# A syntelog-based pan-genome provides insights into rice domestication and de-domestication

Dongya Wu,[1,2,3,#] Lingjuan Xie,[2,#] Yanqing Sun,[2] Yujie Huang,[2,4] Lei Jia,[2] Chenfeng Dong,[2] Enhui Shen,[1,2] Chu-Yu Ye,[1,2] Qian Qian,[4,*] Longjiang Fan[1,2,*]

[1] Hainan Institute of Zhejiang University, Sanya 572025, China
[2] Institute of Crop Science, Zhejiang University, Hangzhou 310058, China
[3] Center for Evolutionary & Organismal Biology, Zhejiang University, Hangzhou, 310058, China
[4] State Key Laboratory of Rice Biology, China National Rice Research Institute, Hangzhou 310006, China

* Correspondence:
Qian Qian (qianqian188@hotmail.com); Longjiang Fan (fanlj@zju.edu.cn)

#these authors contribute equally.

**Abstract**
Asian rice is one of the world's most widely cultivated crops. Large-scale resequencing analyses have been undertaken to explore the domestication and de-domestication genomic history of Asian rice, but the evolution of rice is still under debate. Here, we construct a syntelog-based rice pan-genome by integrating and merging 74 high-accuracy genomes based on long-read sequencing, encompassing all ecotypes and taxa of *Oryza sativa* and *Oryza rufipogon*. Analyses of syntelog groups illustrate subspecies divergence in gene presence-and-absence and haplotype composition and identify massive genomic regions putatively introgressed from ancient Geng/*japonica* to ancient Xian/*indica* or its wild ancestor, including almost all well-known domestication genes and a 4.5-Mb centromere-spanning block, supporting a single domestication event in rice. Genomic comparisons between weedy and cultivated rice highlight the contribution from wild introgression to the emergence of de-domestication syndromes in weedy rice. This work highlights the significance of inter-taxa introgression in shaping diversification and divergence in rice evolution and provides an exploratory attempt by utilizing the advantages of pan-genomes in evolutionary studies.

**Key words**
Rice pan-genome; syntelog; domestication; de-domestication; introgression

## Introduction

As one of the most important calorie sources, Asian rice (*Oryza sativa*), which is considered to be domesticated from its wild progenitor (*Oryza rufipogon*), is widely grown worldwide. Despite its indispensable roles in food supply and fundamental studies about plant biology, the origination, domestication, and subsequent diversification of rice have been under debate for decades, although a considerable amount of archaeological and genetic evidence has been proposed to infer the evolutionary trajectory of rice (Molina et al., 2011; Huang et al., 2012; Civáň et al., 2015; Gross and Zhao, 2014; Choi et al., 2017; Carpentier et al., 2019; Zhang et al., 2021). Multiple taxa or groups in classification, recurrent artificial hybridization during breeding, long-distance dispersal by global trade and other factors have hindered our understanding of rice evolution. The most disputed issue is whether domestication events have happened only once or independently multiple times. Regardless, many domestication- and improvement-related genes have been identified to underlie domestication syndromes, such as plant architecture (*Prog1*), grain shattering (*sh4*), awn length (*An-1* and *LABA1*), pericarp color (*Rc*) and dormancy (*Sdr4*) (Chen et al., 2019). Recently, the issue of rice feralization or de-domestication has attracted great attention in both agricultural production and basic biology because the de-domesticated ecotype of rice (weedy rice, *Oryza sativa* ssp. *spontanea*) has severely threatened rice yield and quality as a commonly seen weed in paddy fields. The genetic resources of weedy rice also show potential for use in enhancing abiotic stress adaptation in rice breeding (Sun et al., 2022). How atavism occurs in rice is an intriguing biological question that has further extended and complicated the evolutionary history of rice (Wu et al., 2021). Previous studies have suggested that the genomes of weedy rice were mostly derived from local cultivated rice by recurrent and independent de-domestication events and that genetic introgression from wild rice may have contributed to weediness (Song et al., 2014; Li et al., 2017; Qiu et al., 2017; Sun et al., 2019; Qiu et al., 2020).

Pan-genomic studies have been conducted in a wide range of crops, including rice (Zhao et al., 2018; Wang et al., 2018b; Qin et al., 2021; Zhou et al., 2020; Zhang et al., 2022; Shang et al., 2022). By comparing *de novo* assemblies, large SVs have been discovered, underlying important traits that could not be explained by small-scale variations. How to utilize pan-genomes in evolutionary studies has been relatively little explored. Here, we integrate high-accuracy rice genomes covering all ecotypes and taxa and high-depth resequenced genomes of wild rice to revisit the origin, domestication, diversification and de-domestication processes of rice based on a syntelog-based pangenome. Whole-genome mosaic haplotype maps intuitively reveal massive introgression footprints of domesticated genomic blocks from the proto-GJ (initial domesticates of subspecies Geng/*japonica*) to XI (subspecies Xian/*indica*) ancestor, strongly supporting the hypothesis of single domestication in rice. Structural variations between weedy and cultivated rice indicate that the introgression events from wild progenitors to different cultivated rice groups probably underlie the parallel

77  convergence in weedy traits (e.g. pericarp and hull color). Briefly, our study
78  comprehensively investigates the complex relationship among different rice taxa and
79  ecotypes from a pan-genome view and highlights the significance of genomic
80  introgression in both rice domestication and de-domestication.
81
82

