

In the format provided by the authors and unedited.

Genomic variation in 3,010 diverse accessions of Asian cultivated rice

Wensheng Wang^{1,17}, Ramil Mauleon^{2,17}, Zhiqiang Hu^{1,3,17}, Dmytro Chebotarov^{2,17}, Shuaishuai Tai^{4,17}, Zhichao Wu^{1,5,17}, Min Li^{6,7,17}, Tianqing Zheng^{1,17}, Roven Rommel Fuentes^{2,17}, Fan Zhang^{1,17}, Locedie Mansueto^{2,17}, Dario Copetti^{2,8,17}, Millicent Sanciangco², Kevin Christian Palis², Jianlong Xu^{1,5,6}, Chen Sun³, Binying Fu^{1,6}, Hongliang Zhang⁹, Yongming Gao^{1,6}, Xiuqin Zhao¹, Fei Shen⁹, Xiao Cui³, Hong Yu¹⁰, Zichao Li⁹, Miaolin Chen³, Jeffrey Detras², Yongli Zhou^{1,6}, Xinyuan Zhang⁵, Yue Zhao³, Dave Kudrna⁸, Chunchao Wang¹, Rui Li³, Ben Jia³, Jinyuan Lu³, Xianchang He³, Zhaotong Dong³, Jiabao Xu⁴, Yanhong Li⁴, Miao Wang⁴, Jianxin Shi³, Jing Li³, Dabing Zhang³, Seunghee Lee⁸, Wushu Hu⁴, Alexander Poliakov¹¹, Inna Dubchak^{11,12}, Victor Jun Ulat², Frances Nikki Borja², John Robert Mendoza¹³, Jauhar Ali², Jing Li³, Qiang Gao⁴, Yongchao Niu⁴, Zhen Yue⁴, Ma. Elizabeth B. Naredo², Jayson Talag⁸, Xueqiang Wang⁹, Jinjie Li⁹, Xiaodong Fang⁴, Ye Yin⁴, Jean-Christophe Glaszmann^{14,15}, Jianwei Zhang⁸, Jiayang Li^{1,10}, Ruairaidh Sackville Hamilton², Rod A. Wing^{2,8*}, Jue Ruan^{5*}, Gengyun Zhang^{4,6*}, Chaochun Wei^{3,16*}, Nikolai Alexandrov^{2*}, Kenneth L. McNally^{2*}, Zhikang Li^{1,6*} & Hei Leung²

¹Institute of Crop Sciences, Chinese Academy of Agricultural Sciences, Beijing, China. ²International Rice Research Institute, Manila, Philippines. ³School of Life Sciences and Biotechnology, Shanghai Jiao Tong University, Shanghai, China. ⁴BGI Genomics, BGI-Shenzhen, Shenzhen, China. ⁵Agricultural Genomics Institute, Chinese Academy of Agricultural Sciences, Shenzhen, China. ⁶Shenzhen Institute for Innovative Breeding, Chinese Academy of Agricultural Sciences, Shenzhen, China. ⁷Anhui Agricultural University, Hefei, China. ⁸Arizona Genomics Institute, School of Plant Sciences, University of Arizona, Tucson, AZ, USA. ⁹China Agricultural University, Beijing, China. ¹⁰Institute of Genetics and Developmental Biology, Chinese Academy of Sciences, Beijing, China. ¹¹DOE Joint Genome Institute, Walnut Creek, CA, USA. ¹²Lawrence Berkeley National Laboratory, Berkeley, CA, USA. ¹³Advanced Science and Technology Institute, Department of Science and Technology, Quezon City, Philippines. ¹⁴UMR AGAP, CIRAD, Montpellier, France. ¹⁵UMR AGAP, Université de Montpellier, Montpellier, France. ¹⁶Shanghai Center for Bioinformatics Technology, Shanghai, China. ¹⁷These authors contributed equally: Wensheng Wang, Ramil Mauleon, Zhiqiang Hu, Dmytro Chebotarov, Shuaishuai Tai, Zhichao Wu, Min Li, Tianqing Zheng, Roven Rommel Fuentes, Fan Zhang, Locedie Mansueto, Dario Copetti. *e-mail: lizhikang@caas.cn; k.mcnally@irri.org; ccwei@sjtu.edu.cn; ruanjue@caas.cn; nalexandrov@inriag.com; zhanggengyun@genomics.cn; rwing@ag.arizona.edu

Genomic variation in 3,010 diverse accessions of Asian cultivated rice

Supplementary Notes

Wensheng Wang^{1,17}, Ramil Mauleon^{2,17}, Zhiqiang Hu^{1,3,17}, Dmytro Chebotarov^{2,17}, Shuaishuai Tai^{4,17}, Zhichao Wu^{1,5,17}, Min Li^{6,7,17}, Tianqing Zheng^{1,17}, Roven Rommel Fuentes^{2,17}, Fan Zhang^{1,17}, Locedie Mansueto^{2,17}, Dario Copetti^{2,8,17}, Millicent Sanciangco², Kevin Christian Palis², Jianlong Xu^{1,5,6}, Chen Sun³, Binying Fu^{1,6}, Hongliang Zhang⁹, Yongming Gao^{1,6}, Xiuqin Zhao¹, Fei Shen⁹, Xiao Cui³, Hong Yu¹⁰, Zichao Li⁹, Miaolin Chen³, Jeffrey Detras², Yongli Zhou^{1,6}, Xinyuan Zhang⁵, Yue Zhao³, Dave Kudrna⁸, Chunchao Wang¹, Rui Li³, Ben Jia³, Jinyuan Lu³, Xianchang He³, Zhaotong Dong³, Jiabao Xu⁴, Yanhong Li⁴, Miao Wang⁴, Jianxin Shi³, Jing Li³, Dabing Zhang³, Seunghye Lee⁸, Wushu Hu⁴, Alexander Poliakov¹¹, Inna Dubchak^{11,12}, Victor Jun Ulat², Frances Nikki Borja², John Robert Mendoza¹³, Jauhar Ali², Jing Li³, Qiang Gao⁴, Yongchao Niu⁴, Zhen Yue⁴, Ma. Elizabeth B. Naredo², Jayson Talag⁸, Xueqiang Wang⁹, Jinjie Li⁹, Xiaodong Fang⁴, Ye Yin⁴, Jean-Christophe Glaszmann^{14,15}, Jianwei Zhang⁸, Jiayang Li^{1,10}, Ruairaidh Sackville Hamilton², Rod A. Wing^{2,8}, Jue Ruan⁵, Gengyun Zhang^{4,6}, Chaochun Wei^{3,16}, Nickolai Alexandrov², Kenneth L. McNally², Zhikang Li^{1,6} & Hei Leung²

¹Institute of Crop Sciences, Chinese Academy of Agricultural Sciences, Beijing, China.

