1 # Pan-genome inversion index reveals evolutionary insights

2 # into the subpopulation structure of Asian rice (*Oryza sativa*)

3

4 Yong Zhou[1a], Zhichao Yu[2a], Dmytro Chebotarov[3a], Kapeel Chougule[4a], Zhenyuan

5 Lu[4], Luis F. Rivera[1], Nagarajan Kathiresan[5], Noor Al-Bader[1], Nahed Mohammed[1],

6 Aseel Alsantely[1], Saule Mussurova[1], João Santos[1], Manjula Thimma[1], Maxim

7 Troukhan[6], Alice Fornasiero[1], Carl D. Green[7], Dario Copetti[8], Dave Kudrna[8], Victor

8 Llaca[9], Mathias Lorieux[10], Andrea Zuccolo[1,11*], Doreen Ware[4,12*], Kenneth

9 McNally[3*], Jianwei Zhang[2,8*], Rod A. Wing[1,3,8*]

10

11 [1]Center for Desert Agriculture (CDA), Biological and Environmental Sciences &

12 Engineering Division (BESE), King Abdullah University of Science and Technology

13 (KAUST), Thuwal, 23955-6900, Saudi Arabia

14 [2]National Key Laboratory of Crop Genetic Improvement, Huazhong Agricultural

15 University, Wuhan 430070, China

16 [3]International Rice Research Institute (IRRI), Los Baños, 4031 Laguna, Philippines

17 [4]Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA

18 [5]Supercomputing Core Lab, King Abdullah University of Science and Technology

19 (KAUST), Thuwal, 23955-6900, Saudi Arabia

20 [6]Persephone Software, LLC, Agoura Hills, California 91301, USA

21 [7]Information Technology Department, King Abdullah University of Science and

22 Technology (KAUST), Thuwal, 23955-6900, Saudi Arabia

23 [8]Arizona Genomics Institute, School of Plant Sciences, University of Arizona,

24 Tucson, Arizona 85721, USA

25 [9]Research and Development, Corteva Agriscience, Johnston, Iowa 50131, USA

26 [10]DIADE, University of Montpellier, CIRAD, IRD. Montpellier, France

27 [11]Institute of Life Sciences, Scuola Superiore Sant'Anna, Pisa, 56127, Italy

28 [12]USDA ARS NEA Plant, Soil & Nutrition Laboratory Research Unit, Ithaca, NY,

29 14853, USA

30

31 Yong Zhou[1a], yong.zhou@kaust.edu.sa, ORCID 0000-0002-1662-9589

32 Zhichao Yu[2a], proyu@webmail.hzau.edu.cn, ORCID 0000-0003-2155-4830

33 Dmytro Chebotarov[3a], d.chebotarov@irri.org, ORCID 0000-0003-1351-9453

34    Kapeel Chougule[4a], kchougul@cshl.edu, ORCID 0000-0002-1967-4246

35    Zhenyuan Lu[4], luj@cshl.edu, ORCID 0000-0003-1758-2636

36    Luis F. Rivera[1], luis.riveraserna@kaust.edu.sa, ORCID 0000-0003-3978-7640

37    Nagarajan Kathiresan[5], nagarajan.kathiresan@kaust.edu.sa,

38    Noor Al-Bader[1], noor.albader@kaust.edu.sa, ORCID 0000-0002-0511-6972

39    Nahed Mohammed[1], nahed.mohammed@kaust.edu.sa, ORCID 0000-0002-8857-

40    3246

41    Aseel Alsantely[1], aseel.alsantely@kaust.edu.sa, ORCID 0000-0002-4990-249X

42    Saule Mussurova[1], saule.mussurova@kaust.edu.sa,

43    João Santos[1], dourado.jns@gmail.com,

44    Manjula Thimma[1], manjula.thimma@kaust.edu.sa, ORCID 0000-0002-1703-8780

45    Maxim Troukhan[6], mtroukhan@persephonesoft.com, ORCID 0000-0002-2671-6002

46    Alice Fornasiero[1], alice.fornasiero@kaust.edu.sa, ORCID 0000-0001-6165-4233

47    Carl D. Green[7], carl.green@kaust.edu.sa, ORCID 0000-0003-2952-3908

48    Dario Copetti[8], dcopetti@email.arizona.edu, ORCID 0000-0002-2680-2568

49    Dave Kudrna[8], dkudrna@email.arizona.edu, ORCID 0000-0002-3092-3629

50    Victor Llaca[9], victor.llaca@corteva.com, ORCID 0000-0003-4822-2924

51    Mathias Lorieux[10], mathias.lorieux@ird.fr, ORCID 0000-0001-9864-3933

52    Andrea Zuccolo[1,11*], andrea.zuccolo@kaust.edu.sa, ORCID 0000-0001-7574-0714

53    Doreen Ware[4,12*], Doreen.ware@usda.gov, ORCID 0000-0002-8125-3821

54    Kenneth McNally[3*], k.mcnally@irri.org, ORCID 0000-0002-9613-5537

55    Jianwei Zhang[2,8*], jzhang@mail.hzau.edu.cn, ORCID 0000-0001-8030-5346

56    Rod A. Wing[1,3,8*], rod.wing@kaust.edu.sa, rwing@ag.arizona.edu, ORCID 0000-

57    0001-6633-6226

58

59    [a]These authors contributed equally to this work.

60    [*]Correspondence and requests for materials should be addressed to: Andrea Zuccolo

61    (email: andrea.zuccolo@kaust.edu), Doreen Ware (email: Doreen.ware@usda.gov),

62    Kenneth McNally (email: k.mcnally@irri.org), Jianwei Zhang (email:

63    jzhang@mail.hzau.edu.cn), or Rod A. Wing (email: rod.wing@kaust.edu.sa;

64    rwing@ag.arizona.edu).

65    **Abstract**

66    Understanding and exploiting genetic diversity is a key factor for the productive and

67    stable production of rice. Utilizing 16 high-quality genomes that represent the

68    subpopulation structure of Asian rice (*O. sativa*), plus the genomes of two close

69    relatives (*O. rufipogon* and *O. punctata*), we built a pan-genome inversion index of

70    1,054 non-redundant inversions that span an average of ~ 14% of the *O. sativa* cv.

71    Nipponbare reference genome sequence. Using this index we estimated an inversion

72    rate of 1,100 inversions per million years in Asian rice, which is 37 to 73 times

73    higher than previously estimated for plants. Detailed analyses of these inversions

74    showed evidence of their effects on gene regulation, recombination rate, linkage

75    disequilibrium and agronomic trait performance. Our study uncovers the prevalence

76    and scale of large inversions ($\geq$ 100 kb) across the pan-genome of Asian rice, and

77    hints at their largely unexplored role in functional biology and crop performance.

78

79    **Keywords**: Asian Rice, PSRefSeqs, Pan-genome, Inversion, Evolution

80    **Main**

81    Asian rice (*Oryza sativa*) is a staple cereal crop that has played an essential role in

82    feeding much of the world for millennia[1,2]. As the population expands to almost 10-

83    billion by 2064[3], the rice community is searching for novel ways to breed new

84    varieties that are sustainable, nutritious and climate resilient[2]. One source of the raw

85    material required to meet this urgent demand is standing natural structure variation

86    (SV), *i.e.* single nucleotide polymorphisms [SNPs], insertions/deletions [INSs/DELs],

87    translocations [TRAs], and inversions [INVs] in the genomes of the more than

88    500,000 accessions of rice and its wild relatives that have been deposited in

89    germplasm banks around the world[2].

90    Inversions are an important subset of this natural variation tool box[4-6] and have

91    been shown to play important roles in genetic recombination (*e.g. Drosophila*[7,8],

92    *Helianthus*[9], yeast[10], bacterial[11]), genome evolution (*e.g.* mouse[12], human[13,14]), and

93    speciation (*e.g. Mimulus guttatus*[15], chimps and humans[16,17]). In rice, inversions are

94    understudied and have been limited to small and mid-size inversions as a

95    consequence of the reliance on short-read data for their detection. For example, Wang

96    et al. (2018) performed a genome scan of inversions in *O. sativa* cv. Nipponbare (*i.e.*

97    IRGSP RefSeq) using re-sequencing data from 453 high-coverage genomes (> 20x )

98    from the 3K Rice Genome Project (3K-RGP) and detected $152 \pm 62$ inversions per

99    genome with a size range of $127.1 \pm 19.4$ kb[18,19]. A phylogenetic analysis of this

100    dataset, including other SV data, demonstrated that SVs can be used to define the

101    population structure of Asian rice[18]. Fuentes et al. (2019) went onto interrogate the

102    entire 3K-RGP dataset in a similar manner and identified 1,255,033 inversions, with

103    the vast majority falling in a size range of 50 bp – 500 kb[20]. A genome scan of the

104    IRGSP RefSeq, plus reciprocal genome alignment to nine Asian rice and two AA-

105    genome wild relatives (*i.e. O. rufipogon* and *O. longistaminata*) confirmed the

106    presence a previously detected ~5 megabases (Mb) inversion spanning the

107    centromere of chromosome 6 in four *Xian-indica* (*XI*) varieties, relative to four *Geng-*

108    *japonica* (*GJ*) varieties, as well as the two outgroup species[21]. A broad phylogenetic

109    study in 13 cultivated and wild *Oryza* genomes using SVs resulted in the

110    identification of 12 large inversions (*i.e.* 60-300 kb) that the authors inferred

111    potentially led to the rapid diversification of the AA genome species within a 2.5

112    million years (MY) span[22].

113    Although these studies contributed to a preliminary understanding of inversions

114     in rice, they are limited due to their reliance on short-read sequencing technology,

115     and the number and quality of genomes analyzed. Of note, a comprehensive analysis

116     of inversions that utilizes ultra-high-quality reference genome sequences, which takes

117     into account the population structure of Asian rice, remains uncharted. To reveal a

118     comprehensive understanding of large inversions ($\geq$ 100 bp) and explore their

119     evolutionary impacts in Asian rice, we used a set of 15 platinum standard reference

120     sequences (PSRefSeqs) that were sequenced with long-read sequencing technology,

121     and assembled, edited and validated with a uniform pipeline[23,24]. When combined

122     with the *O. sativa* IRGSP RefSeq[25], these data can be used as a "pan-genome" proxy

123     to represent the subpopulation structure of cultivated Asian rice, *i.e.* 15-

124     subpopulations from subgroups *Geng/Japonica* (*GJ*), *Xian/Indica* (*XI*), *circum-Aus*

125     (*cA*), *circum-Basmati* (*cB*), plus the largest admixed subpopulation, where K =15[24].

