRA

206

* DEL – deletion, INS – insertion, DUP – duplication, INV – inversion, TRA - translocation

207

208

209 **Table S3.** Summary of structural variants (SV) identified by NGSEP on the complete dataset combining
210 read depth (RD), read pair (RP) and split read (SR) approaches.

211

Event	Analysis	Total events (Millions)	Average events per sample
Deletion	RD CNVnator	23.18	7,741.02
Deletion	RD EWT	29.91	9,985.46
Deletion	RP+SR	55.12	18,405.20
Duplication	RD CNVnator	18.18	6,070.59
Duplication	RD EWT	30.28	10,111.79

212

213

214 **Table S4.** Manually validated SVs using N22 (CX368) reference genome.

	Number of validated SV	T	FP	PT
Deletion	100	70	14	16 ^a
Duplication	40	16	16	8 ^b
Inversion	20	5	15 ^c	0 ^c

215 ^a Partially True: Inaccurate breakpoint, longer or complex event

216 ^b Transposable element

217 ^c Includes inverted repeats and palindromes

218

219

220

221 **Supplemental Methods**

222 **Structural variant discovery**

223 Structural variants can be classified into the following types: deletions, insertions, duplications
224 (tandem and interspersed), inversions, and translocations, and there are five general strategies to detect
225 SVs based on analysis of data from high throughput DNA sequencing technologies: paired-end mapping
226 (RP), split-read mapping (SR), read depth (RD), *de novo* assembly (AS), and a combination of the above
227 approaches (CB). Each of these strategies have different strengths and weaknesses in detection,
228 depending on variant type, sequence length and reference genome quality; hence, applying
229 complementary methods and combining results can overcome some of the limitations inherent to the
230 different approaches (Alkan et al. 2011). To improve performance and reduce false detections, we
231 combined multiple algorithms into a pipeline for SV discovery.

232 Breakdancer (Chen et al. 2009) and GASV (Sindi et al. 2009) are two examples of RP algorithms that
233 analyze relatively discordantly aligned paired reads, while an SR method called Socrates (Schröder et al.
234 2014) considers cases of gapped or broken read alignments to the reference genome. Both RP and SR
235 methods can be used to discover all classes of SVs, but SR has the advantage of single-nucleotide
236 resolution (Alkan et al. 2011). RD methods like GROM-RD (Smith et al. 2015) and CNVnator (Abyzov et al.
237 2011) involve counting reads in windows and segmenting counts; although these RD approaches are able
238 to detect both tandem and interspersed duplications, they only detect losses (deletions) and gains
239 (duplication), rely on good read coverage, and have poor breakpoint resolution (Francia et al. 2015; Tattini
240 et al. 2015). On the other hand, some SV callers that combine multiple signals such as GROM (Smith et al.
241 2017), Pindel (Ye et al. 2009), Delly (Rausch et al. 2012), NGSEP (Duitama et al. 2014) and Lumpy (Layer et
242 al. 2014) provide high precision and sensitivity for insertions, deletions, inversions and tandem
243 duplications but not all are able to identify interspersed duplications.

244 Most of the previous methods mentioned rely on read alignment to a single reference genome; hence,
245 for rice, long novel insertion events relative to the Nipponbare reference genome (IRGSP 1.0) cannot be
246 detected (Schatz et al. 2014). The extensive diversity of rice landraces and varieties exemplified by deep
247 population structure and differences in genome sizes will also impact the ability to detect SVs accurately.
248 *De novo* assembly methods such as Scalpel (Narzisi et al. 2015) and ScanIndel (Yang et al. 2015), assemble
249 reads from suspected breakpoints into contigs and then align these contigs to a reference genome to
250 determine the exact boundaries. AS methods can detect all types of SVs at nucleotide resolution but are
251 computationally intensive and not applicable to low-coverage samples. MindTheGap (Rizk et al. 2014), a
252 variant of AS, implements the Bloom filter, a probabilistic space-efficient data structure, to create a
253 probabilistic de Bruijn graph and requires less memory compared to other AS methods. To further
254 enhance the efficiency of structural variant discovery, newer solutions like SVMerge (Wong et al. 2010)
255 and MetaSV (Mohiyuddin et al. 2015) combine predictions from multiple callers, filter pooled results and
256 refine breakpoints using local assembly. Apart from the more common methods presented above,
257 another tool named forests (Michaelson and Sebat 2012) implements machine learning trained with
258 experimentally-validated structural variant calls, but is currently limited to the human genome.
259 Supplemental Table S2 provides the summary of strategies used and the supported SV types per caller.