²International Rice Research Institute, Manila, Philippines

³School of Life Sciences and Biotechnology, Shanghai Jiao Tong University, Shanghai, China.

⁴BGI Genomics, BGI-Shenzhen, Shenzhen, China.

⁵Agricultural Genomics Institute, Chinese Academy of Agricultural Sciences, Shenzhen, China.

⁶Shenzhen Institute for Innovative Breeding, Chinese Academy of Agricultural Sciences, Shenzhen, China.

⁷Anhui Agricultural University, Hefei, China.

⁸Arizona Genomics Institute, School of Plant Sciences, University of Arizona, Tucson, AZ, USA.

⁹China Agricultural University, Beijing, China.

¹⁰Institute of Genetics and Developmental Biology, Chinese Academy of Sciences, Beijing, China.

¹¹DOE Joint Genome Institute, Walnut Creek, CA, USA.

¹²Lawrence Berkeley National Laboratory, Berkeley, CA, USA.

¹³Advanced Science and Technology Institute, Department of Science and Technology, Quezon City, Philippines.

¹⁴CIRAD, UMR Agap, TA A-108/3, Montpellier, France.

¹⁵AGAP, University of Montpellier, CIRAD, INRA, Montpellier SupAgro, Montpellier, France.

¹⁶Shanghai Center for Bioinformation Technology, Shanghai, China

¹⁷These authors contributed equally: Wensheng Wang, Ramil Mauleon, Zhiqiang Hu, Dmytro Chebotarov, Shuaishuai Tai, Zhichao Wu, Min Li, Tianqing Zheng, Roven Rommel Fuentes, Fan Zhang, Locedie Mansueto, Dario Copetti.

Correspondence and requests for materials should be addressed to Zhikang Li (lizhikang@caas.cn), Kenneth L. McNally (k.mcnally@irri.org), Chaochun Wei (ccwei@sjtu.edu.cn), Jue Ruan (ruanjue@caas.cn), Nickolai Alexandrov (nalexandrov@inariag.com), Gengyun Zhang (zhanggengyun@genomics.cn), and Rod A. Wing (rwing@ag.arizona.edu).

Contents

Overview of the sequencing quality of 3K rice data.....	2
Removal of 14 of the 3,024 accessions	3
<i>k</i> -mer analysis for genome characters.....	3
Mapping and genotype calling rate.....	3
Validation of SNP discovery pipeline	4
Validation of discovered SNPs in 3,010 accessions.....	5
Annotation of SNPs in the context of TEs	7
Inbreeding and heterozygous SNPs	8
SNP discovery and projection of undiscovered fractions	10
Utility of 3K RG for GWAS.....	12
Detection of genomic structural variations (SVs)	14
Correlation of presence/absence of SVs with plant heights.....	15
<i>De novo</i> assembly of 3,010 rice genomes	16
Evaluation of the quality of <i>de novo</i> assembly	16
Construction of the pan-genome sequences	17
Pseudogene detection.....	19
Length distribution of novel genes.....	19
Read mapping to the pan-genome.....	20
Evaluation of gene presence/absence detection	21
Estimating size of the rice pan-genome	23
The average gene/ gene family difference between two accessions	23
Phylogenetic analysis based on gene (or gene family) PAV	23
Inferring gene and gene family age with 446 wild rice genomes	24
References.....	25

1 **Overview of the sequencing quality of 3K rice data**

2 The methods of selecting the accessions for sequencing as well as the sequencing
3 methodology (paired-end sequencing separated by ~450 bp on Illumina HiSeq2000)
4 were described in detail in the 3,000 Rice Genomes data note¹. Updated metadata
5 information is available in **Supplementary Data 1 Table 1**.

6 Regarding data processing, paired-end reads were trimmed to 83 bp, generating
7 205,084,357,762 paired-end reads for 3,024 genomes. For 287 samples, two or more
8 DNA libraries were created; whereas, for the other 2,737 samples, a single library was
9 made. Usually, each sample library was sequenced on several flowcells. Every
10 flowcell lane/sample library combination resulted in a pair of separate fastq files. The
11 total number of fastq files was 51,060 (25,530 pairs). Most samples have 12 files (6
12 pairs sequenced independently). Most flowcell lanes contained 23 different samples
13 that were separated by index sequences.

14 The sequencing depths of the 3,024 genomes ranged widely from 4x to 50x, with
15 a mean of $14.3 \pm 6.3x$, a median of 13.2x, and adjusted mean sequencing depth of
16 $14.9 \pm 6.2x$. Of these, 2,461 (81.4%) rice lines have sequencing depths over 10x and
17 458 (15.1%) have sequencing depths over 20x.

18 Before we carried out further analyses, we examined the quality of our
19 sequencing data with the FastQC software (v0.11.2,
20 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>), with results available at
21 <http://oryzasnp.org/3kfastqc/>.

22 **Removal of 14 of the 3,024 accessions**

23 Fourteen accessions were removed from the subsequent analyses. Samples CX400,
24 CX401, CX402, IRIS_313-11415, and IRIS_313-10729 belong to African cultivated
25 rice (*Oryza glaberrima* L.). Samples IRIS_313-8502, IRIS_313-9233, IRIS_313-8444,
26 IRIS_313-10057, IRIS_313-9184, B014, and IRIS_313-9404 were removed due to
27 significant contamination, and sample B101 was removed due to a very small
28 estimated genome size. IRIS_313-8921 was removed because of its extremely low
29 sequencing depth (~0.3x).

30 ***k*-mer analysis for genome characters**

31 We used the 17 bp *k*-mer analysis with *k*-mer_count (in-house software, written by Jue
32 Ruan) to calculate the distribution of *k*-mer frequency for genome character
33 estimation (genome size, repeat ratio and GC content) with total reads of each
34 accession. The calculation is based only on the statistic of *k*-mers over the first trough
35 point of the 17 bp *k*-mer frequency distribution graph; the data before that were
36 considered to be sequencing errors. The genome size (GS) = total number of *k*-mers
37 /peak value of *k*-mers frequency, repeat ratio = (GS-number of *k*-mers with different
38 sequence + number of *k*-mers before depth of the first trough point)/GS.