83  **Results**

84

85  **High-quality rice genome assemblies**

86

87  To fully capture the genomic diversity and dynamics in rice domestication,
88  improvement and further feralization, we created a panel of high-quality rice genome
89  assemblies, including newly generated assemblies of 11 weedy and one cultivated
90  accession, using PacBio HiFi mode with an average sequencing depth of 32.4×.
91  Contigs were anchored on chromosomes using a reference-guided approach, and Hi-C
92  interaction confirmed the order and orientation accuracy for four accessions
93  (Supplementary Fig. 1). The average contig N50 of the newly assembled genomes
94  was 19.16 Mb (from 10.95 Mb to 30.43 Mb), and the LAI score was on average 21.71
95  (from 20.2 to 23.9), equivalent to those of previous assemblies using PacBio CLR
96  mode or Nanopore sequencing (Supplementary Fig. 2a). Averagely, 97.32% of the
97  4896 core conserved Poales genes (BUSCO) were assembled (Supplementary Table
98  1). The whole-genome synteny against the reference assembly Nipponbare and
99  gapless assembly MH63RS3 (Song et al., 2021) suggested high completeness
100  (Supplementary Fig. 3).
101
102  Before adopting more assemblies into the construction of the rice pan-genome, we
103  systematically evaluated the base-level accuracy and assembly completeness on rice
104  genomes, including recently released assemblies (Qin et al., 2021; Zhou et al., 2020;
105  Zhang et al., 2022). The $k$-mer-based assembly validation results revealed higher
106  assembly consensus quality values (QVs) for HiFi assemblies generated in this study
107  (average QV = 44.16) than those for previous assemblies (Fig. 1a; Supplementary Fig.
108  4a). We quantified the assembly accuracy by calling homozygous single nucleotide
109  polymorphisms (SNPs) and short insertions and deletions (InDels) by mapping
110  available NGS reads for each accession against its own assembly. At the single-base
111  level, averagely the HiFi assemblies showed fewer errors than the PacBio CLR mode
112  and Nanopore sequencing (Fig. 1a; Supplementary Fig. 4b). In terms of InDels, HiFi
113  assemblies showed no obvious differences from CLR mode assemblies, but there
114  were fewer InDels in HiFi than in Nanopore assemblies. The annotation of the
115  potential assembly errors for each accession suggested high (stop loss and gain, start
116  loss, and frame-shift variants) and moderate effects (inframe insertion/deletion and
117  missense variants) in the predicted gene models (Supplementary Fig. 4b). The low
118  assembly quality at the base level directly interferes with the accuracy of haplotype
119  inference, especially for Nanopore-based assemblies despite further polishing using
120  short reads; thus, the recently released Nanopore-based assemblies of cultivated rice

121    were excluded. Given that the available assemblies for wild rice using PacBio
122    sequencing were only for W2014 (Ma et al., 2020) and IRGC106162 (Xie et al.,
123    2020), nine Nanopore-based wild assemblies (Shang et al., 2022) were adopted.
124    Finally, 11 wild, 51 cultivated and 12 weedy rice assemblies were used in the
125    following pan-genomic analysis (Supplementary Table 1).
126    Phylogeny based on 11.6 million whole-genome SNPs of the 74 genomes revealed
127    four *aromatic* (aro), four tropical (trp), two subtropical (subtrp) and 13 temperate (tmp)
128    accessions in the GJ subspecies ($n = 23$ in total) and four *aus*, three XI2, seven XI3,
129    ten XI1A and 16 XI1B accessions in the XI subspecies ($n = 40$) (Fig. 1b).
130    Whole-genomic features, e.g., genome size, number of annotated genes, and
131    transposon element size and proportion, were significantly differentiated between XI
132    and GJ (Supplementary Fig. 2b). The wild population includes two accessions from
133    Or-3 (generally considered to be the ancestral group of GJ), four from Or-2, three
134    from Or-1 (ancestral group of XI), and two from Or-4 (Fig. 1b). The
135    representativeness of the 74 genomes regarding diversity was validated by a kinship
136    analysis using a global panel of approximately seven thousand rice accessions
137    (Supplementary Fig. 5), suggesting that the rice genomes used have covered all major
138    taxa and all ecotypes. This provides a good opportunity to revisit the evolutionary
139    trajectory of rice from a pan-genome view.
140
141
142    **Syntelog-based pan-genome of rice**
143
144    Compared to variation maps obtained by mapping whole-genome resequencing short
145    reads against one reference genome, *de novo* assemblies provide accurate and
146    complete haplotype-resolved genetic and predicted protein sequences as well as
147    genomic coordinates along chromosomes. To incorporate positional information, we
148    constructed a synteny-based pan-genome by clustering approximately 3.10 million
149    genes from the 74 rice genomes with SynPan (see Methods). Pairwise alignments for
150    each pair of genomes were performed first to identify inter-individual syntelogs
151    (syntenic orthologs) and merged together. Compared to the fast aligner Diamond
152    (Buchfink et al., 2021), which is designed for high-performance analysis of big
153    sequence data, pairwise alignment using BLASTP identified more syntelogs between
154    two genomes, with an average addition of 669.1 pairs per genome alignment
155    (Supplementary Fig. 6). Therefore, the syntelog datasets based on the BLASTP
156    approach were used in the pan-genome construction.
157
158    First, a coalescence-free kinship was built based on pairwise whole-genome synteny
159    (Fig. 1b). Generally, the affinity measured by whole-genome synteny among
160    accessions is linearly correlated with that using identity-by-descent (Pearson's
161    correlation = 0.886, $P$ value < 2.2e-16) (Fig. 1c). For GJ groups (aro, trp, subtrp and
162    tmp), each group showed a closer relationship to the other GJ groups, and Or-3 was
163    the closest wild-type group. For XI groups, *aus* exhibited apparent differences from
164    the other XI groups (XI2, XI3, XI1A and XI1B) in genomic arrangements (Fig. 1b).