260

261 **Comparison of structural variant finding programs**

262 To identify a set of SV callers to be integrated into our discovery pipeline, we benchmarked ten callers (i.e.
263 Pindel, Delly, GROM, NGSEP, GROM-RD, Lumpy, Socrates, MindTheGap, MetaSV, Scalpel) in terms of
264 precision and sensitivity (recall) of detecting different types and lengths of simulated SVs. Shown in
265 **Figure S1**, Pindel, Delly and GROM have relatively good F1-scores for detection of deletions and tandem
266 duplications. F1-score computed as $\frac{2(\text{recall} * \text{precision})}{(\text{recall} + \text{precision})}$ is the harmonic average of computed precision and
267 recall. For inversions, Pindel and Lumpy gave the best performance, while Delly and GROM discovered

268 very few events in the first two bins. Scalpel always produced very precise predictions, but was not as
269 sensitive as the others. Lumpy performed best for inversions with low sensitivity for small deletions and
270 tandem duplications.

271 Since long insertions with precise breakpoints and sequences are detected more accurately by
272 assembly-based tools, MetaSV, MindTheGap, and Scalpel were also evaluated. ScanIndel was eliminated
273 from the list due to errors when running larger bins. Only MindTheGap had >50% sensitivity on bin A while
274 MetaSV had the highest sensitivity for all the other bins. Pindel had the highest sensitivity among non-
275 assembly-based tools, but it could not determine sequences for large insertions. Currently, it is still
276 challenging to assess SV callers to run on rice genomes due to the absence of a gold standard list of
277 validated structural variants.

278

279 **Comparison with GATK**

280 We compared the insertion and deletion datasets with the indels deposited in the SNP-Seek database
281 (Mansueto et al. 2017) and found that setting 50 bp as the minimum SV size in our pipeline discarded
282 variants of sizes below 50 bp that the GATK-UG algorithm failed to detect. Reducing the minimum length
283 to 10 bp for DEL and 5 bp for INS allowed the pipeline to report calls that GATK-UG missed due to its
284 algorithmic limitation(s). **Figure S1** shows the range of deletion and insertion sizes GATK-UG and our new
285 pipeline were able to detect. From this figure, we see that GATK-UG can detect most deletions less than
286 30 bp (also see **Figure S1**) and most insertions shorter than 20 bp. For greater sizes, we need to apply
287 specialized methods designed for detecting large SVs.

288

289 **SV detection and genotyping of CNVs with NGSEP**

290 Copy number variants (CNVs) and large deletions were identified independently in the 3K RG samples for
291 both read depth (RD) and read pair/split read (RP+SR) analyses in NGSEP using default parameters. The

292 RD analysis used both the CNVnator algorithm and the EWT algorithm. Between 13,000 and 20,000 events
293 were detected on average per sample using each approach (**Table S3**).

294 The distribution of lengths for deletions and duplications called by each approach shows that the
295 CNVnator algorithm tends to call both deletions and duplications with average lengths of around 10 kbp,
296 whereas the EWT algorithm and the RP+SR analysis call mainly events of lengths below 2 Kbp
297 (**Figure S13a**). The events called by CNVnator (both deletions and duplications) cover on average 60Mbp
298 of the genome for each sample, whereas the events called by EWT cover on average 10Mbp and always
299 less than 40Mbp (**Figure S13b**). Interestingly, the amount of the genome covered by deletions called by
300 the RP+SR analysis is on average 80 Mbp of the genome. This is mainly explained by the tail of long events
301 called with this analysis (Last bar of **Figure S13a**).

302 Since the RP+SR approach is also implemented by other tools compared in the analysis, and that
303 Pindel and Delly reported better accuracy using simulated data, the rest of this analysis focused on CNVs
304 predicted with the CNVnator algorithm. The distribution of lengths and genome coverage shown in **Figure**
305 **S13** suggests that this analysis provides deletions and duplications that cannot be identified by other
306 approaches, complementing the results provided by the other tools. It is well known that the accuracy of
307 the RD analysis is mainly affected by the average read depth at which each sample is sequenced, the
308 distribution of reads across the genome, and the amount of repetitive content within each particular
309 region (Teo et al. 2012; Duitama et al. 2014). To account for the first two issues, a subset of 938 samples
310 was selected if the average read depth was 15× and if the total genome covered by duplications was below
311 100 Mb. The latter filter assumes that samples for which CNV predictions cover a larger percentage of the
312 genome are likely to contain more false positive calls due to an overall non-uniform distribution of reads.
313 Calls from these samples were merged into a consolidated set of 365,761 regions affected by CNVs using
314 a heuristic procedure described in Lobaton et al. (2018) for a similar study in common bean . The repetitive
315 content of each of these CNV regions was calculated and annotated to allow filtering of regions using

316 different thresholds based on the percentage of repetitive content depending on the particular
317 downstream analysis of these CNVs.