39 **Mapping and genotype calling rate**

40 The alignment statistics table (**Supplementary Data 2 Table 1**) shows that the
41 highest mapping rate of reads for an accession occurs when the accession and

42 reference genome belong to the same variety group (grouping is determined by
43 population analysis described in subsequent sections), which also indicates that more
44 genotyping calls are detected when the accession and reference genome belong to the
45 same variety group. It is also notable that there is no appreciable increase in
46 genotyping sensitivity from sequencing coverage 30x and beyond (**Supplementary**
47 **Data 2 Figure 1**).

48 **Validation of SNP discovery pipeline**

49 To validate the GATK UnifiedGenotyper (GATK-UG)² SNP discovery pipeline, we
50 compared the SNPs discovered by this pipeline against a standard SNP discovery
51 method on two selected subset accessions. For the standard method, we used the
52 MUMmer³ pipeline to first discover SNPs between two published reference genomes,
53 Nipponbare RefSeq (Nipponbare IRGSP 1.0 genome) and IR 64
54 (os.ir64.cshl.draft.1.0). Reciprocal whole-genome alignment was carried out using
55 nucmer using default settings. SNPs from this alignment were extracted from each
56 reciprocal alignment using show-snps with the parameters “-C -l -r -T”. The common
57 SNPs from each reciprocal alignment (same positions, same allele states) were
58 extracted and 776,622 high-quality SNPs were discovered between the Nipponbare
59 and IR64 references (which we call the Nipponbare-IR 64 reference SNPs).

60 We then selected the two re-sequenced accessions, CX140 (Nipponbare) and
61 CX403 (IR 64a), from the 3K RG that have the same name but were of different
62 provenance as the reference/published genomes. SNPs were discovered for these

63 accessions for IR 64 ref vs. CX140, and Nipponbare ref vs. CX403, using the
64 GATK-UG pipeline with high-quality SNPs selected by strict filtering criteria (not
65 flagged as LowQual, no missing alleles, homozygous only). Of the 689,797
66 Nipponbare ref: CX403 and 699,429 IR 64 ref: CX140 SNPs that shared common
67 positions with the Nipponbare-IR 64 reference SNPs, 99.9% of the alternate allele
68 calls were concordant, indicating that the GATK-UG pipeline performed well in SNP
69 calling.

70 **Validation of discovered SNPs in 3,010 accessions**

71 The ~27 million biallelic SNPs from 3K RG were also compared with those
72 discovered by previous projects: (1) the 44k SNP project⁴, (2) the Rice50 SNPs
73 project⁵, (3) SNP-Seek release 1⁶, and (4) dbSNP rice release 147⁷
74 (<https://www.ncbi.nlm.nih.gov/projects/SNP/index.html>).

75 Using 44,100 SNPs from the 44k SNP project, most (36,775) of the original 44k
76 SNPs could be re-mapped to the Nipponbare RefSeq. Of the re-mapped 44k SNPs,
77 94.5% are in the same position as the 3K RG biallelic SNPs, and in this common set,
78 99.8% have the same alternate allele calls.

79 SNPs from the 50 resequenced rice genomes⁵ (6,496,456 high-quality Rice50
80 SNPs) were compared for concordance with the 3K RG SNPs. Since the Rice50 SNPs
81 were anchored to the IRGSP4 Nipponbare genome assembly, these were mapped to
82 the Nipponbare RefSeq, and the overwhelming majority (>99.9%) of the SNPs
83 (6,496,018) from Rice50 were anchored to this reference version. The intersection

84 with the 3K RG SNPs showed that ~56% of the Rice50 SNPs (3,669,353) mapped to
85 the same position as the 3K RG SNPs. However, since many of these Rice50 SNPs
86 had ambiguous or heterozygous alternate allele calls (1,559,842), these were excluded
87 in the concordance comparison, leaving 2,109,511 common SNP positions having
88 definite alternate calls. For these common SNP positions, ~2.09 million were
89 concordant in alternate allele calls (99.1%).

90 The full 3K RG Nipponbare SNP dataset (~29 million SNPs) was contrasted to
91 the ~20.3 million SNPs in the first release⁶ (2015) of SNP-Seek to determine the
92 effect of earlier software versions of BWA (ALN vs MEM)⁸ and GATK (2 vs 3)² on
93 SNP discovery. Although the majority of the SNPs in the new set are common to
94 those in the 2015 SNP-Seek release (~17.54 million, 86% of the release 1 set), ~14.5
95 million SNPs are unique to the new 3K RG SNPs (45.3% of the release 2 set) and
96 ~2.77 million SNPs (13.7% of the release 1 set) are unique to the 2014 SNP-Seek
97 release, highlighting the effects of using updated software on variant detection.
98 Therefore, it is worthwhile to update SNP discovery analyses as newer (and better)
99 SNP calling software becomes available.

100 We also compared the ~27 million biallelic SNPs from 3K RG with the NCBI
101 Reference cluster ID (rs#) SNPs from chromosomes 1 to 12 of build 147 of dbSNP -
102 rice (9,922,318 SNPs and 230,253 small indels of ≤50 bp). All rs# SNPs were
103 re-mapped (using the flanking sequence information) to the current build of the
104 Nipponbare RefSeq prior to comparison using mega-BLAST of NCBI, and
105 alignments where the entire flanking sequence aligns with one mismatch on the SNP

106 position were selected for the SNP comparison. A total of ~8.5 million rice SNPs in
107 dbSNP remapped to Nipponbare RefSeq, and 51.5% of the remapped dbSNPs had
108 common positions with the 3K RG SNPs. Of this common position SNP set,
109 4,199,667 remapped dbSNPs and 3K RG SNPs had the same alternate allele call (96%
110 concordance). In all, ~4.3 million remapped rice dbSNPs are unique from the 3K RG
111 set, and 25.4 million 3K RG SNPs are unique from rice dbSNPs. In comparison with
112 the newest submitted NCBI Assay ID (ss#) SNPs in dbSNP 147 (Q4 2014 submission
113 by McCouch et al.⁹, 700,000 SNPs with ss#), there were relatively higher numbers of
114 intersecting and concordant SNPs with the 3K RG SNPs (538,887 of 700k
115 intersecting the 3K RG SNPs at 99% concordance).