165 Although *aus* is sister to other XI groups and nested within wild group Or-1 in the
166 phylogenetic tree (Fig. 1b), the kinship between *aus* and Or-1 measured by synteny
167 was distant and even larger than that between *aus* and Or-4 (the basal wild outgroup),
168 which implied that the wild ancestor of *aus* is distinct from that of the other XI groups
169 (Supplementary Fig. 7a). Notably, some XI accessions suggested closer affinity with
170 GJ, which reflected the inter-subspecies hybridization in modern breeding (e.g., Y58S
171 and its offspring J4115) (Fig. 1b; Supplementary Fig. 7b). Y58S is an XI-type
172 photothermosensitive genic-male-sterile (PTGMS) line with the characteristics of
173 high-light-efficiency use and disease and stress resistance, and it is widely used for
174 the breeding of two-line hybrid rice varieties, especially super hybrids. The GJ
175 accession Lemont is one of the parental lines used in the breeding of Y58S (China
176 Rice Data Center, https://ricedata.cn/).
177
178 Based on whole-genome synteny, 175,528 syntelog groups (SGs) were clustered, with
179 13,908 core (present in the genomes of all accessions), 14,423 soft-core (present in
180 the genomes of >90% accessions), 62,425 dispensable (present in the genomes of less
181 than 90% of all but at least two accessions) and 84,772 private SGs (only present in a
182 single genome) (Fig. 2a; Supplementary Fig. 8a). The SG size is 1.62 times larger
183 than the number of orthogroups or ortholog groups (OGs) clustered by the Markov
184 clustering (MCL) algorithm ($n = 67,080$ when the inflation parameter is 1.5, including
185 15,749 core, 11,725 soft-core, 18,901 dispensable and 20,705 private OGs). Even
186 though the inflation parameter was set as 2.5, the OG size increased to $n = 72,226$,
187 which included only 41.1% of the SG size (Supplementary Fig. 9). In a perfect
188 ortholog group, one gene from one genome is expected, despite individual-specific
189 duplication. The MCL method does not perform well in distinguishing paralogs
190 against orthologs, especially in the core and soft-core groups (Supplementary Fig. 8b).
191 By taking advantage of the genomic coordinates of genes in assemblies, the SGs
192 provide a more precise and accurate ortholog classification.
193
194 Approximately 34.5%, 34.9%, 27.7% and 2.8% of genes were assigned to core,
195 soft-core, dispensable and private SGs, respectively (Supplementary Fig. 10a). Protein
196 domains could be identified using InterProScan (Jones et al., 2014) in a total of 81.5%
197 and 68.9% of core and soft-core genes, which is nearly twice as high as 36.9% and
198 28.7% for dispensable and private genes, respectively (Supplementary Fig. 10b).
199 Protein domain gain-and-loss variation was found within a single SG, indicating
200 functional diversification in rice evolution. In total, 1.10% of core genes from 14.1%
201 of core SGs suggested domain gain-and-loss, and these percentages were lower than
202 those for soft-core and dispensable types (2.0% of soft-core genes in 19.5% of
203 soft-core SGs and 5.0% of dispensable genes in 18.9% of dispensable SGs)
204 (Supplementary Fig. 10c). For example, two adjacent genes, *SaM* and *SaF,* encode a
205 ubiquitin-like modifier E3 ligase-like protein and an F-box protein, respectively, and
206 their interactions are responsible for XI-GJ hybrid male sterility (Long et al., 2008).
207 The SGs of *SaM* and *SaF* were both soft-core genes present in 68 and 72 genomes,
208 respectively. The domains of all *SaF* syntelogs are completely annotated, while the

209   domains of 29 syntelogs from *SaM* SG were not found, which suggested the dynamics
210   domain gain-and-loss in conserved SGs.
211
212   Only 49.5% ($n$ = 86,974) and 61.3% ($n$ = 107,565) of SGs were present in the GJ and
213   XI subspecies, respectively, indicating the large genomic diversity of wild accessions
214   and genetic bottlenecks due to artificial selection or genetic drift (Fig. 2a). In total,
215   4,662 SGs showed presence-and-absence (PAV) frequency biases between GJ and XI,
216   with a frequency difference greater than 0.6, including 1,918 SGs absent from the
217   reference genome Nipponbare. For example, the key genes in a casbene-derived
218   diterpenoid biosynthetic gene cluster DGC7 on chromosome 7 suggested differential
219   PAVs between GJ and XI, where *CYP71Z21*, *TPS28* and *CYP71Z2* were almost
220   absent in XI but fixed in GJ (Fig. 2b), which may be related to the differential
221   responses of subspecies to biotic stress (Zhan et al., 2020). *RePRP1* and *RePRP2* are
222   functionally redundant suppressors of root cell expansion (Tseng et al., 2013). Both
223   *RePRP1* and *RePRP2* were present in GJ accessions, while only *RePRP2* was found
224   in most XI and wild accessions, which implied that the copies of RePRP genes may
225   underlie the differential root development in different subspecies (Fig. 2b).
226
227
228   **Rice pan-NLRome**
229
230   Pan-genomes provide an opportunity to uncover the diversity of highly variable gene
231   families, such as those encoding nucleotide-binding leucine-rich repeat (NLR)
232   proteins related to disease resistance, across species (so-called pan-NLRome). A total
233   of 37,079 NLR genes in 74 rice genomes were identified, ranging from 452 (W2014
234   from wild group Or-4) to 532 (FH838 from group XI1B) NLRs per genome. Distinct
235   from the NLR composition in the *Arabidopsis thaliana* pan-NLRome, no TIR-NLR
236   (TNL) genes were found in the rice genomes. Rice NLRs were categorized into three
237   types: CNL (including CC-NB-LRR or CC-NB), NL (including NB-LRR or NBS)
238   and null (NLR genes identified by syntelogs whose encoding proteins contain no
239   canonical NBS domain). The sizes of NLRs in cultivated and weedy genomes were
240   both significantly larger than those in the wild group ($P$ = 2.8e-5 and 4.3e-5, Student's
241   *t* test) (Supplementary Fig. 11a), which could be the consequence of
242   disease-resistance gene aggregation during domestication and improvement. However,
243   the relatively low assembly quality of the wild genomes may also be related to this
244   difference, considering that the completeness of assemblies was significantly
245   correlated with the NLR size (Pearson's correlation = 0.53, $P$ value = 9.7e-7)
246   (Supplementary Fig. 11d). At the subspecies level, although NLR gene numbers were
247   similar between XI and GJ, the XI genomes contained more NLs than those in GJ ($P$
248   = 4.3e-10, Student's *t* test), and GJ had more CNLs ($P$ = 3.6e-6, Student's *t* test)
249   (Supplementary Fig. 11b). Null NLRs without canonical NBS domains were more
250   abundant in GJ than in XI ($P$ = 1.6e-10, Student's *t* test), implying that more NLRs in
251   GJ may degenerate functionally by losing domains (Supplementary Fig. 11b).
252