318

319 Copy number for each CNV region for each of the 3K RG samples was based on the read-depth within the
320 region as compared to the average read-depth across the genome for a particular sample taking into
321 account its variance (see Lobaton et al. (2018) for details). This procedure led to a matrix of predictions of
322 copy number, which we term CNV genotype calls, having as many rows as CNVs and as many columns as
323 analyzed samples. A dataset of 669 million genotype calls with minimum genotyping quality score of 10
324 was assembled, having a percentage of missing data of 39.5%. If only the 938 samples selected for the
325 merging step are taken into account, the number of genotype calls reduces to 243 million, with a
326 percentage of missing data of 29.1%. This dataset represents a raw catalog of copy number variation
327 events that, similar to the database of SNPs (Mansueto et al. 2017) can be filtered in different ways
328 depending on the expected level of precision and sensitivity and on the desired downstream analysis. For
329 example, since the confounding effects produced by misalignments around repetitive regions is the main
330 source of false positive predictions, application of stringent filtering removed regions that overlap by even
331 1bp with repetitive regions and reduced the dataset to 5,351 CNV regions and 8.7 million genotype calls
332 over the 3K RG samples having a percentage of missing data of 46.3%. Within the samples selected for
333 the merging step, the number of genotype calls becomes 3.2 million (36.4% missing data).

334

335 **Validation of CNVs**

336 We performed two different approaches for validation of CNVs. First, we used predictions of copy number
337 variation as alleles of genetic markers and performed clustering of samples provided by the structure
338 software (Pritchard et al. 2000) on different datasets obtained by applying filters of number of samples
339 genotyped, intersection with repeat regions and quality score of the SV calls (see results in the main text).

340 Second, we tried to validate the predictions of copy number of the sample N22 by performing BLAST
341 searches in the recently published *de novo* assembly of this sample and checking if the number of copies
342 found in the assembly corresponds with the predictions copy number based on RD. Unfortunately, the
343 read depth distribution for this sample deviated severely from a normal distribution (**Figure S19a**). We
344 followed this procedure with other Aus samples assuming that the Aus samples would share a common
345 core of structural variation relative to Nipponbare. The precision obtained with the raw dataset of calls
346 was about 20%. However, we noticed that most false positive calls were those with predictions of copy
347 number equal to 1 (heterozygous deletion) or 3 (heterozygous duplication). Hence, we decided to filter
348 out these calls and then remove CNVs that do not have predictions of copy number below 1 or above 3
349 on the 938 selected samples. This filter produced the final dataset of CNVs for this study having 207,927
350 CNVs. **Figure S19b** shows that the estimated precision in this dataset increases close to 60%, even taking
351 into account that the assembly used for comparison corresponds to a sample different from the tested
352 samples. Further filtering based on quality score only reduced recall without increasing precision.

353

354 **References**

- 355 Abyzov A, Urban AE, Snyder M, Gerstein M. 2011. CNVnator: An approach to discover, genotype, and
356 characterize typical and atypical CNVs from family and population genome sequencing. *Genome*
357 *Research* **21**: 974-984.
- 358 Alkan C, Coe BP, Eichler EE. 2011. Genome structural variation discovery and genotyping. *Nature reviews*
359 *Genetics* **12**: 363-376.
- 360 Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, McGrath SD, Wendl MC, Zhang Q, Locke
361 DP et al. 2009. Breakdancer: an algorithm for high resolution mapping of genomic structural
362 variation. *Nature Methods* **6**: 677-681.
- 363 Duitama J, Quintero JC, Cruz DF, Quintero C, Hubmann G, Foulquié-Moreno MR, Verstrepen KJ, Thevelein
364 JM, Tohme J. 2014. NGSep: An integrated framework for discovery and genotyping of genomic
365 variants from high-throughput sequencing experiments. *Nucleic acids research* **42**: e44-e44.
- 366 Francia E, Pecchioni N, Policriti A, Scalabrin S. 2015. CNV and Structural Variation in Plants: Prospects of
367 NGS Approaches. doi:10.1007/978-3-319-17157-9, pp. 211-232. Springer International
368 Publishing.
- 369 Layer RM, Chiang C, Quinlan AR, Hall IM. 2014. LUMPY: A probabilistic framework for structural variant
370 discovery. *Genome biology* **15**: R84-R84.