126     This Asian rice pan-genome was then scanned for inversions $\geq$ 100 bp, all anchored

127     within a phylogenetic context, using two additional *de novo* assembled (to a similar

128     quality) genomes from a representative species of the progenitor of Asian rice (*O.*

129     *rufipogon*) and the BB genome species - *O. punctata,* as outgroups.

130        In this study, we comprehensively interrogated this pan-genome dataset to detect

131     and analyze the inversion landscape of Asian rice at the population structure level, the

132     results of which revealed salient evolutionary insights into the genome biology of

133     Asian rice:

134      1.   We created a novel Asian rice pan-genome, including 16 PSRefSeqs that

135         represent its K=15 population structure, plus PSRefSeqs from two close wild

136         relatives (*O. rufipogon* and *O. punctata*).

137      2.   A pan-genome inversion index of 1,054 non-redundant inversions was

138         generated and independently validated with physical maps (*i.e.* Bionano

139         optical maps) and resequencing data (*i.e.* 3K-RGP).

140      3.   A novel "pan-genome inversion rate" was estimated at 1,100 inversions per

141         million years in Asian rice, which is 37 to 73 times higher than previous

142         estimated in plants.

143      4.   Biological functions *via* gene disruption, recombination rate, and linkage

144         disequilibrium (LD) were investigated where we found that, on average, 8

145         genes were disrupted per genome; the genome recombination rate of a RIL

146         population decreased from 6.98 to 4.00 cM/Mb; and 88.6% of the inversions

147         tested may contain traces of recombination.

148    5.  The biological consequences of a ~400 kb inversion cluster (*i.e.*
149        INV030400/INV030410/INV030420) on chromosome 3, that arose in the
150        *Xian*/*Indica* (*XI*) and *circum-Aus* (*cA*) subgroups, was shown to be under
151        positive selection, and was associated with delayed flowering with respect to
152        standard genotypes.

153 **Results**

154 **The 18-genome Data Package**

155 To investigate the inversion landscape of Asian rice from a population structure
156 perspective, we first combined a set of 16 previously published high-quality genomes
157 that represent the K=15 population structure of *O. sativa*, plus the largest *Xian*/*indica*
158 (*XI*) admixed subpopulation (*XI*-adm: Minghui 63 (MH63)) to create a "pan-genome"
159 of Asian rice. To anchor this novel pan-genome within a phylogenetic context, we
160 long-read sequenced, *de novo* assembled and validated two additional genomes from
161 both a representative species of the progenitor of Asian rice - *i.e. O. rufipogon* [AA],
162 and the African BB genome outgroup species - *O. punctata* (Table 1, Supplementary
163 Table 1, Extended Data Fig. 1, and Supplementary Note 1).
164     All Asian rice assemblies were annotated using a uniform annotation pipeline to
165 minimize methodological artifacts (Table 1, Supplementary Tables 2-4, Extended
166 Data Fig. 2, and Supplementary Note 1), except for the *XI*-adm: MH63 and *XI*-1A:
167 Zhenshan 97 (ZS97) genomes, whose annotations were previously published[23].
168 Lastly, we integrated and compared all annotations with that of the *GJ*-temp: IRGSP
169 RefSeq[25].
170     All 18 genomes, and their annotations, are henceforth referred to as the "18-
171 genome data package" (See "18-genome data package" in the Supplementary Note 1
172 section for a complete description of this data set).
173

174 **Creation of a Pan-genome Inversion Index for Asian Rice**
175 We pairwise compared 17 reference genome assemblies with the IRGSP RefSeq[25]
176 and identified a total of 2,915 inversions (≥ 100 bp) (Table 2, Supplementary note 2),
177 of which, 1,054  were non-redundant (Supplementary Table 5). As expected, more
178 inversions were observed when we compared the Asian rice pan-genome with both
179 the AA and BB genome outgroups to the IRGSP RefSeq: 194 (total length = 13.05
180 Mb) and 316 (total length = 17.85 Mb) for *O. rufipogon* and *O. punctata,* respectively
181 (Table 2). On average, each *O. sativa* genome was found to contain 160 inversions,
182 ranging from 88 (*GJ*-subtrp: CHAO MEO) to 187 (*XI*-3B1: KHAO YAI GUANG)
183 (Table 2). We found a larger number of inversions (172 to 187) when comparing the
184 *O. sativa XI*-subgroup genomes to *GJ*-temp: IRGSP RefSeq, than when comparing

185    the *O. sativa GJ*-subgroup genomes (88 to 112 inversions) to the same reference

186    (Table 2). The total length of the inverted regions, per genome, ranged from 7.73 Mb

187    (*GJ*-subtrp: CHAO MEO) to 14.95 Mb (*cA*2: NATEL BORO) (Table 2). When

188    chromosome location was taken into account, these inversions appeared to be evenly

189    distributed genome-wide (Kolmogorov-Smirnov test, *P* value 0.02-0.95)

190    (Supplementary Table 6, Extended Data Fig. 3, Supplementary note 3).

191

192    **Species and Subpopulation Specific Inversions**

193    Of the 1,054 non-redundant inversions detected (Supplementary table 5), we

194    classified them into different categories: *i.e.* species-specific (the inversion could be

195    only observed in *O. punctata*, *O. rufipogon* or *O. sativa* genomes); group-specific

196    (only observed either in *GJ*, *XI*, *cA* or *cB* subgroup in *O. sativa*); and genome-specific

197    (only observed in one of 16 *O. sativa* genomes). As a result, 968 (91.8%) appeared to

198    be species-specific (*i.e. O. sativa*: 550 (totaling 50.21 Mb), *O. rufipogon*: 105 (11.05

199    Mb), and *O. punctata*: 313 (17.84 Mb)) (Fig. 1). The remaining 86 were found in two

200    or more species and totaled to about 2.06 Mb in size.

201    Two hundred and forty-five of the 550 *O. sativa* specific inversions were specific

202    to one of the 16 *O. sativa* genomes, while 305 were shared with more than one

203    genome (Fig. 1). The frequency of 100 randomly selected *O. sativa* genome-specific

204    inversions were further investigated in different subpopulations using a subset of 3K-

205    RGP data set (*i.e.* 192 highly re-sequenced (> 20 ×) accessions) (Supplementary table

206    7). With the exception of 18% of the inversions, for lack of evidence within the 3K-

207    RGP dataset, the remaining 82% could be classified in four groups: 15% genome-

208    specific inversions, 23% subpopulation specific inversions, 33% near-subpopulation

209    specific, and 11% subpopulation shared inversions (See supplementary online

210    methods for category definitions) (Supplementary table 8).

211    Of the 305 inversions present in more than one of the 16 Asian rice genomes

212    included in our dataset, we identified 29 that were shared among closely related

213    populations (Fig. 1), which we defined as "group specific". Four inversions were

214    shared among 4 *GJ* genomes, 3 were shared among 9 *XI* genomes, 11 were shared

215    among 2 *cA* genomes, and 11 were identified by comparing 5 *GJ* and *cB* genomes to

216    11 *XI* and c*A* genomes (Fig. 1). These 29 inversions were also studied in different

217    subpopulations across a high-coverage subset of 3K-RGP dataset. Excluding two

218   inversions that couldn't be tested (*i.e.* because no reads were observed at the

219   breakpoints), we found that 11 (38%), 13 (44%) and 3 (10%) inversions were group

220   specific, near-group specific and group shared inversions at the subpopulation level,

221   respectively (Supplementary table 8). The remaining 278 inversions appear to be

222   shared across different genomes or subpopulations reflecting the substantial

223   admixture in the evolution of subpopulations in Asian rice and mixed ancestry in the

224   pedigrees of some accessions used (*e.g.* IR8, IR64, MH63, and ZS97)[26,27].

225        Altogether, > 85% of the genome specific inversions and > 85% of the *O. sativa*

226   group specific inversions could be validated with the high-coverage subset of the 3K-

227   RGP dataset, and appear to be Asian rice subpopulation(s)- or subgroup(s)-specific.

228   This analysis validates the accuracy in detecting inversion boundaries, provides initial

229   estimates of inversion frequencies in rice subpopulations, and patterns of shared

230   inversions between subpopulations.

231

232   **Five Largest Inversions**

233   We identified 5 inversions greater than 1 Mb relative to the IRGSP RefSeq on

234   chromosomes 1, 6, 8, and 10 (INV010130 [2.00 Mb], INV010560 [1.82 Mb],

235   INV060390 [4.57 Mb], INV080710 [1.12 Mb] and INV100690 [1.33 Mb]), two of

236   which (INV060390[28] and INV080710[29]) were previously reported.

237        INV010130, INV010560 and INV060390 appear to be Asian rice specific (Fig. 2,

238   Supplementary Table 5 & 9, Supplementary note 4), and could not be found in the

239   outgroup genomes. To determine if these inversions are subpopulation(s) specific, we

240   interrogated the high-coverage subset of the 3K-RGP data set and found that

241   INV010130 was specific to the *XI*-3A subpopulation, INV010560 to the c*A* and *XI*-

242   adm subpopulations, and INV060390 to the *GJ*-tmp and *GJ*-subtrp subpopulations

243   (Supplementary Table 9, Extended Data Fig. 4, Supplementary note 4).