253 Adopting the definition by Wang et al. (2019) of an NLR cluster containing more than
254 two NLR genes distributed within a 300-kb genomic region, 43.2%-59.9% of NLR
255 genes in rice were located in such clusters, where GJ showed more clusters in both
256 numbers and proportions across all NLRs than XI (Supplementary Fig. 11c).
257 Head-to-head pairing of NLR genes is highly associated with disease resistance in
258 plants, with one NLR acting in effector recognition (known as a sensor) and the other
259 acting in signaling activation (known as a helper). We found 28 to 54 such paired
260 NLRs per genome. A total of 6,008 NLRs encoded at least one non-canonical NLR
261 domain (NBS, LRR and CC), also known as the integrated domain (ID), representing
262 116 distinct Pfam domains and 16.2% of the total NLRs, which was much higher than
263 that in *Arabidopsis thaliana* (5.0%). We identified 480 distinct architectures in the
264 pan-NLRome, of which only 67 were found in the Nipponbare reference genome
265 (IRGSP v1.0). Fewer than 3% of architectures, 12, correspond only to different
266 configurations of the canonical CC, NBS and LRR domains, even though they
267 accounted for the majority (83.8%) of NLRs.
268
269 As expected, the haplotype diversity of NLR SGs was significantly higher than that of
270 non-NLR SGs (Fig. 2c). In total, 0.64% ($n = 238$) of all NLRs were present in only
271 one accession, representing 238 private SGs, while the remaining 10,878 (29.3%),
272 14,760 (39.8%) and 11,203 (30.2%) NLRs grouped into 147 core, 208 soft-core and
273 407 dispensable SGs, respectively. Although the NLR syntelogs among individuals
274 were well defined by genome synteny, within a single SG, the functional types (CNL,
275 NL or null) and structural types (e.g., in clusters or pairs) were diversified,
276 particularly for the core NLR SGs (Supplementary Fig. 12a). Nucleotide diversity for
277 NLRs in core and soft-core SGs was lower than that in dispensable SGs but was not
278 significant. Tajima's *D* values, which indicate balancing and purifying selection,
279 showed no significant differences across different NLR classes, with all classes
280 containing extremes in both directions (Supplementary Fig. 12b).
281
282
283 **Rice origin based on the mosaic genomic map of syntelog haplotypes**
284
285 Although whole-genome single nucleotide variants provide a comprehensive variation
286 landscape in genome evolution, millions of markers could be somewhat redundant
287 because of synonymous mutations. Here, we used the predicted amino acid sequences
288 of genes to dissect the genomic ancestry of each gene in all rice genomes, given that
289 the protein sequences function directly and are degenerated with high tolerance to
290 synonymous mutations. Haplotypes were first assigned for each of the 49,438 SGs
291 whose syntelog members were present in at least ten genomes. The haplotype
292 complexity for each SG based on protein sequences was highly reduced compared
293 with those using full-length gene sequences and coding sequences, as indicated by
294 haplotype diversity (the average haplotype differences between any two members
295 from a single SG), haplotype N100 (total number of unique haplotype sequences) and
296 N90 (the least haplotype number that needs to be included for covering 90% of

297    sequences in an SG) (Supplementary Fig. 13a; Supplementary Fig. 14a). The
298    haplotype number and diversity were higher for core and soft-core SGs than
299    dispensable SGs. On average, 11.2 haplotypes were found for each SG, where
300    haplotype numbers for core ($n = 12.0$) and soft-core SGs ($n = 12.8$) were higher than
301    that for dispensable SGs ($n = 9.5$). Specifically, 147 core SGs were extremely
302    conserved with the fully identical protein sequences of housekeeping genes involved
303    in fundamental biological processes, such as *LEA5* in late embryogenesis (He et al.,
304    2012), *eIF-4A* and *Os-eIF6;1* in translation initiation (Kato et al., 2010), *OsFd1* in
305    photosynthesis (He et al., 2020) and *OsAtg8* in autophagy (Izumi et al., 2015). As
306    expected, the haplotype diversity of XI was significantly higher than that of GJ, with
307    average haplotype diversity values of 0.578, 0.362 and 0.441 for all rice accessions,
308    GJ and XI, respectively (Supplementary Fig. 13b).
309
310    Using a semi-supervised approach, haplotypes were reassigned by comparing their
311    abundance in different groups for each SG, labeled as hapI to hapV, hapR (all other
312    rare haplotypes) and absence (represented by red, blue, orange, yellow, green, dark
313    gray and light gray blocks in Fig. 3a and Supplementary Fig. 15) to represent the
314    ancestral sources of rice haplotypes (Supplementary Fig. 13b). Apparently, the
315    whole-genome haplotype maps visually reflected the shared haploblocks at the
316    individual and group levels. For example, mosaic genomes of SN265, DHX2 and
317    02428 from the GJ tmp group and Y58S, J4115 and FH838 from the XI1B group
318    suggested large introgressed regions from the other subspecies (Fig. 3a;
319    Supplementary Fig. 15). Haplotypes in two Mb-level genomic regions at the head of
320    chromosomes 6 and 7 showed that only XI1B in XI groups was shared with GJ, which
321    was consistent with the fact that group XI1B was mainly composed of modern
322    cultivars that were frequently bred by utilizing genetic resources from other
323    subspecies (Supplementary Fig. 15).
324
325    As observed from the mosaic genomic map, *aus* showed obvious differences from the
326    other XI groups; thus, the *aus* group was excluded from XI in the following analyses
327    (Fig. 3; Supplementary Fig. 15). We used inter-subspecies diversity to quantify
328    haplotype divergence (HDG) between GJ and XI (Supplementary Fig. 13a). The
329    average haplotype divergence between GJ and XI was 0.667, and 46.0% of SGs ($n =$
330    22,745) suggested high haplotype divergence with HDG > 0.8, indicating great
331    divergence between subspecies at the translation level. Most (14/21) well-known
332    improvement genes (Chen et al., 2019) showed high divergence, such as *DEP1* (HDG
333    = 1.000), *Sd1* (0.940*)*, *TAC1* (0.930), *GS5* (0.913), *GW6a* (0.925), *TGW6* (0.962),
334    *GW7* (0.954), *GW8* (0.976), *Ghd7* (0.963), *Ghd8* (0.984), *Hd1* (0.979), *NRT1.1B*
335    (0.977), *DRO1* (0.850), and *Chalk5* (0.933), implying independent selection for yield-
336    and flowering-related genes in the improvement of XI and GJ. Notably, 720 SGs
337    showed weak divergence with HDG < 0.2, including essential domestication genes
338    *Prog1* (HDG = 0.025), *GAD1* (0.110) and *sh4* (0.155). The SGs with low divergence
339    (HDG < 0.5) were clustered into local genomic blocks, with a total length of 23.38
340    Mb (6.24% of Nipponbare assembly) in 73 blocks, including 2,786 SGs (6.90% of all