371 Lobaton JD, Miller T, Gil J, Ariza D, de la Hoz JF, Soler A, Beebe S, Duitama J, Gepts P, Raatz B. 2018.
372 Resequencing of Common Bean Identifies Regions of Inter-Gene Pool Introgression and Provides
373 Comprehensive Resources for Molecular Breeding. *Plant Genome* **11**.

374 Mansueto L, Fuentes RR, Borja FN, Detras J, Abriol-Santos JM, Chebotarov D, Sanciangco M, Palis K,
375 Copetti D, Poliakov A et al. 2017. Rice SNP-seek database update: new SNPs, indels, and queries.
376 *Nucleic acids research* **45**: D1075-D1081.

377 Michaelson JJ, Sebat J. 2012. forestSV: structural variant discovery through statistical learning. **9**: 819-
378 821.

379 Mohiyuddin M, Mu JC, Li J, Bani Asadi N, Gerstein MB, Abyzov A, Wong WH, Lam HYK. 2015. MetaSV: an
380 accurate and integrative structural-variant caller for next generation sequencing. *Bioinformatics*
381 (*Oxford, England*) doi:10.1093/bioinformatics/btv204: 1-4.

382 Narzisi G, Rawe JAO, Iossifov I, Lee Y-h. 2015. Accurate de novo and transmitted indel detection in exome-
383 capture data using microassembly. *Nature Methods* **11**: 1033-1036.

384 Pritchard JK, Stephens M, Donnelly P. 2000. Inference of population structure using multilocus genotype
385 data. *Genetics* **155**: 945-959.

386 Rausch T, Zichner T, Schlattl A, Stütz AM, Benes V, Korbel JO. 2012. DELLY: structural variant discovery by
387 integrated paired-end and split-read analysis. *Bioinformatics (Oxford, England)* **28**: i333-i339.

388 Rizk G, Gouin A, Chikhi R, Lemaitre C. 2014. MindTheGap : integrated detection and assembly of short and
389 long insertions. *Bioinformatics*: 1-7.

390 Schatz MC, Maron LG, Stein JC, Wences AH, Gurtowski J, Biggers E, Lee H, Kramer M, Antoniou E, Ghiban
391 E et al. 2014. Whole genome de novo assemblies of three divergent strains of rice, *Oryza sativa*,
392 document novel gene space of aus and indica. *Genome Biology* **15**: 506-506.

393 Schröder J, Hsu A, Boyle SE, Macintyre G, Cmero M, Tohill RW, Johnstone RW, Shackleton M, Papenfuss
394 AT. 2014. Socrates: identification of genomic rearrangements in tumour genomes by re-aligning
395 soft clipped reads. *Bioinformatics (Oxford, England)* doi:10.1093/bioinformatics/btt767: 1-9.

396 Sindi S, Helman E, Bashir A, Raphael BJ. 2009. A geometric approach for classification and comparison of
397 structural variants. *Bioinformatics* **25**: 222-230.

398 Smith SD, Kawash JK, Grigoriev A. 2015. GROM-RD: resolving genomic biases to improve read depth
399 detection of copy number variants. *PeerJ* **3**: e836-e836.

400 Smith SD, Kawash JK, Grigoriev A. 2017. Lightning-fast genome variant detection with GROM. *GigaScience*
401 **6**: 1-7.

402 Tattini L, D'Aurizio R, Magi A. 2015. Detection of Genomic Structural Variants from Next-Generation
403 Sequencing Data. *Frontiers in bioengineering and biotechnology* **3**: 92-92.

404 Teo SM, Pawitan Y, Ku CS, Chia KS, Salim A. 2012. Statistical challenges associated with detecting copy
405 number variations with next-generation sequencing. *Bioinformatics* **28**: 2711-2718.

406 Wong K, Keane TM, Stalker J, Adams DJ. 2010. Enhanced structural variant and breakpoint detection using
407 SVMerge by integration of multiple detection methods and local assembly. *Genome biology* **11**:
408 R128-R128.

409 Yang R, Nelson AC, Henzler C, Thyagarajan B, Silverstein KaT. 2015. ScanIndel: a hybrid framework for
410 indel detection via gapped alignment, split reads and de novo assembly. *Genome medicine* **7**: 127-
411 127.

412 Ye K, Schulz MH, Long Q, Apweiler R, Ning Z. 2009. Pindel: a pattern growth approach to detect break
413 points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics*
414 (*Oxford, England*) **25**: 2865-2871.

415