244        The remaining two inversions (*i.e.* INV080710 and INV100690) were detected

245   only in *O. punctata* and *O. rufipogon*, respectively (Supplementary Table 5 & 9), and

246   thus appeared to be species specific. To test this hypothesis, we investigated the

247   presence or absence of these inversions in high-quality genomes of 5 additional

248   *Oryza* species (*i.e. O. nivara* [AA], *O. glaberrima* [AA], *O. barthii* [AA], and the

249   distantly related subgenomes of *O. coarctata* [KKLL] and *O. alta* [CCDD]

250   (unpublished data). Results showed that neither of these inversions could be detected

251    in these five species. Thus, we conclude that INV080710 (*O. punctata*) and

252    INV100690 (*O. rufipogon*) are species specific.

253

**254    Characterization of Transposable Element Content within Inversions and**

**255    Breakpoints**

256    Transposable elements (TEs) are known to be associated with inversions[18,20,30], thus

257    we analyzed the TE content across the inversion index, and at their breakpoints. The

258    total amount of TE related sequences within these inversions ranged from 64% (*GJ-*

259    subtrp: CHAO MEO) to 73% (*XI*-adm: MH63) (Supplementary Table 10), which is

260    significantly (student's test, $p < 0.01$) higher than the average content of TEs across

261    all 16 *O. sativa* genomes at 51.3% (Table 1, Supplementary table 10 & 11). These

262    results demonstrate that TEs are enriched within inversions.

263        Analysis of breakpoints revealed that both long terminal repeat retrotransposons

264    (LTR-RTs, *i.e.* Ty3-*gypsy* and Ty1-*copia*) and DNA TE Mutator-like elements

265    (MULEs) were significantly (student's test, $p < 0.01$) enriched, when the frequency

266    of their presence at the 2,108 breakpoints was compared to 21,080 randomly selected

267    genomic locations (*i.e.* 10 replicates) (Fig. 3a). We further studied TEs at the

268    breakpoint of each inversion that were shared across all Asian rice genomes. In doing

269    so, we identified 17 TE families (*i.e.* 13 Ty3-*Gypsy*, 1 Ty1-*Copia*, 2 *CACTA*, and 1

270    *Mutator*) present at the breakpoints of more than 10 inversions (Fig. 3b &

271    Supplementary Table 12). An example of an inversion enriched in TEs, including the

272    internal and LTR portions of at least three different LTR-RTs, is shown in Fig. 3c.

273        Together, our results reveal an enrichment of TE related sequences both within

274    inversions and at their breakpoints.

275

**276    Characterization of Gene Content within Inversions and Breakpoints**

277    Based on the pan-genome inversion index we identified a total of 15,530 genes

278    (~1,035/genome) within or at inversion breakpoints (Supplementary Table 13). To

279    investigate the effect of inversions on the expression of genes located within inverted

280    regions, we interrogated a transcriptome dataset derived from a subset of the 18-

281    genome data packaged including *O. sativa cv*. *XI*-adm: MH63, *XI*-1A: ZS97 and *GJ-*

282    temp: Nipponbare (*i.e.* dataset#2 - see online methods). Based on 284 and 356

283    expressed orthologous genes between the reference (*GJ*-temp: IRGSP RefSeq) and

284    two queries (*XI*-adm: MH63 and *XI*-1A: ZS97), we detected 10.9% (31) genes from

285    *XI*-adm: MH63 and 7.3% (26) from *XI*-1A: ZS97 that were differentially expressed

286    (DEG, fold change > 2, *P* value < 0.01) (Supplementary table 14) relative to the *O.*

287    *sativa cv. GJ*-temp: Nipponbare genome.

288        To investigate the effect of inversions on the transcription of genes located at

289    inversion breakpoints - *i.e.* about 55 genes per genome (Supplementary table 13), we

290    interrogated both our baseline RNA-Seq dataset (dataset#1- see online methods) and

291    dataset#2 for changes in transcript abundance. On average, 28 of the 55 genes per

292    genome were found to be expressed in the tissues tested (Supplementary table 14). Of

293    these, transcript abundance of an average of 20 genes per genome did not change due

294    to the presence of duplicated genes at both ends of their inversion breakpoints

295    (Supplementary table 14). An example of this observation is represented by two

296    *OsNAS* genes (*NAS1* and *NAS2*) located at the breakpoint of INV030200 (~4.3 kb)

297    (Fig. 4a & a). The remaining ~8 genes/genome were single copy and were disrupted

298    by the inversion events, leading to the absence of transcript evidence (Supplementary

299    table 14). As an example, transcripts of the Nipponbare Fbox gene (Os11g0532600)

300    can be detected in the four tissues tested. However, the first exon of this gene is

301    disrupted in MH63 by INV110960, resulting in transcript ablation (Fig. 4C & D).

302

**Recombination Rate and Genomic Inversions**

304    To evaluate the effect of inversions on recombination frequency, a previously

305    published recombinant inbred line (RIL-10) population of 210 inbred lines[31] derived

306    from a cross between *O. sativa cv. XI*-adm: MH63 and *XI*-1A: ZS97 was

307    investigated. We detected 78 inversions between MH63 and ZS97, totaling 3.58 Mb

308    and 3.51 Mb based on the MH63RS2 and ZS97RS2 genome assemblies, respectively

309    (Supplementary table 15). The recombination rate along each chromosome was

310    assessed by comparing genetic and physical distances between neighboring bins. The

311    average recombination rate for each chromosome ranged from 5.95 (chromosome 6)

312    to 9.92 (chromosome 12) cM/Mb, and varied from 0 to 153.93 cM/Mb across the

313    genome with an average of 6.98 cM/Mb (Extended Data Fig. 5A). The average

314    recombination rate over the 78 inverted regions was 4.00 cM/Mb (0 - 23.26 cM/Mb),

315    which is significantly lower (Student's t-test, $p = 0.0002$) than that observed genome-

316    wide (Extended Data Fig. 5B). These results indicate that a marked suppression of

317    genetic recombination is associated with inversions.

318

**Effect of Large Inversions on Population SNP Variation**

The occurrence of inversions can affect DNA polymorphism at the population level in several ways, including increased divergence in the inverted region and changes in linkage disequilibrium (LD) patterns[32]. The latter is particularly interesting as it can affect SNPs that are mapped to positions megabases apart, and can be a confounding factor in LD-based analyses. To determine whether large *O. sativa* inversions left a trace in patterns of LD along the IRGSP RefSeq, we used the 3K-RGP dataset to examine LD blocks near inverted regions (> 100 kb). First, inversions having a reciprocal overlap of more than 80% of their length were clustered and considered as putative unique inversions. In doing so, we considered 53 clusters including from 1 to 6 inversions each (Supplementary Table 16). An inversion fixed in a population may lead to the disruption of LD blocks, in which some SNPs flanking the inversion on one side are in LD with SNPs on the distal part of the inversion, but not on the adjacent part (Fig. 5), due to the reversed order of SNPs inside the inverted region in samples that carry the inversion allele. By an LD block we mean only a set of SNPs in high LD ($r^2$ > 0.8 in this analysis).

Next, we examined the entire 3K-RGP variation data set and searched for LD blocks that connect the flanking regions of inversions, having no SNPs in the proximal parts of each inversion. Such blocks (Fig. 5), were found in nearly all large inversions (63 out of 81 [75.3%] alignment-based inversions, or 47 out of 53 [88.7%] inversion clusters) (Supplementary Table 16) with only two classes of exceptions: *i.e.* inversions in regions of complex chromosomal rearrangements (INV080210-INV080250, INV080510-INV080530, INV110600-INV110660), and three putative "recent" inversions (INV020230, INV100080, INV100320), each of which were found in single genomes and may lack sufficient frequencies in a population to contain traces of recombination. Some of the disrupted LD blocks contained a particularly large number of SNPs and were seen as a distinctive checkered pattern on LD heat maps (Fig. 5). This comparatively large number of SNPs along with low haplotype diversity, despite the presence of recombination, could be a consequence of selective pressure.

**Phenotypic consequences of inversions: Inversion Cluster 92**

To investigate the phenotypic consequences of inversions (if any) on agronomic traits in rice, we correlated ten known phenotypes catalogued in SNP-Seek (https://snp-

353   seek.irri.org/) across the high-coverage subset of the 3K-RGP data set (*i.e.* same

354   subset as mentioned above) with our pan-genome inversion index. In this analysis,

355   we used a linear model function in R to assess association between phenotype and

356   inversion status, controlling for population structure, and in doing so, we identified an

357   inversion cluster (*i.e.* INVCluster92), 400 kb in length, which was significantly

358   (linear model test, $p < 0.01$) associated with delayed flowering time of 13 days, on

359   average (Fig. 6, Supplementary table 17).

360        INVCluster92 was found to be composed of three inversions (*i.e.*

361   INV030400/INV030410/INV030420), with a minor difference at the breakpoints

362   (Supplementary table 5), and is shared among twelve of the sixteen genomes in our

363   Asian rice pan-genome data set (*i.e.* 9 *XI*, 2 *cA* and 1 *GJ* genomes (*GJ*-trop1:

364   Azucena)) (Fig. 6a). Analysis of the 3K-RGP subset revealed that 137 (71.4%)

365   contained the inversion (INV) genotype, while 55 (28.6%) did not (*i.e.* standard

366   genotype (STD)) (Fig. 6a). Comparative LD analysis of STD vs. INVCluster92

367   genotypes is shown in Fig. 6b, and showed that the LD block of STD genotypes was

368   disrupted by INVCluster92. In addition, the INV genotypes showed a significantly

369   higher (student's test, $p < 0.01$) nucleotide diversity ($\pi$), Watterson's theta ($\theta$), and

370   Tajima's *D* than the STD genotypes (Fig. 6c). This evidence supports the hypothesis

371   that the regions immediately surrounding INVCluster92 may be under positive

372   directional selection.

373        To find clues as to why INVCluster92 is associated with delayed flowering time

374   (Fig. 6d, left), we compared gene content, structure and expression differences for

375   eight genes present in both the STD and INV genotypes. Of note, four of these genes

376   (*i.e.* LBD37[33], EXO70H2[34], WD40-76[35], shl4[36]) have been functionally characterized

377   (Supplementary table 18), and 4 are unknown (Fig. 6e). Among these eight genes, a

378   total of 3 SNPs and a 12-bp insertion were observed within the coding sequence

379   (CDS) regions of two genes - LBD37 (Os03g0445700) and WD40-76

380   (Os03g0448600). For LBD37, a single 'G' to 'T' SNP (SNP-1) was observed at base

381   pair 487 resulting into a non-conservative amino acid change from glycine to cysteine

382   (Fig. 6f). For gene WD40-76, 4 'GCC' repeats were inserted (in frame) 6 bp after the

383   start codon, resulting in the addition of 4 alanine residues at the beginning of the

384   protein. In addition, two additional SNPs were detected in this gene, a synonymous

385   'A' to 'G' SNP (SNP-2) at position 468 bp and an 'A' to 'C' SNP (SNP-3) at position

386   557bp that changed a charged histidine amino acid into a non-polar proline amino

387   acid (Extended Data Fig. 6). We also validated these SNPs across the 3K-RGP high-

388   coverage subset and found that all these natural variants were absent in all (55) STD

389   genotypes (Supplementary table 19). However, SNP-1 and the 12 bp InDel were

390   present in 97.8%, SNP-2 was found in 90.5%, and SNP3 was found in all (137) INV

391   genotypes.