SGs in Nipponbare), implying putative genomic introgression between XI and GJ (Fig. 3a and 3b; Supplementary Fig. 15; Supplementary Tables 2 and 3). Typically, a total of 18 blocks were beyond 300 kb, and the largest three blocks were on chromosomes 5, 8 and 4, spanning over 4.48, 2.22 and 1.75 Mb, respectively. Introgression and incomplete lineage sorting (ILS) would both result in haplotype similarity and low divergence between sequences from two lineages. To distinguish introgression from ILS, which is more randomly distributed along chromosomes (Wu et al., 2022b; Edelman et al., 2019), we tested the significance of lowly divergent SG clustering by 100000-times random sampling. The significant nonrandom distribution of the lowly divergent SGs in these blocks implied that inter-subspecies introgression caused low GJ-XI divergence, rather than ILS (Fig. 3c; Supplementary Fig. 15). Additionally, as expected, the relative divergence indicated by the synonymous substitution rate ($K_s$) between XI and GJ genes in putative introgression blocks was significantly lower than that of genes in adjacent genomic regions (Supplementary Fig. 16).

Most domestication genes (9/10) underlying key domestication syndromes (Chen et al., 2019) were found in the introgression blocks (Supplementary Fig. 15; Supplementary Table 2). On chromosome 4, four introgressed blocks larger than 300 kb were found, including domestication genes *LABA1* and *An-1* responsible for awn presence-and-absence and length in blocks 4#1 and 4#2 (Luo et al., 2013; Hua et al., 2015), *sh4* for grain shattering in block 4#3 (Li et al., 2006), and *Bh4* for hull color in block 4#4 (Zhu et al., 2011)(Fig. 3a). On chromosome 5, *GW5* from block 5#3 is a QTL for grain width and weight (Liu et al., 2017). On chromosome 7, *Prog1* was located in block 7#1, underlying the transition from prostrate plant architecture to erectness during domestication (Jin et al., 2008). *GAD1*, encoding a secreted awn development-related peptide, and *IPA1*, which is considered a typical improvement gene controlling ideal plant architecture and immunity, were both located in block 8#1 on chromosome 8 (Jiao et al., 2010; Wang et al., 2018a). Interestingly, except for *IPA1*, no other improvement and diversification genes were found in the introgression blocks, which suggested that introgression events occurred in the initial period of domestication (Supplementary Table 2).The largest introgression block, 5#1, with a length of 4.48 Mb, spanned the centromere region on chromosome 5 (Supplementary Fig. 15). Although more than three hundred genes were annotated in this region, no known domestication-related genes were found. The presence of block 5#1 with low divergence between XI and GJ was probably due to suppressed recombination over the centromere region rather than linkage by key domestication genes, which provides an ideal clue to subspecies introgression.

We further utilized 184 wild genomes with high whole-genome sequencing depth (averagely >8×), encompassing four groups Or-1 to Or-4, to confirm the introgression and trace the spread routes of domestication haplotypes (Supplementary Fig. 17). Totally the introgression inference in 64 blocks (96.2% in genomic length) were supported by phylogeny (Supplementary Table 3). Phylogenetic trees of these blocks (including all 18 blocks longer than 300 kb) indicated that GJ and XI

accessions were nested within the Or-3 group of wild rice, indicating that Or-3 was the shared wild progenitor group of GJ and XI in most introgression blocks, and that the domesticated alleles in XI were likely to be derived from proto-GJ by introgression with local wild rice from Or-1 (Fig. 3d; Supplementary Table 3). Notably, between the major clade of domestication haplotypes and the Or-3 clade, wild accessions from the Or-1 or Or-2 group were commonly observed (Fig. 3d), which implied that these haplotypes contained relict ancient domesticated alleles, although some introgression events from cultivated to wild rice were observed. We speculate that a gene pool under early domestication was introduced into South Asia and partially maintained in the genomes of present Or-1 or Or-2 wild rice.

Statistical ABBA-BABA tests were also performed to confirm the introgression inference. In a total of 57 blocks (21.2Mb in length), high $f_d$ values were observed in the models with introgression direction from tmp(GJ) to XI groups (Supplementary Figs. 18 and 19). Thus combing haplotype inference, phylogeny, and ABBA-BABA test, 65 blocks (22.6 Mb in length) were finally determined as introgression regions from GJ to XI (Supplementary Table 3). Auxin related pathways were significantly enriched (Supplementary Table 4). Besides well-known domestication genes related awn presence, shattering, and tiller angle, many seed dormancy or germination related genes were observed in introgression blocks (Supplementary Table 5). Block 3#5 from chromosome 3 included a seed dormancy-related gene *OsG* , which has been parallelly selected in multiple crop families (Wang et al., 2018a), and *qLTG3-1*, a major quantitative trait locus controlling low-temperature germinability (Fujino et al., 2008). *OsC1*, which regulates hull pigmentation and pre-harvest sprouting, was located in block 6#3 on chromosome 6. Yield related genes include at least *HOX3* (1#2), *SPL6* (3#6), *GPA3* (3#6), *OsNAC2* (4#4), *D11* (4#7)and *FZP* (7#5).