116 **Annotation of SNPs in the context of TEs**

117 The repeat content of the Nipponbare RefSeq is 48% of the total genome, with most
118 of the discovered SNPs occurring in LTR retro-elements and TIR DNA transposon
119 repeat types, the bulk of repeats in the rice genome. Repeats associated with telomeres
120 and centromeres, although comprising a very small portion of the genome, exhibit the
121 highest SNP densities observed (163 and 177 SNPs/kb); however, the number of
122 SNPs in these repeat types is too low to impact the number of SNPs discovered. The
123 average genome-wide SNP density (combined repetitive and non-repetitive regions) is
124 lower (72 SNPs/kb) than the average SNP density in repetitive regions only (103
125 SNPs/kb) and even lower in the non-repeat portion of the genome (44 SNPs/kb)
126 **(Supplementary Data 2).**

127 **Inbreeding and heterozygous SNPs**

128 Given that (a) rice is a selfing species, (b) accessions were subjected to single seed
129 descent before sequencing, and (c) there is a high degree of differentiation between
130 groups of rice varieties, one thus expects the number of heterozygous calls per SNP to
131 be much lower than postulated by Hardy-Weinberg equilibrium. Indeed, we observe
132 this systematic deviation from HW equilibrium in the 3K dataset (**Extended Data**
133 **Fig.1a**) as well as in major subpopulations.

134 To facilitate the discussion, let us introduce the following notation:

135 **Hobs**: observed SNP heterozygosity, the proportion of heterozygous calls in all
136 non-missing calls of a SNP.

137 **Hexp**: expected SNP heterozygosity given by Hardy-Weinberg equilibrium, i.e.
138 “ $Hexp=2pq$ ”, where p and q are the two allele frequencies.

139 The ratio $Hobs/Hexp$ is expected to be distributed around $1-F$, where F is Wright’s
140 inbreeding coefficient. We use the distribution of $Hobs/Hexp$ to estimate F and
141 remove outlier SNPs that might represent alignment errors.

142 We analyzed $Hobs/Hexp$ distribution in the whole 3K SNP dataset, as well as in
143 two subsets: XI (1,789 samples) and GJ (772 samples). In all three datasets, the
144 distribution of $Hobs/Hexp$ is bimodal, with one peak harboring most of the common
145 SNPs and that corresponds to a higher F (~ 0.95) and the other peak at around
146 “ $Hobs/Hexp=1$ ” caused mostly by rare SNPs with low numbers of homozygous

147 alternate calls and an excess of heterozygotes.

148 There is an excess number of points along the upper left boundary of the
149 scatterplot (**Extended Data Fig. 1a**) that corresponds to a maximal value of
150 “Hobs=2p”, where p is the minor allele frequency, indicating an excess of SNPs with
151 no or very few homozygous alternate calls. We hypothesize that these SNPs are due to
152 alignment errors caused by duplications that do not occur in a reference but are
153 present in certain genotypes.

154 We estimate the inbreeding coefficient F for XI and GJ datasets, as well as
155 (effective) inbreeding coefficient in the whole 3K dataset as the median value of
156 “1-Hobs/Hexp” for SNPs where “Hobs/Hexp <1” and the minor allele frequency
157 is >5%. Doing so, we ignore the peak near “Hobs/Hexp =1” as it is likely an
158 alignment artifact where reads map to multiple regions that are duplicates. The
159 resulting estimates are F=0.954 for the whole 3K, F=0.925 for XI, and F=0.969 for
160 GJ.

161 We use the estimates of F to introduce a quality cutoff for number of observed
162 heterozygotes:

$$163 \text{Hobs}_{\text{max}} = 10 (1-F) \text{Hexp}$$

164 That is, a SNP whose heterozygosity is >10x higher than the most likely value for
165 a given frequency and the dataset’s inbreeding rate will be deemed as having an
166 excessive number of heterozygotes and will be filtered out. The cutoff values for
167 different datasets are thus 0.4795082 for 3K, 0.7485704 for XI, and 0.3106786 for GJ
168 datasets.

169 We remove SNPs that violate this quality criterion in the 3K XI and GJ datasets
170 from the set of biallelic SNPs (**Extended Data Fig.1b**). We call the resulting set of
171 16,874,733 SNPs polymorphic in 3010 genomes as the **Base SNP set**.

172 Of note is that, although we expect that these might be erroneous calls, the fact
173 they occur preferentially in the third base of codons and exhibit properties for coarse
174 classification indicates that there is some biological significance meriting further
175 investigation.

176 **SNP discovery and projection of undiscovered fractions**

177 *Estimated proportion of IRRI Genebank SNPs discovered in 3,010 samples for a*
178 *given allele frequency.*

179 For a SNP that has a frequency f in the genebank, the probability that it has been
180 observed in 3,010 samples can be approximated by

$$181 \text{Prob (observed in 3K)} = 1 - (1-f)^{3010}$$

182 Note that this approximation is robust due to the large size of the genebank; using the
183 exact hypergeometric formula gives a probability estimate that differs from this by at
184 most 0.4% at any point. Note also that this is a conservative estimate, since we treat
185 each sample as haploid when in reality each is diploid.

186 Using this function, one can estimate that 3,010 samples capture more than 99.9%
187 of genebank SNPs of frequency greater than 0.25%, and virtually 100% of SNPs of
188 frequency $>1.1\%$ (**Extended Data Fig. 1c**).

189 However, this function alone does not allow estimation of the proportion of the

190 total number of SNPs (of all frequencies) captured by the 3K, since the distribution of
191 SNP frequencies is not known. This distribution depends on past demographic events
192 and selection. To estimate the total number of undiscovered SNPs, we adopted the
193 simulation approach outlined below.

194 **Estimation of the total number of Nipponbare RefSeq-based SNPs occurring in**
195 **genebank samples**

196 We computed SNP discovery rates based on 6,000 random permutations of samples
197 from the **Base** SNP set in the following way. For each permutation of sample order,
198 we computed the number of new SNPs added by each consecutive sample (i.e. SNPs
199 not seen in previous samples either as HOM or as HET). We then computed the mean
200 number of additional SNPs when an Nth sample is added across all permutations and
201 projected the mean for the range [3,010, 120,000] using regression between
202 $\log(\text{mean_new_SNP})$ and $\log(\text{sample})$.

203 The fitted model is

$$204 \quad \log(\text{mean_new_SNP}) = -0.75 * \log(\text{sample}) + 5.74$$

205 with $R^2 = 0.995$.

206 Then, we estimated the number of SNPs as a sum of the estimated mean added
207 SNPs at each value from 3,009 to 120,000.

208 As a result, we find that ~27M new SNPs are estimated to be discovered upon
209 genotyping the rest of the genebank, leading to ~44M total Nipponbare-based SNPs.