392       Preliminary transcript abundance analysis of the eight genes were compared using

393   deep RNA-Seq data from leaves, roots and mixed stage panicles (*i.e.* RNA

394   dataset#2). The only difference in transcript abundance that could be observed was

395   for gene LBD37 in panicle and young leaves, respectively. The LBD37 gene

396   appeared to be up-regulated in panicle tissue, but down-regulated in young leaf

397   tissues in INV genotypes, as compare with STD genotypes (Fig. 6g and

398   Supplementary table 18), which is compatible with the inversion genotype and SNP-

399   1. The phenotype variation analysis based on SNP-1 is also congruent with LBD37

400   over-expression (Fig. 6d, right), as previously reported in rice, *i.e.* a delay in heading

401   date[33]. These results suggest that SNP-1 within LBD37 is under positive selection and

402   may contribute to the observed phenotypic variation.

## Discussion

Inversions are an important class of structural variations that have been shown to play important roles in the suppression of recombination that can lead to the selection of adaptive traits, reproductive isolation and eventual speciation, and are quite common in plants[22,32,37]. For example, over the 50-60 MY history of the *Poaceae*, where gene order has been largely conserved, Ahn and Tanksley (1993) showed (using molecular genetic maps) that multiple inversions and translocations occurred during the evolution of maize and rice from a common ancestor[38].

Here, we present, to our knowledge, the first comprehensive analysis of the inversion landscape of any cereal at the population structure level with the discovery of 1,054 non-redundant inversions that range from 8 Mb to 25 Mb in cumulative size (Table 2). It is estimated the AA genomes of the *Oryza* diverged from the BB genome type about 2.5 million years ago (MYA)[2], which equates to an inversion rate of 63.2 inversions per million years (*i.e.* 316 inversions/ (2*2.5 MY)) - about 2 to 4 times higher than that recently estimated in plants (*i.e.* 15 to 30 inversions/MY)[32]. However, Huang and Rieseber[32] (2020) noted that this earlier estimate should be considered an underestimate and is dependent on the quality of the genomes analyzed, and other factors. If we use the implied AA genome diversification rate of ~0.50 net new species/million years[2,22], then we calculate an inversion rate of 194 *O. sativa* inversions/MY - about 6.5 to 13 times higher than recently estimated in plants[32]. Taken one step further, by taking into account that Asian rice is estimated to have been domesticated 10,000 years ago[2], and using only divergence within GJ group genomes not shared by any other groups (22 out of 88 inversions in KETAN NANGKA) which is likely post domestication, we can arrive at an estimated inversion rate > 1,100 inversions/MY (*i.e.* > 37 to 73 times previous estimates in plants[32]). Such a rapid pan-genome inversion rate over such a short time period may be reflective of high fixation rates of rearrangements in plants[6], high chromosomal evolution rates in annual plants[39-41], and intense human selection since the dawn of agriculture[18,20,42].

Although an inversion genome scan for rice has been previously published[20], when we extracted the inversion coordinates (*i.e.* 2,402 inversions, average length 43.3 kb), from the same set of 15 accessions used to generate our pan-genome mapped to the IRGSP-RefSeq, we found that only 200 could be validated with dot plots (a 91.7% false positive rate) (Supplementary Table 20), 194 of which

437    overlapped with our inversion index. The 6 remaining contained 2 that overlapped,

438    and only 4 that were not present in our inversion index (Supplementary Table 21).

439    These analyses combined reveal the limitations of inversion callers with short read

440    data and provide a cautionary note as to the validity of many of the inversions

441    catalogued to date.

442        Several key factors led to our ability to generate a definitive gene inversion index

443    for cultivated Asian rice. The first was our use of a set of 16-ultra high-quality

444    reference genomes that represented the K = 15 population structure of Asian rice[24],

445    and 2 phylogenetically anchored wild AA and BB genome species (Table 1 and

446    Supplementary Table 1). Secondly, we did not computationally collapse this 18-

447    genome data package into pan-genome (*e.g.* genome graph), but maintained all 18

448    genomes in their native state. This was key to our ability to precisely compare all

449    genomes one-by-one. Lastly, we interrogated a high-sequence coverage subset of the

450    3K-RGP data to estimate and validate the population genetics of each inversion. As

451    sequencing costs continue to plummet, the ease at which ultra-high-quality genomes

452    can be generated, and with computational power exceeding current limits[43-45], we

453    predicted that there will no longer be a need to computationally generate pan-

454    genomes to perform similar analyses across much larger genomes, such as wheat

455    (genome size = 15 Gb)[46].

456        The Asian rice pan-genome inversion index is the first step on our quest to

457    precisely discover all standing natural variation that exists in Asian rice and

458    eventually the genus *Oryza* as a whole. The next step will be the generation of a

459    digital genebank for Asian rice whereby resequencing data from > 100,000

460    accessions will be mapped to our pan-genome - *i.e.* all 16 genomes. Preliminary data

461    (unpublished) shows that we can now easily call SNPs with resequencing data from >

462    3,000 individuals in 5 days per genome or less using supercomputer workflows

463    optimized for GATK4[47] software. Such call rates will undoubtedly increase over the

464    next year with a targeted rice digital genebank release date of January 1st, 2025.

465

466    **Online Methods**

467    All the methods are available in the supplementary information.

481   **AUTHOR CONTRIBUTIONS**

482   A.Z., D.W., K.M., J.Z., and R.A.W. designed and conceived the research. K.M. and

483   IRRI provided seed and/or tissue for all *Oryza* accessions. D.K., N.M., D.Ch., and

484   M.L. performed DNA extractions and genome sequencing. Y.Z., Z.Y., and J.Z.

485   performed sequence assembly, GPM edit and validation of 18 genome sequences.

486   Y.Z., N.M., D.K. and V.L. carried out the optical map sequence and analysis. D.Ch.,

487   Y.Z., J.S. and K.M. performed population genetic analysis. K.C., Z.L., and D.W.

488   performed the genome annotation and validation. Y.Z., Z.Y., J.Z. and A.Z. performed

489   the gene expression analysis and TE annotation. Y.Z., A.Z., J.S., A.A., S.M., Z.Y.,

490   and J.Z. carried out the inversion identification and population level validation. L.R.,

491   N.K., M.T., M.T., C.D., and K.C. managed the computing platforms. Y.Z., Z.Y.,

492   D.Co., K.C., N.A., A.F., A.Z., K.M., J.Z., and R.A.W wrote and edited the paper. All

493   authors read and approved the final manuscript.

494   **Competing Interests statement**

495   The authors declare that there is no conflict of interest regarding the publication of

496   this article.

497 **Figs legends**

498 **Fig. 1** Inversion landscape across 17 PSRefSeqs all relative to the IRGSP RefSeq shows

499 species-specific, group-specific and genome-specific inversions, *i.e.,* species-specific

500 inversions (*O. punctata* and *O. rufipogon*) are shown by black rectangles, group-specific

501 inversions are shown by red rectangles, and genome-specific inversions are shown by yellow

502 rectangles, respectively. On the left, accessions are phylogenetically ordered[24], on the bottom,

503 the tree are the clustering of inversions, and on the right, the numbers and lengths of the

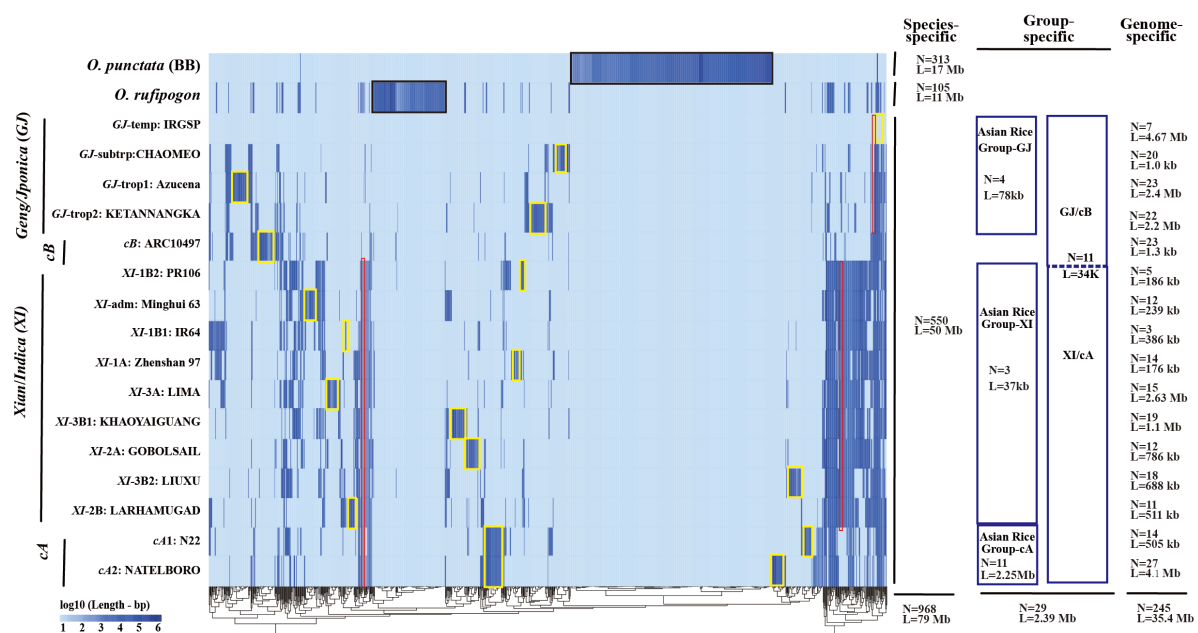504 specific inversions are presented.