We noticed that the *aus* group exhibited differences from the GJ and XI groups in most blocks (e.g., 1#1, 2#2 and 5#1). Phylogenetic analysis revealed that XI and *aus* belonged to two separate subclades, although both were nested within Or-1 group (Supplementary Fig. 17a). There were 32 blocks where the closest wild group to *aus* was Or-3, similar to XI and GJ, while in the other 24 and 8 blocks, the Or-1 and Or-4 groups were the wild progenitors of *aus*, respectively (Fig. 3e). In ABBA-BABA analysis, only 23 blocks showed gene flow signals from GJ to *aus* (Supplementary Fig. 18; Supplementary Table 3). Therefore, the *aus* group shared some domesticated alleles, suggesting that the introgression from proto-GJ or Or-3 to XI probably occurred after the divergence between XI and *aus*, and some domesticated alleles or blocks were subsequently transferred from XI to *aus* (e.g., 4#1) (Supplementary Fig. 15). Combining the evidence from whole-genome pairwise synteny (Fig. 1b), the *aus* group differs from other XI groups and has a novel evolutionary process.


**Structural variations during rice de-domestication**

429    Previous studies have found a 0.5-Mb de-domestication genomic island on
430    chromosome 7 from 6.0 Mb to 6.5 Mb that plays essential roles in rice feralization
431    (Qiu et al., 2020). Within this region, the key gene *Rc* regulating the red pericarp and
432    a cluster of seed storage-related genes (six *RAL* and three *LtpL* genes) contribute to
433    the fitness of weedy rice. The genomic synteny in this region was investigated
434    between weedy and cultivated rice. Although weedy rice originated independently and
435    repeatedly from cultivated rice, the genomic landscape of XI weedy accessions in this
436    region showed distinct patterns against their corresponding closest cultivated genomes
437    (Fig. 4a). This weedy pattern was prevalent in the genomes of wild accessions.
438    Notably, a remarkable 10-kb translocation (including one gene encoding an RNA
439    polymerase II transcription subunit) was found in all XI weedy rice when aligned to
440    the Nipponbare assembly but was absent from all XI cultivated rice (except cultivar
441    IRGC34749 with a red pericarp) and all GJ accessions (except Basmati1) (Fig. 4a).
442    This translocation was found in 7 of 11 wild accessions. Thus, it was speculated that
443    genomic introgression of the de-domestication genomic island from wild rice
444    contributed to the feralization of XI cultivated rice and the emergence of XI weedy
445    rice. A phylogenetic tree based on SNPs around the *Rc* region further supported the
446    introgression and suggested that the Or-1 group from Southeast Asia was the ancestral
447    progenitor of the de-domestication genomic island in XI weedy rice (Fig. 4b). For GJ
448    weedy rice, no obvious signals of wild introgression in this de-domestication island
449    were found (Fig. 4a). However, from the phylogeny of *Rc*, weedy accessions were
450    clustered together and nested within some Or-3 wild accessions, which implied that
451    Or-3 may be the donor of the *Rc* haplotype underlying the red pericarp (Fig. 4b).
452    However, it should be noted that some cultivated rice also showed a red pericarp (e.g.,
453    LJ from GJ), suggesting another potential origin from local landraces of red rice.
454
455    To gain insight into detailed structural variations during de-domestication, we
456    compared the genomic sequences between weedy accessions and their closest
457    corresponding cultivated accessions based on phylogeny (Fig. 1b). Given that
458    de-domestication occurred recently (Qiu et al., 2020; Sun et al., 2019), no large
459    chromosome rearrangements were observed from genome synteny, except for the
460    comparison between Tetep and PI653432, owing to the low assembly quality of Tetep
461    (Supplementary Fig. 20). On average, 73,111 small insertions and deletions (InDels,
462    ≤50 bp) and 2,810 large structural variants (SVs, >50 bp) were identified in ten
463    weed-cultivar pairs, spanning genomic regions from 2.2 to 13.8 Mb in total
464    (Supplementary Fig. 21). The number of structural variants was lower for the GJ
465    weedy-cultivated pairs than for the XI and *aus* pairs. This could be a result of the
466    recent origin of GJ weedy rice but more ancient origin of XI and *aus* (Qiu et al., 2020)
467    and more introgression from other taxa or (sub)-species into XI and *aus* cultivated
468    rice. The XI cultivar NJ11 and its weedy descendant CX20 were reported as a typical
469    case of recent de-domestication events (Qiu et al., 2020). When the HiFi-based
470    genome assemblies were compared, a minimum number of SVs ($n = 1,299$) were
471    detected, where six SVs larger than 100 kb were close to peri-centromeric regions,
472    and the largest SV spanned 823 kb and harbored 66 genes on chromosome 7