210 **Extended Data Fig. 1d** shows the number of SNPs added with each sample, on a
211 log-log scale, with a linear regression fit. One can also see from this graph, if the

212 trend continues as shown, it would be possible to extrapolate the number of SNPs that
213 may occur in the entire population of *O. sativa* including those not yet conserved in
214 genebank(s). This analysis is based on a subset of the **Base** SNP set consisting of
215 3,006 samples (in addition to the previously identified problematic samples, we
216 removed 4 samples that were outliers in terms of many private SNPs and made
217 extrapolation harder).

218 Similar analyses were done for XI and GJ separately, estimating 14M and 13M
219 new SNPs, respectively. However, one needs to model the overlap, and there is some
220 sharing of even the rarest non-singleton SNPs. Since this makes the modeling unduly
221 complicated, we report only the analysis based on the whole set.

222 **Utility of 3K RG for GWAS**

223 As a test case, a genome-wide association study (GWAS) was conducted using
224 historical phenotypic data for grain length (GRLT) and grain width (GRWD) and
225 newly acquired data for bacterial blight (BLB) resistance, with an LD pruned subset
226 of the 3K RG.

227 **1. Sample and SNP filtering**

228 We performed quality control measures by filtering low quality samples and markers
229 for use in GWAS. The samples were filtered to remove those for which missing data
230 or call rates (CR) were below 80% and with an over- and under-abundance of
231 heterozygous SNPs in the interquartile range (IQR), calculated as the difference
232 between the upper and lower quartiles ($IQR = Q3 - Q1$). We also excluded markers

233 with an over-abundance of heterozygous alleles (number of alleles >2), and those
234 markers with low call rates (CR <0.9) and minor allele frequencies (MAF <0.05).
235 Further, we pruned the dataset based on LD using the Composite Haplotype Method
236 (CHM)¹⁰ algorithm, with the following parameters: window size of 35 SNPs, window
237 increment of 15 SNPs, and r^2 threshold of 0.5. The LD pruned datasets had 2,012
238 samples and 223,743 markers for GRLT and GRWD and 381 samples and 148,999
239 markers for BLB.

240 **2. Phenotypic data**

241 We performed GWAS for source accessions of the sequenced genetic stocks with
242 historical phenotypic data for grain length (GRLT) and grain width (GRWD) and
243 newly-quantified bacterial leaf blight (BLB) scores. Trait data for GRLT and GRWD
244 were collected from unreplicated trials as genebank characterization data over many
245 seasons of trials; data for the source accession of the sequenced genetic stock were
246 used as a proxy for those of the derived genetic stock that underwent one or more
247 cycles of single seed descent from the source accession. BLB resistance was
248 measured as lesion length after infection with C5 Chinese strain of *Xanthomonas*
249 *oryzae* on the genetic stocks that were used for sequencing.

250 **3. GWAS methods**

251 We implemented an EMMAX (Efficient Mixed-Model Association eXpedited)¹¹
252 single-locus mixed linear model in SNP & Variation Suite v8.4.0 software
253 (<http://www.goldenhelix.com>) for GWAS. EMMAX allows correction for cryptic
254 relatedness and other fixed effects using a kinship matrix (as random effect) and

255 population stratification using the top four principal components (as fixed effect).
256 Both the kinship matrix and the principal components were generated from the LD
257 pruned datasets. We used the False Discovery Rate (FDR <0.01) multiple testing
258 correction to identify significant markers, and generated Manhattan and QQ plots
259 from the EMMAX output using the qqman package¹² in R.

260 **Detection of genomic structural variations (SVs)**

261 We tested BreakDancer¹³, DELLY¹⁴, and novoBreak¹⁵
262 (<https://sourceforge.net/projects/novobreak/?source=navbar>) for SV calling against
263 the Nipponbare RefSeq and with several SVs inserted into the Nipponbare genome
264 (**Supplementary Data 3 Tables 5 and 6**). novoBreak was found to have the lowest
265 false positive rate and comparatively good resolution and was therefore selected for
266 all subsequent SV detection. Briefly, novoBreak employed a *k*-mer (contiguous
267 nucleotide sequence of length *k*) targeted local assembly algorithm to detect structural
268 variation breakpoints in single base pair resolution. When applied to the 3K RGs for
269 discovering SVs, novoBreak first constructed a hash table of all the reads of a sample.
270 Next, any *k*-mers matching the Nipponbare RefSeq were removed. Then, novoBreak
271 employed a counting bloom filter to calculate the occurrence of all *k*-mers. At this
272 step, low-frequency *k*-mers reflecting sequencing errors were removed. Next,
273 high-frequency *k*-mers and their associated read pairs were clustered by a modified
274 union-find algorithm, ensuring that each cluster represented a single breakpoint. Then,
275 for each cluster, an assembler with a greedy algorithm was applied to assemble the

276 read pairs spanning the breakpoint into optimal and sub-optimal contigs. By
277 comparing the assembled contigs with the Nipponbare RefSeq, novoBreak inferred
278 the exact breakpoints of all types of SVs. Finally, novoBreak scored each SV based on
279 alignment and assembly evidence and a filter was applied to generate a
280 high-confidence SV list. In novoBreak, the detected translocations were referred to as
281 ‘inter-chromosomal breakpoints’. We detected SVs in the 3,010 accessions. In order
282 to minimize the probability of false positives, SVs detected in fewer than 6 accessions
283 (the number of such SVs = 207,879) or in more than 80% (the number of such SVs =
284 446) of the 3,010 accessions were removed. We analyzed the SVs detected in 453
285 well-sequenced accessions. Translocations and deletions account for 74.3% and 21.4%
286 of all SVs, respectively. Inversions and duplications account for only 1.7% and 2.4%,
287 respectively. The percentage of SVs detected may reflect both the real number of
288 different type SVs in the genomes and the false positives and negatives in SV
289 detection. Genes interrupted by SVs or inside SV regions were identified and we also
290 checked the co-localization between TEs and SVs (**Supplementary Data 3 Table 7**),
291 with 1 kb, 5 kb, and 10 kb windows with one breakpoint as the background. TE
292 annotation was from RGAP 7¹⁶.

293 **Correlation of presence/absence of SVs with plant heights**

294 This was calculated using ‘cor.test’ (method="spearman") in R. The SV with the
295 highest correlation was a ~385 bp deletion, located in the *sd1*¹⁷ gene
296 (LOC_Os01g66100) (rho= -0.40, *P*-value = 2.48E-10). The average height of

297 accessions with the deletion was 84.96 cm, while it was 126.50 cm for accessions
298 without the deletion.