505



506

507 **Fig. 2** Bionano validation of inversions larger than 1Mb. a. INV010130, b. INV010560 and c.

508 INV060390. In each panel, the top line has the optical map used as a reference, the bottom

509 line has the genome assembly of the variety with the inversion. Gray lines connect restriction

510 sites that are aligned (blue regions), while yellow segments are unaligned regions. Black

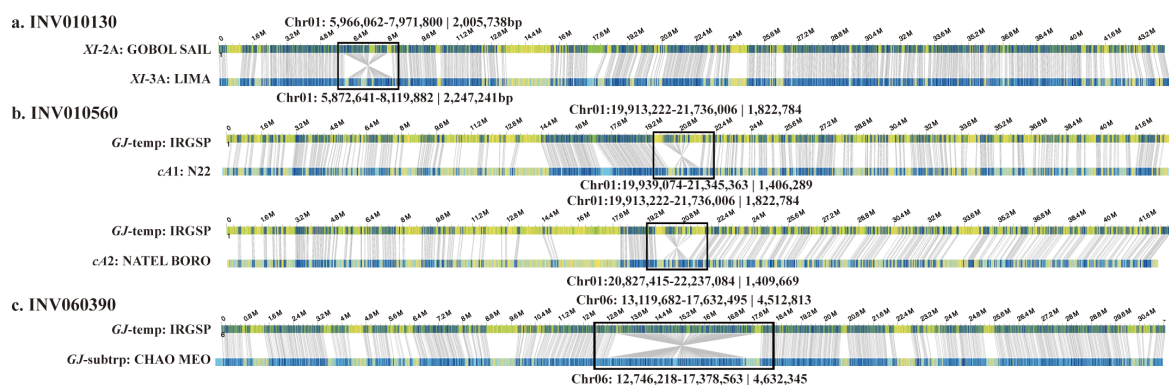511 boxes highlight the position of the inversion.

512

513   **Fig. 3** Transposable elements (TEs) are associated with inversions.

514   a. The amount (y-axes) of different TE families (x-axis) show that three TE families (*i.e.*

515   LTR-RT Ty1-*copia*, Ty3-*gypsy* and DNA-TE MULE) were observed in higher frequencies at

516   the breakpoints of the pan-genome inversion index than the resampled control tests. Box-

517   plots and bar-plots show the frequencies of TEs observed at the breakpoints of 10 random

518   resampling regions and the pan-genome inversion index, respectively.

519   b. Enrichment/depletion of 17 TEs present at the inversion breakpoints with more than 10

520   copies.

521   c. Details of Ty3-*gypsy* Os0025 presence at inversion breakpoints, with support from CCS
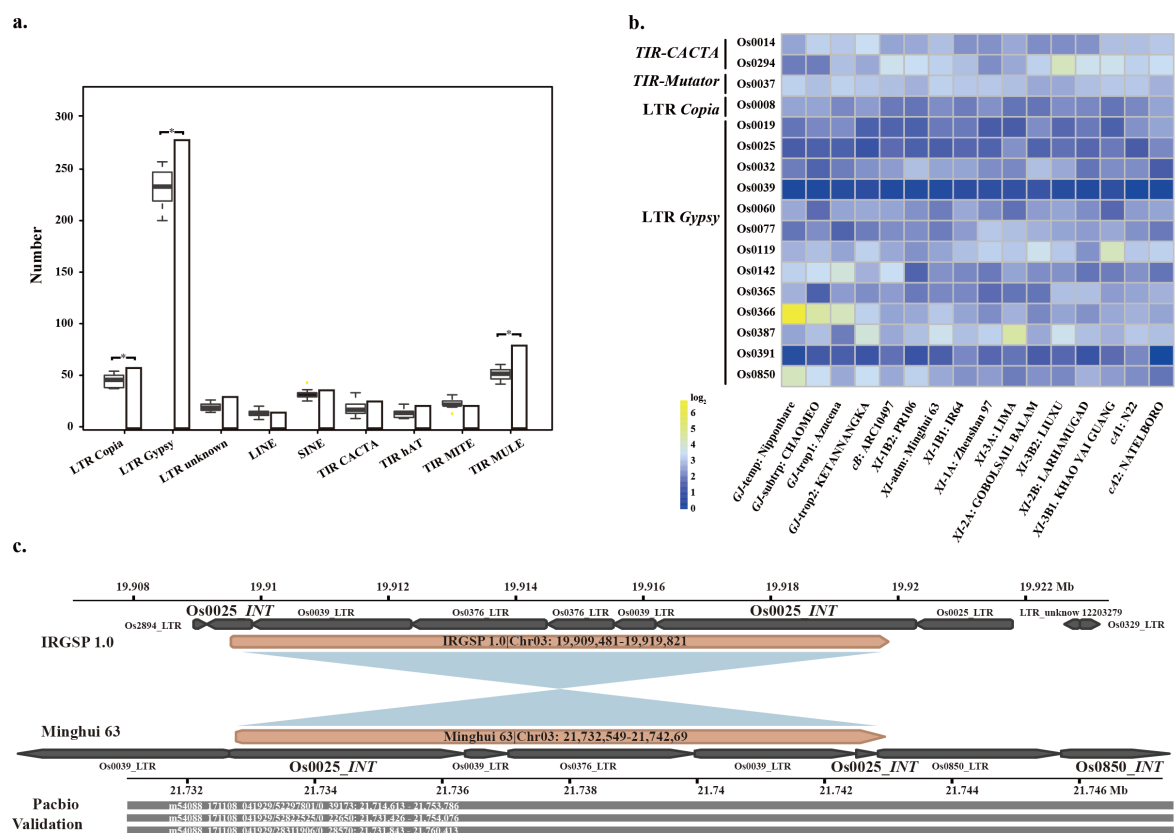
522   PacBio long reads.



523

524    **Fig. 4** Transcript abundance of genes located at inversion breakpoints.

525    a. Two copies of the *OsNAS* gene lie at the end of an inversion in the MH63RS2 (*XI*-adm)

526    genome. This inversion disrupted the 5' UTR regions.

527    b. *OsNAS* gene transcript abundance in root tissue.

528    c. The coding sequence (CDS) of a Fbox gene was disrupted by an inversion in the

529    MH63RS2 (*XI*-adm) genome.

530    d. Fbox transcript abundance was suppressed in all tissues tested in the MH63RS2 (*XI*-adm)
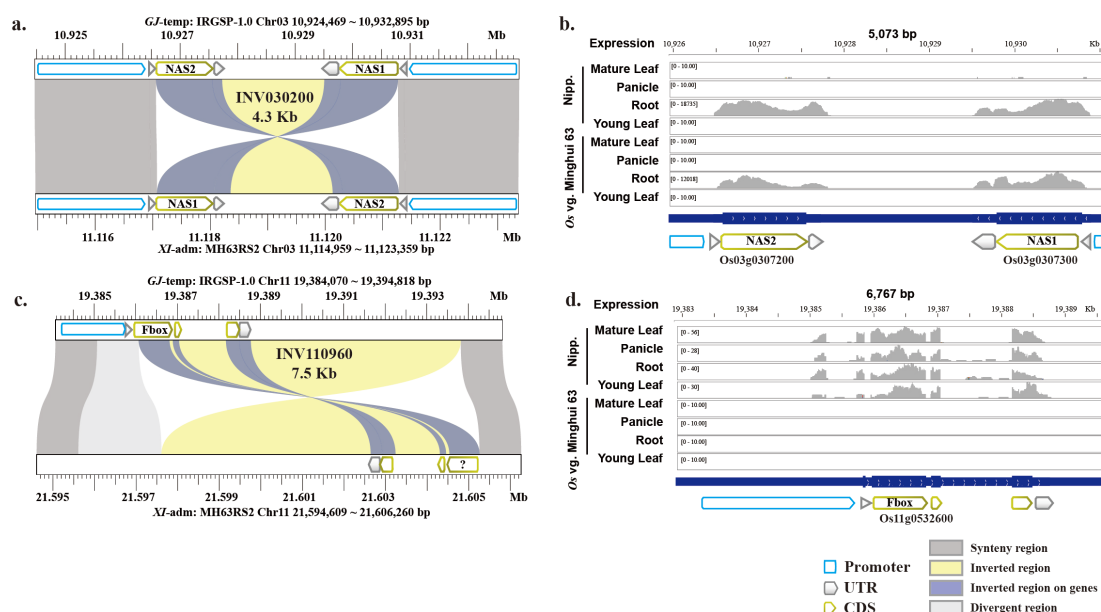
531    genome.



532

533 **Fig. 5** Population level SNP variation across large inversions. A schematic diagram of LD

534 block disruption arising from the presence of an inversion, as shown in A and B.

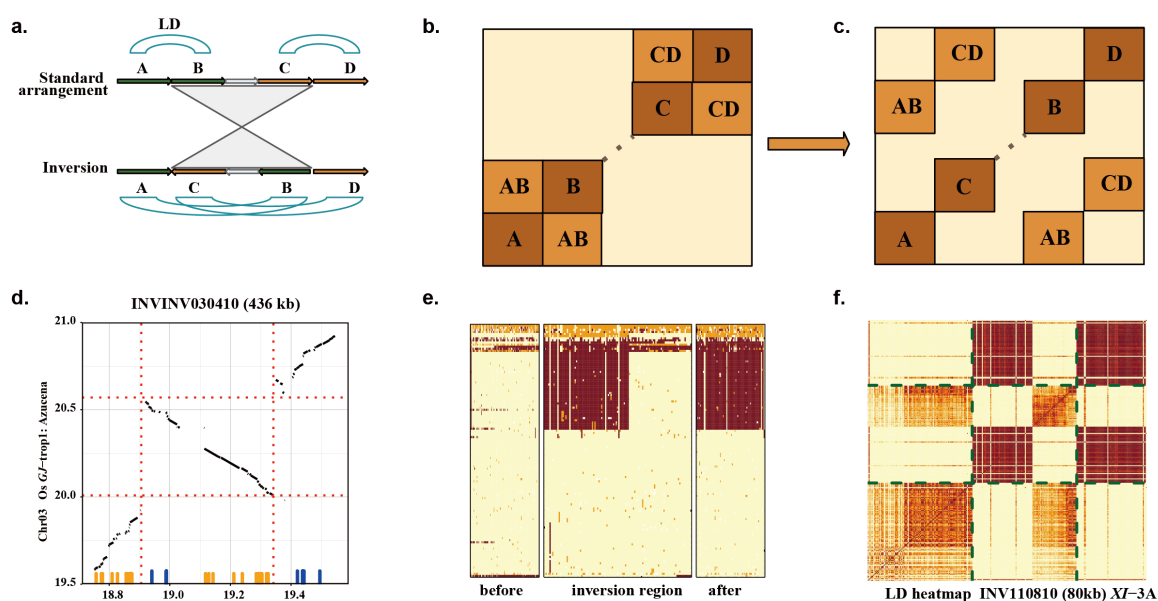535 a. Cartoon view of an inversion with breakpoints disrupting two LD blocks.