473 (Supplementary Fig. 22). *OsGWD1*, which is involved in transitory starch degradation
474 in source tissues and is also a positive regulator of rice seed germination, was lost in
475 the weedy CX20 genome, which may be related to the difference in rice quality
476 between cultivated and weedy rice (Wang et al., 2021). Additionally, we noticed the
477 absence of the rice blast resistance gene *Pid2* in weedy CX20. Compared with NJ11,
478 equivalent insertions ($n = 19,340$ for small insertions and $n = 666$ for large insertions)
479 and deletions ($n = 19,410$ for small deletions and $n = 633$ for large deletions) were
480 found in CX20, spanning a total gain-and-loss length of 2.91 and 2.97 Mb,
481 respectively. Gene Ontology analysis suggested that SV-associated genes were
482 enriched in the biological process of reproductive system development (adjusted $P =$
483 0.00065).
484
485 Despite the independent and recurrent de-domestication events observed from the
486 phylogeny, we found 3,614 SGs in which at least one SV was detected in at least four
487 de-domestication lineages from six groups (tmp, *aus*, XI2, XI3, XI1A and XI1B),
488 implying potential convergent genetic mechanisms underlying feralization. Within
489 them, a 14-bp PAV in *Rc* was identified in all XI and GJ weedy lineages (Wu et al.,
490 2021). SVs were found in other domestication-related genes regulating seed shattering,
491 hull color and seed dormancy or germination, of which most were located in
492 regulatory and intron regions (Supplementary Table 6). For shattering, a 2-bp
493 insertion in exon 1 of *sh4* and a 12-bp deletion in exon 1 and a 146-bp deletion in
494 exon 6 of *SHAT1* were found in the *aus* and XI weedy genomes. Despite *Rc*'s role in
495 regulating seed dormancy and germination, SVs or InDels in *OsC1* and *Sdr4* were
496 also found in the GJ and XI weedy genomes, respectively (Fig. 4c; Supplementary
497 Table 6). *OsC1*, a rice R2R3-MYB transcriptional regulator that interacts with *Rc* and
498 *OsVP1*, plays an important role in regulating preharvest sprouting tolerance in red
499 pericarp rice (Wang et al., 2020). Two different variants in *OsC1* were found in GJ
500 weedy accessions, which were a 983-bp deletion resulting in incompleteness of exon2
501 and loss of exon 3 in accession 18WR-118 from Korea and 13-65 from Italy and a
502 3-bp insertion combined with a 2-bp deletion in accessions YCW03 and WR04-6
503 from China, which led to MYB domain loss (Fig. 4c; Supplementary Fig. 23). Brown
504 hull, which is mainly regulated by *Bh4,* is also a characteristic of feralization for some
505 weedy rice (Zhu et al., 2011). A 22-bp insertion in *Bh4* was found in both the GJ
506 weedy genome (WR04-6) and the *aus* weedy genome (PI653439), which
507 corresponded to their seed hull phenotypes (Fig. 4d; Supplementary Fig. 23). The
508 22-bp PAV was under selection during rice domestication (Zhu et al., 2011). The
509 phylogeny of *Bh4* confirmed that the 22-bp PAV in weedy rice was likely derived
510 from different groups of wild rice (Or-1 for *aus* PI653439 and Or-3 for GJ WR04-6).
511 The results also revealed that the brown or black hull color in weedy and cultivated
512 rice was due to wild introgression (Supplementary Fig. 24). Briefly, although
513 structural differences could be found between weedy and cultivated rice, almost no
514 single causative mutation could explain the convergent phenotypic change for all
515 different weedy rice groups, with the only exceptional being *Rc* for red pericarp and
516 seed dormancy; in other words, weedy rice from different groups may have

517  experienced independent evolution after the acquisition of *Rc* from wild rice or local
518  landraces of red rice.
519
520

521  **Discussion**
522

523  Two major competing hypotheses (single domestication or multiple domestication)
524  have been proposed to describe the origin of different subspecies of cultivated Asian
525  rice according to previous archaeological and genetic evidence (Molina et al. 2011;
526  Huang et al., 2012; Civáň et al., 2015; Choi et al., 2017; Choi and Purugganan, 2018;
527  Wang et al., 2018b; Carpentier et al., 2019; Zhang et al., 2021). There is no dispute
528  that rice subspecies have different wild progenitors at whole-genome level, especially
529  for XI and GJ, which means subspecies have multiple origins for most genomic
530  regions (Huang et al., 2012). However, the sources of domestication alleles are
531  debated. For well-known domestication genes, phylogeny analysis has supported
532  single domestication (e.g. *Prog1*, *sh4*, *LABA1*, *Bh4*, *OsC1*)(Huang et al., 2012; Choi
533  and Purugganan, 2018), while some wild accessions with domesticated alleles were
534  located within cultivated rice, which is also observed in this study (Fig. 3d). Gene
535  flow from cultivated to wild rice could be used to explained the sporadic distribution
536  of wild rice within cultivated sub-trees (Wang et al., 2017). However, some studies
537  interpreted these as ancestral domesticated alleles in wild rice prior to domestication
538  and highlighted possibility of independent acquisition of domestication alleles from
539  wild populations (Civáň and Brown, 2017; Civáň and Brown, 2018). To this point of
540  view, this inference seems not reasonable. If these wild alleles emerged before
541  domestication, their phylogenetic positions would not been within cultivated
542  accessions.
543

544  Wang et al. (2018b) analyzed haplotypes of nine domestication and improvement
545  genes using about 3000 rice genomes, found that many XI domesticated alleles were
546  absent in GJ and hence concluded as a support for the hypothesis of independent
547  domestication in XI, rather than GJ-to-XI introgression. This conclusion is somewhat
548  arbitrary. Firstly, the filtering and selection to variants will somewhat impact the
549  haplotype inference. Second, the method using haplotype-wise genetic distance has
550  simplified the identification criterion of introgression and is not reliable as phylogeny
551  approaches, or ABBA-BABA test. Third, the effects by genetic drift should be taken
552  into consideration. Cultivated rice genomes we have sampled and sequenced now are
553  only a subset of ancestral proto-GJ or proto-XI genetic pool, and can not fully
554  represent the diversity in domestication. GJ has suffered a dramatic genetic bottleneck
555  around three thousand years ago, while XI shows no obvious decline in effective
556  population size in the last few thousand years (Qiu et al., 2020; Gutaker et al., 2020).
557  Here, in our phylogeny analysis on putative introgression blocks, some Or-1 or Or-2
558  accessions were located between domestication clade and Or-3 clade, as relict alleles
559  of proto-GJ and genomic footprints of ancient introgression, which are absent in
560  current GJ gene pool but present in wild Or-1 or Or-2 (Fig. 3d). Therefore, the genetic