299 ***De novo* assembly of 3,010 rice genomes**

300 In order to gain better assemblies, we compared the performance of several assembly
301 tools developed for NGS including SOAPdenovo version r240¹⁸, Velvet version
302 2.2.5¹⁹ and SPAdes version 3.0.0²⁰. Finally, a method with iterative use of
303 SOAPdenovo was selected for the 3K rice assembly that had better performance than
304 SPAdes and relatively good speed (~3.94 times running time of default SOAPdenovo).
305 The key idea of this variant method was to select the best *k*-mer for each sample.

306 QUAST version 2.3²¹ was used for evaluation of the assemblies, including 1)
307 comparison of assemblies among SOAPdenovo, Velvet and SPAdes and 2)
308 comparison among all rice accessions with our variant method described above with
309 parameter “-t 16 --min-contig 500 -o output --no-plots -R IRGSP-1.0.fa”. The
310 Nipponbare RefSeq (IRGSP-1.0) genome was used for all the evaluations. The
311 Nipponbare RefSeq was downloaded from the Rice Annotation Project (RAP)²².
312 Several important indices, including N50, assembled size, genome fraction (how
313 much of the Nipponbare RefSeq can be covered with the assembled contigs), and
314 unaligned contig size, were selected to evaluate the assembly performance.

315 **Evaluation of the quality of *de novo* assembly**

316 In order to evaluate the quality of *de novo* assembly, we first developed a pipeline to

317 correct/remove the possible misassembled contigs from read mapping. Reads (used
318 for the assembly) are mapped to the assembled contigs and we broke down the contigs
319 at positions with no evidence of connections supported by read alignments.
320 Remaining fragments shorter than 500 bp were removed.

321 Next, we assessed the assembly results based on accessions CX140 (compared
322 with the Nipponbare RefSeq), CX133 (compared with the Zhenshan 97 genome²³),
323 and CX145 (compared with the Minghui 63 genome²³). For CX140, we first applied
324 the correction procedure to both raw SOAPdenovo contigs and GapCloser contigs.
325 For CX145 and CX133, only raw SOAPdenovo contigs were evaluated. The
326 assemblies were assessed by QUILT with the Nipponbare RefSeq as a gold standard.
327 Results are shown in **Supplementary Data 3 Tables 8 and 9**.

328 The mis-assembly rate of SOAPdenovo is quite low (~0.1%). Obviously,
329 GapCloser improved the assembly indices dramatically, including total length, N50,
330 and genome fraction. However, it introduced >20 times the amount of mis-assemblies.
331 A large part of these mis-assemblies can be corrected, but a significant proportion still
332 remained. Therefore, we decided to remove the GapCloser step for construction of the
333 pan-genome sequence. In addition, mis-assemblies in the SOAPdenovo results could
334 not be removed by the correction procedure, indicating that most of them are reliable
335 assembled contigs. Hence, we concluded that assembly errors should be very low.

336 **Construction of the pan-genome sequences**

337 We constructed the pan-genome of rice with the Nipponbare RefSeq and

338 non-redundant novel sequences in the assembly of 3,010 rice accessions. The
339 Nipponbare RefSeq was selected because (1) this genome has relatively good
340 annotation and (2) it is widely used for current rice studies, which enables our
341 pan-genome outcome to be integrated easily with current and historical rice studies
342 that have employed the same, giving researchers the ability to easily identify
343 presence/absence of their gene(s) of interest in the rice accessions involved in their
344 study, as well as the attributes of these genes (core, distributed, or GJ-specific, etc.).

345 Blast was used to evaluate the pan-genome. The global identity (G_{iden}) of contig C
346 was calculated as follow:

$$347 \quad G_{iden}(C) = \sum_{i=1}^N w_i \cdot M_i \cdot P_i / L_c$$
$$348 \quad w_i = \begin{cases} 1, & HSP_i \text{ doesn't overlap with any of } HSP_1, \dots, HSP_{i-1}; \\ 0, & HSP_i \text{ overlaps with at least one of } HSP_1, \dots, HSP_{i-1}; \end{cases}$$

349

350 where M_i is length of the HSP_i ; P_i is the percent identity of HSP_i ; L_c is the length of
351 the contig C ; and w_i is a weight indicating whether HSP_i overlaps with previous HSPs.

352 Using this method, we can retrieve novel sequences with an identity cutoff at any
353 value (0.3, 0.5, or 0.7, etc.) in comparison to the Nipponbare RefSeq. Similar methods
354 were used to remove redundant sequences. Like CD-HIT, we use a 'longest sequence
355 first' list removal algorithm to remove sequences above a given identity (see CD-HIT
356 software for details). The difference is that we calculate the global identity of two
357 contigs (A and B) based on NCBI-blast similar to the above method:

$$358 \quad G_{iden}(A, B) = G_{iden}(B, A) = \sum_{i=1}^N w_i \cdot M_i \cdot P_i / \min(L_A, L_B)$$

359
$$w_i = \begin{cases} 1, & \text{HSP}_i \text{ doesn't overlap with any of } \text{HSP}_1, \dots, \text{HSP}_{i-1}; \\ 0, & \text{HSP}_i \text{ overlaps with at least one of } \text{HSP}_1, \dots, \text{HSP}_{i-1}; \end{cases}$$

360

361 where M_i is length of the HSP_i ; P_i is the percent identity of HSP_i ; L_A and L_B are the
362 lengths of the contigs A and B ; and w_i is a weight indicating whether HSP_i overlaps
363 with previous HSPs. The longest contig within each cluster is selected as the
364 representative. Combining these two steps, we can retrieve non-redundant novel
365 sequences at any identity cutoff of P (these sequences have global identities below P
366 in comparison with IRGSP genome sequences, as well as among the sequences
367 themselves).

368 **Pseudogene detection**

369 The sequences of predicted single-exon genes from the novel sequences were
370 extracted and aligned to all other gene sequences using NCBI-blastn with E-value =
371 $1e-5$. A novel single-exon gene with >90% of its full length similar to another gene (a
372 multi-exon novel gene or a Nipponbare gene) is defined as a candidate pseudogene.
373 As a result, 1,030 of the 12,465 genes might be candidate pseudogenes.