536 b. Expected features of the corresponding LD heat map.

537 c. Example of SNP blocks in high LD that are disrupted by an inversion.

538 d. The panel shows alignments, with the inversion marked by dotted lines. Small vertical

539 lines above the horizontal axis mark the location of SNPs constituting a disrupted LD block.

540 Orange and blue colors delineate two LD blocks that are contiguous in the of *GJ*-trop1

541 population, but appear as split when aligned to the IRGSP RefSeq (*GJ*-temp). Disruption of

542 Azucena (*GJ*-trop1) haplotype blocks along the IRGSP RefSeq in the region of INV030410,

543 as shown in e and f.

544 e. Genotype heat map of the *GJ*-trop1 subpopulation (samples in rows, SNPs in columns;

545 light yellow: reference call, orange: heterozygous, brown: homozygous variant).

546 f. LD heat map of the same subpopulation. Dotted lines show the inversion region. Darker

547 colors show larger $r^2$. Note that the scaling of X-axis in the genotype heat map is not uniform,
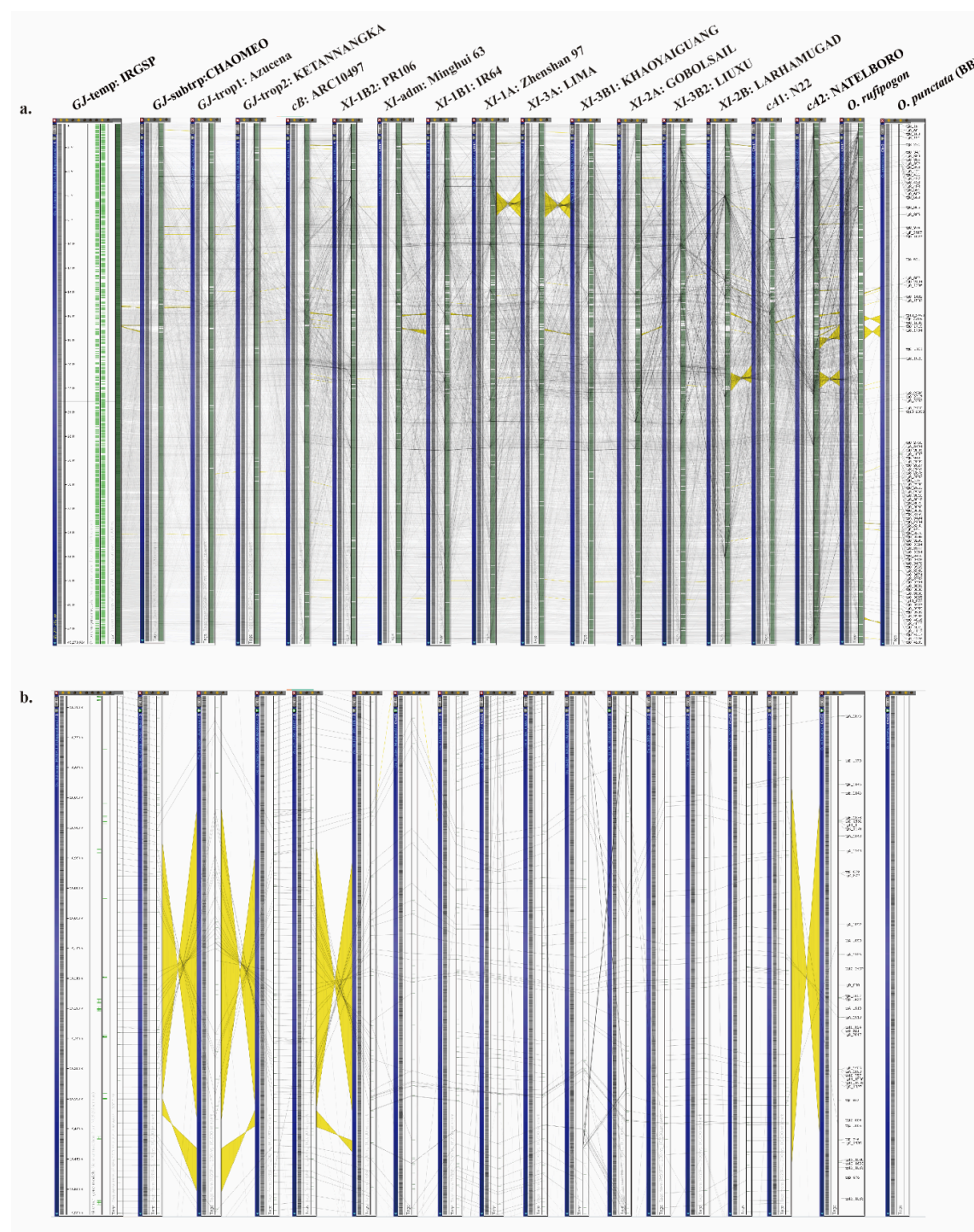
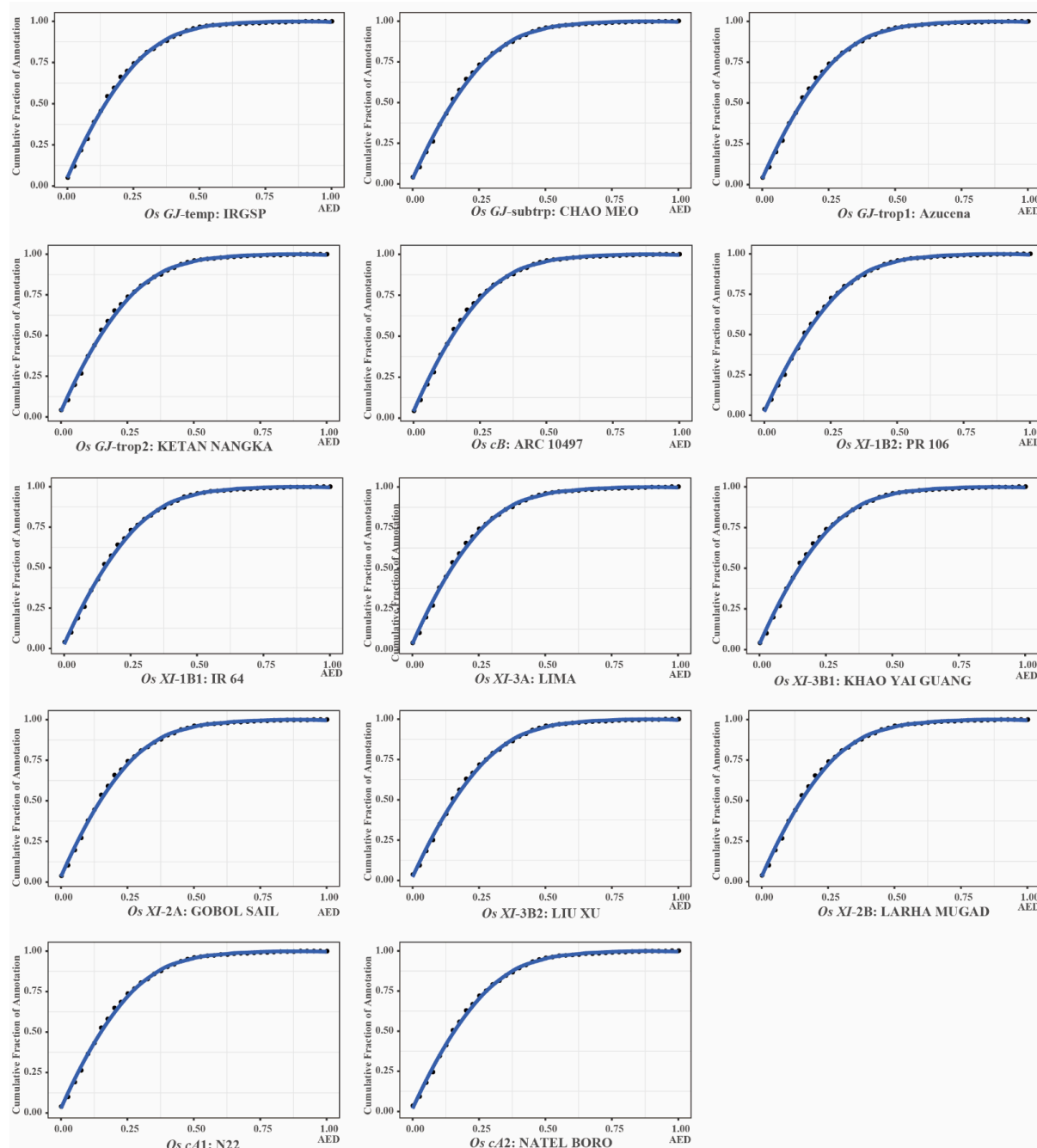548 allotting half of X axis space to the inverted region.



549

550 **Fig. 6** An integrated study of ClusterINV92.

551 a. Genome comparisons identified that inversion ClusterINV92 arose in *XI* and *cA* genomes

552 and are nearly fixed in *XI* and *cA* subpopulations. Note: [a] means 192 deep re-sequenced (> 20

553 ×) samples from the 3K-RGP were used for validation. b. LD block disruption arising from

554 the presence of ClusterINV92in populations with INV genotypes compared to standard

555 (STD) genotypes. c. Population genetic variation of nucleotide diversity, theta, and Tajima's

556 *D* of inversion cluster 92 genotypes compared to standard (STD) genotypes. The gray vertical

557 lines delimit the inversion coordinates on the IRGSP reference. d. Phenotyping test shows a

558 significant difference (linear model test, $p < 0.01$) in flowering time between ClusterINV92

559 genotypes and standard (STD) genotypes (left), and SNP variation in the LBD gene (right).

560 e. A total of 8 genes, including 4 reported genes, were observed in ClusterINV92. f. A single

561 SNP (G to T) caused a non-conservative amino acid change from the hydrophobic glycine

562 (Gly) to hydrophilic cysteine (Cys) within the LBD gene. g. Expression of the LBD gene was

563 up-regulated in panicle tissue, and suppressed in young leaf tissue of INV genotypes *XI*-adm:

564 MH63 and *XI*-1A: ZS97, compared with the standard IRGSP RefSeq genotype.



565

566 **Extended Data Figs**

567 **Extended Data Fig. 1** The 18-genome dataset (18 PSRefSeqs) was input into Persephone

568 (https://web.persephonesoft.com/) and made publicly available.

569 a. A panel shows overall alignments of the 18 maps (genomes) using chromosome 1 as an

570 example. The gray lines show the alignments of sequence tags and the yellow ribbons show

571 the inversions.

572 b. A panel shows an 800 kb region that includes INVcluster92 with yellow ribbons.



573

574    **Extended Data Fig. 2** Cumulative AED distributions of 13 genomes and their annotations
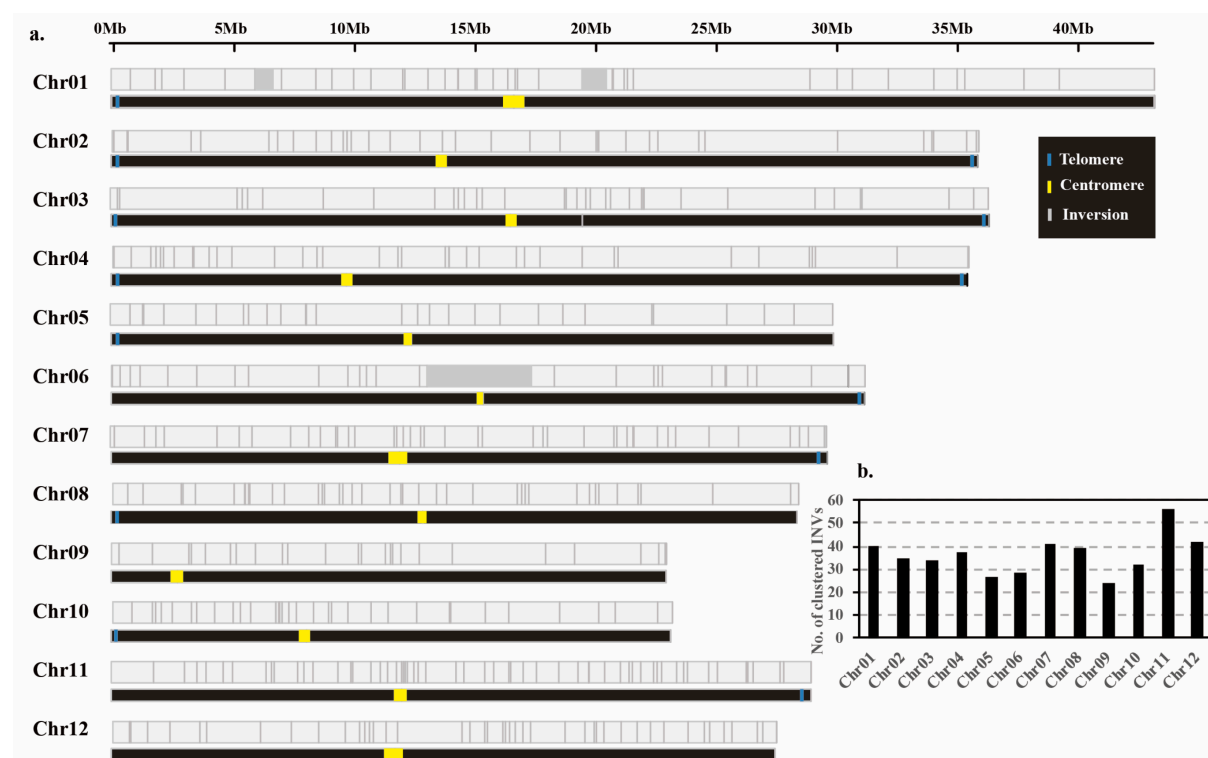
575    were plotted.



576

577 **Extended Data Fig. 3** Chromosome distribution and amount of inversions.

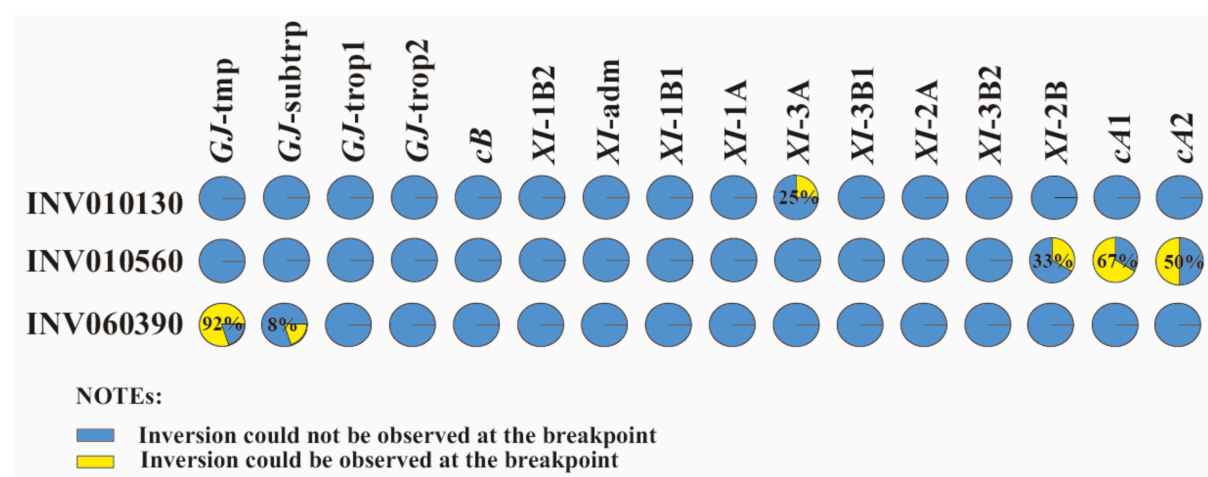578 a. Genome-wide distribution of inversions based on the *GJ*-temp: IRGSP-1.0 genome.

579 b. Number of inversions on 12 chromosomes based on the *GJ*-temp: IRGSP-1.0 genome.

580

581

582 **Extended Data Fig. 4** Validation of the 3 large (> 1 Mb) inversions in Asian rice,

583 (INV010130, INV0101306 and INV060390) using NGS read data having at least 20×

584 coverage from 192 3K-RGP samples.

585

586

587 **Extended Data Fig. 5** Recombination rate variation and inversion distribution in two *O.*

588 *sativa* PSRefSeqs (*XI*-adm: MH63RS2 and *XI*-1A: ZS97RS2).

589 a. Recombination rate of a RIL mapping population (Minghui 63 × Zhenshan 97). Dot points

590 indicate the recombination rate and gray boxes indicate inversions.

591 b. Comparison of recombination rate across inversions vs. genome-wide.



592

593 **Extended Data Fig. 6** Two SNPs and a 12 bp InDel variant were observed within WD40-76

594 gene.



595

596   **Tables**

597   **Table 1**. Assembly and annotation statistics of 18 *Oryza* reference genomes (18-genome data

598   package).

599   **Table 2**. Summary of inversions identified across the 18-genome dataset.

600

601   **Supplemental Tables**

602   **Supplementary Table 1**. Sequencing, data statistics of genome features, and BUSCO

603   evaluation of *de novo* assemblies for 2 new wild *Oryza* genomes, *i.e. O. rufipogon* [AA] and

604   *O. punctata* [BB].

605   **Supplementary Table 2**. Genome annotation statistics of 13 genomes.

606   **Supplementary Table 3**. BUSCO assessments of genome annotations using both (a)

607   transcriptome, and (b) protein model evidence.

608   **Supplementary Table 4**. Amount of PacBio Iso-Seq and RNA-Seq transcripts used for

609   genome annotation.

610   **Supplementary Table 5**. Pan-genome inversions across the 18-genome data package, by

611   comparing 15 *O. sativa* accessions, and 2 close relative genomes to the IRGSP-1.0. RefSeq.

612   **Supplementary Table 6**. Kolmogorov-Smirnov (KS) tests for genome-wide inversion

613   distribution.

614   **Supplementary Table 7**. Asian rice subpopulation validation for specific inversions.

615   **Supplementary Table 8**. A summary of subpopulation specific and group specific

616   inversions.

617   **Supplementary Table 9**. Details of 5 large inversions (> 1 Mb).

618   **Supplementary Table 10**. Summary of TE content of fine-scaled inverted regions based on

619   16 Asian rice genomes.

620   **Supplementary Table 11**. TE annotation of 16 *Oryza sativa* genomes.

621   **Supplementary Table 12**. 17 TEs from 4 superfamilies were observed with a higher amount

622   (> 10 in this study) at inversion breakpoints of 16 Asian genomes.

623   **Supplementary Table 13**. Gene content analysis of inversions across 16 Asian rice genomes.

624   **Supplementary Table 14**. Comparison of transcript abundance levels for genes that were

625   located within inversions, or at inversion breakpoints.

626   **Supplementary Table 15**. Seventy-eight inversions identified between *XI*-adm: MH63RS2

627   and *XI*-1A: ZS97RS2 genomes, the parents of a RIL-10 population. Note: To identify

628   recombination rates, we only focused on inversions > 1 kb.

629   **Supplementary Table 16**. Cluster for inversions larger than 100 kb.

630 **Supplementary Table 17**. Investigation of 10 phenotypes. FT: flowering time from sowing.