374 **Length distribution of novel genes**

375 We checked the gene length differences of the novel genes by comparing them with
376 Minghui 63 (MH63)/Zhenshan 97 (ZS97) genes²³. The MH63 or ZS97 protein
377 sequences (for only the longest ORF of a gene with alternative transcripts removed)
378 were first aligned to all pan-genome proteins; then, each MH63 or ZS97 protein with

379 its best hit (measured by E-value and global identity defined as
380 “ $2 * \text{aligned_length} / (\text{query_length} + \text{target_length})$ ”) was considered as a pair. Those
381 pairs with the pan-genome proteins in multiple pairs were further removed, forming a
382 “single-copy” gene pair set. This set includes 2,474 MH63 genes and 2,441 ZS97
383 genes, each of which formed a pair with its corresponding novel gene; then, the length
384 ratio was calculated and the density was plotted using a log₂ scale (**Extended Data**
385 **Fig. 5**). Generally, the distribution should be symmetric: ratio >0 means the novel
386 gene is longer while a ratio <0 means the novel gene is shorter.

387 **Read mapping to the pan-genome**

388 Mapping raw reads of all rice accessions to the pan-genome sequences is an essential
389 step in our pan-genome analysis. With the mapping results, we can (1) evaluate the
390 sequencing quality based on the percentage of mapped reads and the comparison of
391 sequencing and mapping depths; (2) determine presence/absence of pan-genome
392 contigs in each rice accession; and (3) determine presence/absence of each gene in
393 each rice accession.

394 We compared several mapping tools including SOAP version 2.21²⁴, Bowtie2
395 version 2.2.3²⁵, and BWA version 0.7.10⁸ (both ‘bwa aln’ usage and ‘bwa mem’ usage)
396 based on simulating reads from the 93-11 genome
397 (<http://rice.genomics.org.cn/rice2/link/download.jsp>) and mapping reads to the
398 Nipponbare RefSeq. A gold standard of alignments of the 93-11 genome and
399 Nipponbare RefSeq was built with MUMmer³. Finally, ‘bwa mem’ was selected for

400 all genomic mapping tasks in our pan-genome analysis. The mapping depth and
401 mapping coverage for the Nipponbare RefSeq of each rice line were calculated with
402 Qualimap version 2.0²⁶ and bamUtil (<https://github.com/statgen/bamUtil>),
403 respectively.

404 **Evaluation of gene presence/absence detection**

405 Gene presence/absence detection is a necessary step in high-resolution pan-genome
406 analyses. Previous pan-genome studies in bacteria and rice assembled and annotated
407 each individual genome separately. Here, we compared our method to a recent rice
408 pan-genome study with three representative rice accessions, including the Nipponbare
409 RefSeq²⁷. In their study, they assembled the genomes from deep sequencing and
410 multiple sequencing libraries. They were able to assemble 81.3~82.5% of non-N
411 bases of each genome (81.8% for Nipponbare). This number increased to 88.5~91.4%
412 (91.4% for Nipponbare) if Ns are considered, and they predicted 39,083 genes for the
413 Nipponbare genome. Therefore, the total gene number should be 42,852~47,779 if the
414 entire genomes are assembled and the gene density remains the same (actually, there
415 should be fewer genes in the remaining sequences, which are mostly composed of
416 repetitive and low-complexity sequences). Nevertheless, we can estimate the
417 sensitivity and specificity of gene presence-absence detection approximately. The
418 gold standard annotation for the Nipponbare RefSeq has 35,633 annotated genes.
419 Assuming that all genes on the assembled sequences are correctly predicted (as we
420 also used reference annotations directly in our analyses), the sensitivity can be

421 estimated (as the assembled fraction of the genome) to be about 81.3% or 91.4%. If
422 all the reference genes can be predicted, the specificity should be
423 $35,633 / (>39,083 / 0.914) = 83.1\%$.

424 We then evaluated the accuracy of our method based on the result of CX140 (an
425 independent accession of Nipponbare with sequencing depth at 19x). We detected
426 41,039 genes present in the CX140 *de novo* assembly, including 34,759 reference
427 genes and 6,280 novel genes. The sensitivity is $34,759 / 35,633 = 97.5\%$. The mapping
428 coverage of the CX140 genome is 98.4%; therefore, we think we can capture almost
429 all genes for which mapping evidence exists. The specificity can be estimated as
430 $34,759 / 41,039 = 84.7\%$, which is higher than with traditional pan-genome methods²⁷.

431 However, there are still a significant number of falsely discovered genes. In-depth
432 studies suggest that all of these gene regions show high similarity (>90%) to the
433 reference genome, indicating that these corresponding regions in the reference
434 genome contained no genes, but that we are able to predict genes on similar sequences.
435 Such gene calls might be the false positives from gene predictions, as in the previous
436 work²⁷ in which 39,083 genes were predicted on the incomplete sequences, ~3,000
437 genes more than the RAP annotation²². Alternatively, this might be partially attributed
438 to gene loss in the reference due to SNPs and small indels. This might be a
439 short-coming of our mapping-based pan-genome study. Nevertheless, these false
440 positives are still based on gene sequences and are not random, and sequences of the
441 genes are indeed present in the genome. We therefore concluded that our
442 mapping-based method has relatively good accuracy with very high sensitivity and

443 reasonable specificity.

444 **Estimating size of the rice pan-genome**

445 The sizes of the pan-genome, core gene families, and candidate core gene families
446 were estimated based on simulations. We randomized the order of the 453 rice
447 accessions for 500 times. Each time, we counted the number of core gene families and
448 pan-genome for the first i accessions ($i = 1, 2, \dots, 453$) based on the predefined order.

449 **Fig. 4c** shows the simulation results. The lighter lines stand for the results from 500
450 times randomization and the dark lines stand for the mean values. This showed that
451 the total number of gene families of the rice pan-genome stabilized when the number
452 of accessions was larger than 100.

453 **The average gene/ gene family difference between two accessions**

454 The average gene / gene family difference between two accessions (**Fig. 4e** and
455 **Extended Data Fig.7e**) were calculated as the average of all combinations of each 2
456 of 453 accessions. The average proportions were calculated as the number of such
457 differentiating gene families adjusted by the average gene / gene family numbers held
458 in common by the two major groups.

459

460 **Phylogenetic analysis based on gene (or gene family) PAV**

461 The core genes (or gene families) present in all rice accessions by definition provide
462 no variation. Only the distributed genes (or gene families) were used for the

463 phylogenetic study. The gene (or gene family) presence/absence information of the
464 453 rice accessions was arranged as a 0-1 matrix with each line representing a gene
465 (or gene family) and each column representing a rice accession. The PARS program
466 within PHYLIP (<http://evolution.genetics.washington.edu/phylip.html>) was used to
467 infer the phylogenetic relationship from the presence/absence matrix. The
468 phylogenetic tree was subsequently plotted with the APE²⁸ package in R.