631 GWE: 100 grain weight (g). PL: panicle length (cm). PS: panicle shattering. SF: spikelet

632 fertility. LS: leaf senescence. LW: leaf width. LL: leaf length. SIEC12: Salt injury at EC12.

633 SIEC18: Salt injury at EC18.

634 **Supplementary Table 18**. Transcript abundance (FPKM value) of 8 genes located within

635 ClusterINV92 (INV030400/INV030410/INV030420).

636 **Supplementary Table 19**. Validation of inversion ClusterINV92 and variants (3 SNPs and

637 one 12 bp InDel) among different subpopulations by using 192 deep re-sequenced samples (>

638 20×).

639 **Supplementary Table 20**. Validation of 2,042 predicted inversions, of the same accessions

640 used to create the Asian rice pan-genome, collected from the 3K-RGP data set.

641 **Supplementary Table 21**. A list of 6 inversions that were identified from the 3K-RGP data

642 set, but were missed in the Asian rice pan-genome inversion index. INV3KSNPSEEK71124

643 and INV3KSNPSEEK71127, and INV3KSNPSEEK117247 and INV3KSNPSEEK117248

644 were found to be overlapping, and thus, were considered as a single inversion, resulting in 4

645 inversions were undetected in the pan-genome inversion index.

# References

1. Hossain, M. & Fischer, K. Rice research for food security and sustainable agricultural development in Asia: achievements and future challenges. *GeoJournal* **35**, 286-298 (1995).

2. Wing, R.A., Purugganan, M.D. & Zhang, Q. The rice genome revolution: from an ancient grain to Green Super Rice. *Nat Rev Genet* **19**, 505-517 (2018).

3. Vollset, S.E. *et al.* Fertility, mortality, migration, and population scenarios for 195 countries and territories from 2017 to 2100: a forecasting analysis for the Global Burden of Disease Study. *The Lancet* **396**, 1285-1306 (2020).

4. Kirkpatrick, M. How and why chromosome inversions evolve. *PLoS Biol* **8**, e1000501 (2010).

5. Wellenreuther, M. & Bernatchez, L. Eco-evolutionary genomics of chromosomal inversions. *Trends in ecology evolution* **33**, 427-440 (2018).

6. Hoffmann, A.A. & Rieseberg, L.H. Revisiting the impact of inversions in evolution: from population genetic markers to drivers of adaptive shifts and speciation? *Annual review of ecology, evolution, systematics* **39**, 21-42 (2008).

7. Sturtevant, A.H. A case of rearrangement of genes in Drosophila. *Proc. Natl. Acad. Sci. USA* **7**, 235 (1921).

8. Dobzhansky, T. & Sturtevant, A.H. Inversions in the chromosomes of Drosophila pseudoobscura. *Genetics* **23**, 28 (1938).

9. Barb, J.G. *et al.* Chromosomal evolution and patterns of introgression in Helianthus. *Genetics* **197**, 969-979 (2014).

10. Volkert, F.C. & Broach, J.R. Site-specific recombination promotes plasmid amplification in yeast. *Cell* **46**, 541-550 (1986).

11. Johnson, R.C. Bacterial site-specific DNA inversion systems. in *Mobile DNA II* 230-271 (American Society of Microbiology, 2002).

12. Hammer, M.F., Schimenti, J. & Silver, L.M. Evolution of mouse chromosome 17 and the origin of inversions associated with t haplotypes. *Proc. Natl. Acad. Sci. USA* **86**, 3261-3265 (1989).

13. Flores, M. *et al.* Recurrent DNA inversion rearrangements in the human genome. *Proc. Natl. Acad. Sci. USA* **104**, 6099-6106 (2007).

14. Hellen, E.H. Inversions and evolution of the human genome. *eLS*, 1-6 (2015).

15. Coughlan, J.M. & Willis, J.H. Dissecting the role of a large chromosomal inversion in life history divergence throughout the Mimulus guttatus species complex. *Molecular ecology* **28**, 1343-1357 (2019).

16. Hey, J. Speciation and inversions: chimps and humans. *Bioessays* **25**, 825-8 (2003).

17. Levy-Sakin, M. *et al.* Genome maps across 26 human populations reveal population-specific patterns of structural variation. *Nat Commun* **10**, 1-14 (2019).

18. Wang, W. *et al.* Genomic variation in 3,010 diverse accessions of Asian cultivated rice. *Nature* **557**, 43-49 (2018).

19. RGP, K. The 3,000 rice genomes project. *Gigascience* **3**, 7 (2014).

20. Fuentes, R.R. *et al.* Structural variants in 3000 rice genomes. *Genome Res* **29**, 870-880 (2019).

21. Kou, Y. *et al.* Evolutionary genomics of structural variation in Asian rice (Oryza sativa) domestication. *Molecular biology evolution* **37**, 3507-3524 (2020).

691 22. Stein, J.C. *et al.* Genomes of 13 domesticated and wild rice relatives highlight genetic
692 conservation, turnover and innovation across the genus Oryza. *Nat Genet* **50**, 285-
693 296 (2018).

694 23. Zhang, J. *et al.* Extensive sequence divergence between the reference genomes of
695 two elite indica rice varieties Zhenshan 97 and Minghui 63. *Proc. Natl. Acad. Sci. USA*
696 **113**, E5163-71 (2016).

697 24. Zhou, Y. *et al.* A platinum standard pan-genome resource that represents the
698 population structure of Asian rice. *Sci Data* **7**, 113 (2020).

699 25. Kawahara, Y. *et al.* Improvement of the Oryza sativa Nipponbare reference genome
700 using next generation sequence and optical map data. *Rice (N Y)* **6**, 4 (2013).

701 26. Huang, J. *et al.* Identifying a large number of high-yield genes in rice by pedigree
702 analysis, whole-genome sequencing, and CRISPR-Cas9 gene knockout. *Proc Natl Acad*
703 *Sci U S A* **115**, E7559-E7567 (2018).

704 27. Xie, F. & Zhang, J. Shanyou 63: an elite mega rice hybrid in China. *Rice (N Y)* **11**, 17
705 (2018).

706 28. Kou, Y. *et al.* Evolutionary genomics of structural variation in Asian rice (Oryza sativa)
707 and its wild progenitor (O. rufipogon). *BioRxiv* (2019).

708 29. Kim, H. *et al.* Comparative physical mapping between Oryza sativa (AA genome type)
709 and O. punctata (BB genome type). *Genetics* **176**, 379-90 (2007).

710 30. Qin, P. *et al.* Pan-genome analysis of 33 genetically diverse rice accessions reveals
711 hidden genomic variations. *Cell* **184**, 3542-3558 (2021).

712 31. Longbiao, G. *et al.* Genetic Analysis and Utilization of the Important Agronomic Traits
713 on Zhenshan 97* Minghui 63 Recombinant Inbred Lines (RIL) in Rice (Oryza sativa L.).
714 *Zuo wu xue bao* **28**, 644-649 (2002).

715 32. Huang, K. & Rieseberg, L.H. Frequency, Origins, and Evolutionary Role of
716 Chromosomal Inversions in Plants. *Frontiers in Plant Science* **11**, 296 (2020).

717 33. Li, C. *et al.* OsLBD37 and OsLBD38, two class II type LBD proteins, are involved in the
718 regulation of heading date by controlling the expression of Ehd1 in rice. *Biochemical*
719 *Biophysical Research Communications* **486**, 720-725 (2017).

720 34. Žárský, V., Sekereš, J., Kubátová, Z., Pečenková, T. & Cvrčková, F. Three subfamilies of
721 exocyst EXO70 family subunits in land plants: Early divergence and ongoing
722 functional specialization. *Journal of experimental botany* **71**, 49-62 (2020).

723 35. Ouyang, Y., Huang, X., Lu, Z., Yao, J. Genomic survey, expression profile and co-
724 expression network analysis of OsWD40 family in rice. *BMC Genomics* **13**, 100
725 (2012).

726 36. Itoh, J.I., Kitano, H., Matsuoka, M. & Nagato, Y. Shoot organization genes regulate
727 shoot apical meristem organization and the pattern of leaf primordium initiation in
728 rice. *The Plant Cell* **12**, 2161-74 (2000).

729 37. McClintock, B. *Cytological observations of deficiencies involving known genes,*
730 *translocations and an inversion in Zea mays*, (University of Missouri, College of
731 Agriculture, Agricultural Experiment Station, 1931).

732 38. Ahn, S. & Tanksley, S. Comparative linkage maps of the rice and maize genomes.
733 *Proc. Natl. Acad. Sci. USA* **90**, 7980-7984 (1993).

734 39. Burke, J.M. *et al.* Comparative mapping and rapid karyotypic evolution in the genus
735 helianthus. *Genetics* **167**, 449-57 (2004).

736 40. Husband, B. Chromosomal variation in plant evolution. (JSTOR, 2004).

737    41.    Levin, D.A. & Donald, A. *The role of chromosomal change in plant evolution*, (Oxford
738          University Press, USA, 2002).
739    42.    Crow, T. *et al.* Gene regulatory effects of a large chromosomal inversion in highland
740          maize. *PLoS Genet* **16**, e1009213 (2020).
741    43.    Gage, J.L., Monier, B., Giri, A. & Buckler, E.S. Ten years of the maize nested
742          association mapping population: impact, limitations, and future directions. *The Plant*
743          *Cell* **32**, 2083-2093 (2020).
744    44.    The Computational Pan-Genomics Consortium. Computational pan-genomics: status,
745          promises and challenges. *Brief Bioinform* **19**, 118-135 (2018).
746    45.    Vernikos, G.S. A Review of Pangenome Tools and Recent Studies. in *The Pangenome:*
747          *Diversity, Dynamics and Evolution of Genomes* (eds. Tettelin, H. & Medini, D.) 89-112
748          (Cham (CH), 2020).
749    46.    Athiyannan, N. *et al.* Long-read genome sequencing of bread wheat facilitates
750          disease resistance gene cloning. *Nat Genet* **54**, 227-231 (2022).
751    47.    Bathke, J. & Lühken, G. OVarFlow: a resource optimized GATK 4 based Open source
752          Variant calling workFlow. *BMC bioinformatics* **22**, 1-18 (2021).
753