469 **Inferring gene and gene family age with 446 wild rice genomes**

470 Most of the more than 10,000 novel genes were assigned an age of PS13. It is unlikely
471 that such a large number of genes arose within less than the 10,000 years of rice
472 domestication. We therefore inferred that these genes already existed in the wild
473 progenitors of rice; to check this point, we interrogated the whole-genome sequencing
474 data of 446 wild rice accessions²⁹. In Huang et al.'s paper²⁹, 446 wild accessions were
475 previously classified into three groups: Or-I (OR-XL), Or-II (OR-Int), and Or-III
476 (OR-GL). We used the “map-to-pan” strategy³⁰ to study if *O. sativa* genes exist in
477 Or-Int, Or-XL, and Or-GL. To overcome the shortcoming of insufficient sequencing
478 depth of each accession, sequencing data of the same group (Or-Int, Or-XL and
479 Or-GL) were merged and then mapped to the pan-genome sequences of *O. sativa* L.,
480 and genes with both gene body coverage >0.95 and CDS coverage >0.95 were
481 considered as present. As a result, we found that 98.95% of *O. sativa* genes could be
482 detected in wild rice, including >99.9% of the core genes and 96.9% of the distributed
483 genes. Moreover, 437 of the 528 XI-private genes, 110 of the 132 GJ-private genes,

484 56 of the 61 cA-private genes, and 41 of the 48 cB-private genes could be detected in
485 wild rice. Genes found in wild rice that were previously labeled as PS13 were
486 assigned an age of PS12.

487

488

489 **References**

- 490 1 The 3K RGP. The 3,000 rice genomes project. *GigaScience* **3**, 1-6 (2014).
- 491 2 McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing
492 next-generation DNA sequencing data. *Genome Res.* **20**, 1297-1303 (2010).
- 493 3 Delcher, A. L., Salzberg, S. L. & Phillippy, A. M. Using MUMmer to identify similar regions
494 in large sequence sets. *Curr. Protoc. Bioinf.* 10.13. 11-10.13. 18 (2003).
- 495 4 Zhao, K. *et al.* Genome-wide association mapping reveals a rich genetic architecture of
496 complex traits in *Oryza sativa*. *Nat. Commun.* **2**, 467 (2011).
- 497 5 Xu, X. *et al.* Resequencing 50 accessions of cultivated and wild rice yields markers for
498 identifying agronomically important genes. *Nat. Biotechnol.* **30**, 105-111 (2012).
- 499 6 Alexandrov, N. *et al.* SNP-Seek database of SNPs derived from 3000 rice genomes. *Nucleic*
500 *Acids Res.* **43**, D1023 (2015).
- 501 7 Sherry, S. T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**, 308
502 (2000).
- 503 8 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform.
504 *Bioinformatics* **25**, 1754-1760 (2009).
- 505 9 McCouch, S.R. *et al.* Open access resources for genome-wide association mapping in rice. *Nat.*
506 *Commun.* **7**, 17532 (2016).
- 507 10 Zaykin, D. V., Meng, Z. & Ehm, M. G. Contrasting linkage-disequilibrium patterns between
508 cases and controls as a novel association-mapping method. *Am. J. Hum. Genet.* **78**, 737-746
509 (2006).
- 510 11 Kang, H. M. *et al.* Variance component model to account for sample structure in genome-wide
511 association studies. *Nat. Genet.* **42**, 348-354 (2010).
- 512 12 Turner, S. qqman: QQ and manhattan plots for GWAS data. *R package version 0.1* **2** (2014).
- 513 13 Chen, K. *et al.* BreakDancer: an algorithm for high-resolution mapping of genomic structural
514 variation. *Nat. Methods* **6**, 677-681 (2009).
- 515 14 Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end and split-read
516 analysis. *Bioinformatics* **28**, i333-i339 (2012).
- 517 15 Chong, Z. *et al.* novoBreak: local assembly for breakpoint detection in cancer genomes. *Nat.*
518 *Methods* **14**, 65-67 (2017).
- 519 16 Kawahara, Y. *et al.* Improvement of the *Oryza sativa* Nipponbare reference genome using next
520 generation sequence and optical map data. *Rice* **6**, 4 (2013).
- 521 17 Sasaki, A. *et al.* Green revolution: a mutant gibberellin-synthesis gene in rice. *Nature* **416**,
522 701-702 (2002).

523 18 Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-read de novo
524 assembler. *GigaScience* **1**, 18 (2012).

525 19 Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly using de
526 Bruijn graphs. *Genome Res.* **18**, 821-829 (2008).

527 20 Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to
528 single-cell sequencing. *J. Comput. Biol.* **19**, 455-477 (2012).

529 21 Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for
530 genome assemblies. *Bioinformatics* **29**, 1072-1075 (2013).

531 22 Ohyanagi, H. *et al.* The Rice Annotation Project Database (RAP-DB): hub for *Oryza sativa*
532 ssp. japonica genome information. *Nucleic Acids Res.* **34**, D741-D744 (2006).

533 23 Zhang, J. *et al.* Extensive sequence divergence between the reference genomes of two elite
534 indica rice varieties Zhenshan 97 and Minghui 63. *Proc. Natl. Acad. Sci. USA* **113**,
535 E5163-E5171 (2016).

536 24 Li, R. *et al.* SOAP2: an improved ultrafast tool for short read alignment. *Bioinformatics* **25**,
537 1966-1967 (2009).

538 25 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,
539 357-359 (2012).

540 26 García-Alcalde, F. *et al.* Qualimap: evaluating next-generation sequencing alignment data.
541 *Bioinformatics* **28**, 2678-2679 (2012).

542 27 Schatz, M. C. *et al.* Whole genome de novo assemblies of three divergent strains of rice,
543 *Oryza sativa*, document novel gene space of aus and indica. *Genome Biol.* **15**, 506 (2014).

544 28 Paradis, E., Claude, J. & Strimmer, K. APE: analyses of phylogenetics and evolution in R
545 language. *Bioinformatics* **20**, 289-290 (2004).

546 29 Huang, X. *et al.* A map of rice genome variation reveals the origin of cultivated rice. *Nature*
547 **490**, 497 (2012).

548 30 Hu, Z. *et al.* EUPAN enables pan-genome studies of a large number of eukaryotic genomes.
549 *Bioinformatics*, btx170 (2017).

550

551