561    drift played unneglectable roles in the presence of non-GJ domesticated alleles in XI,
562    and the hypothesis of GJ-origin of domesticated alleles in XI can not be rejected just
563    based on haplotype presence or absence. Indeed, the non-GJ alleles in XI
564    domestication genes have been observed in this study. *Rc* (HDG=0.803), *An-1* (0.886)
565    and *GW5* (0.875) have high haplotype divergence between XI and GJ, but the
566    genomic regions they are located in showed robust introgression signals, supported by
567    haplotype similarity or sharing, phylogeny and ABBA-BABA tests (Supplementary
568    Tables 2 and 3). Haplotype analysis on a specific gene only sometimes will mislead
569    because of gene diversifying after domestication, and additional approaches should be
570    performed to verify inference in the meantime.
571
572    By taking advantage of 74 high-quality rice genomes, we revisited the origin of
573    different rice subspecies and groups from a pan-genome view (Fig. 5). Different from
574    previous studies using selective sweep and nucleotide diversity to infer domestication
575    genomic regions and then reconstruct their evolutionary trajectory (Huang et al., 2012;
576    Civáň et al., 2015), we investigated the core issues of the dispute, whether
577    introgression from GJ to XI has happened and whether the introgressed blocks
578    harbored domesticated alleles. Compared to whole-genome resequencing, genome
579    assembly enables us to analyze high-throughput full-length sequences of genetic
580    elements, which eliminates systematic errors caused by read mapping and sequencing
581    depth. By comparing the absolute differences among predicted protein sequences to
582    assign haplotypes, we compressed the variations within a gene, which provided us
583    with a visual landscape of haplotype similarity among taxa in each syntelog group.
584    Similar approaches using haplotypes of single genes or blocks has been adopted in
585    recent studies (Zhang et al., 2021; Wang et al., 2022). A haplotype map based on
586    SNPs from coding regions was constructed for 3,010 cultivated and 15 wild rice
587    accessions (Zhang et al., 2021). The inferred proto-ancestors of different groups
588    correlated strongly with wild rice from the same geographic regions, which was
589    considered to support a multi-domestication model of rice. However, in spite of
590    extremely insufficient number in wild rice, the whole-genome haplotype similarity
591    among populations can not indicate how domesticated alleles come from, thus the
592    data present in Zhang et al. (2020) could not lead to the conclusion of multi-origin
593    domestication model in rice. Here, by comparing the identity of predicted protein
594    sequences and the frequency of each haplotype in different groups, we assigned
595    ancestral or dominant haplotypes for each SG and determined a total of >20 Mb
596    genomic regions as putative introgression blocks, which avoided potential
597    mis-assignment by hard thresholds used in conventional clustering based on genetic
598    distance or phylogenetic relationship (Supplementary Fig. 13b). Such an ancestral
599    genomic haploblock dissection method has also been employed in tracing the origin
600    of polyploid wheat, and mosaic genomic graphs have suggested dispersed emergence
601    and protracted domestication in wheat (Wang et al., 2022). In the shared haplotype
602    blocks between GJ and XI, the majority of their phylogenetic trees, combined with the
603    relatively low divergence or genetic distances within these regions between XI and GJ,
604    supported the introgression of domestication genes from proto-GJ to Or-1 wild group.

605 Interestingly, the low-recombination centromere region of chromosome 5, which
606 suggested clear genomic affinity between subspecies, remained a relict footprint of
607 ancient introgression (Supplementary Fig. 15). Statistical ABBA-BABA tests using
608 high-depth sequencing wild genomes and larger genome sampling also confirmed the
609 reliability of putative introgression blocks (Supplementary Fig. 18).
610
611 Utilizing domestication alleles from other geo-isolated populations or species has
612 facilitated the generation of local domesticates. In wheat, dispersal domestication
613 events generated domesticated alleles or haplotypes for different genes in different
614 locations. Genomic introgression among different populations or species by human
615 activities has gathered domestication alleles of different genes together, leading to the
616 emergence and popularity of hexaploid wheat (Wang et al., 2022). Here, our results
617 strongly confirmed the previous hypothesis that genomic introgression of
618 domestication alleles from proto-GJ to wild Or-1 led to the emergence of XI (Fig. 5).
619 Despite limited sampling of the *aus* type, we inferred that the ancestral wild
620 population of *aus* was different from that of XI, although they were both clustered
621 within the Or-1 wild group. Genomic introgression from local wild rice and
622 domesticated XI rice may have led to the birth of *aus* rice. More genomes and detailed
623 analysis in the future will uncover the complex evolutionary process of the *aus* group.
624
625 De-domestication is an atavistic process in domesticates that has been studied in crops
626 (e.g., rice, wheat and sunflower) and livestock (e.g., chicken and dog) (Wu et al.,
627 2021). How de-domestication evolves in rice genomes has been investigated in recent
628 years (Song et al., 2014; Li et al., 2017; Sun et al., 2019; Qiu et al., 2020). A genomic
629 island on chromosome 6 that potentially contributes to rice de-domestication
630 syndromes (mainly red pericarp and seed dormancy or germination) has been defined
631 (Li et al., 2017; Qiu et al., 2020). Here, our analysis combining both structural
632 comparison and phylogenetic analysis highlighted the influences of wild introgression
633 on the emergence of weedy rice, despite independent introgression events for GJ and
634 XI (Fig. 5). For the brown hull of some weedy accessions, the causative structural
635 mutation was indicated to be derived from the corresponding wild groups Or-1 and
636 Or-3 for XI and GJ, respectively. Although pairwise whole-genome comparison
637 identified thousands of structural variants between weedy and cultivated rice,
638 structural convergence in different weedy-cultivated lineages was seldom observed,
639 which implied the recurrent independent emergence of weedy traits. The mechanism
640 underlying high shattering in weedy rice was still not well resolved from the view of
641 structural variation in known domestication genes (except *sh4* and *SHAT1* for
642 shattering). Given that our known shattering-related genes are almost all transcription
643 factors, variations in other regulatory elements or even epigenetic factors may lead to
644 high shattering in weedy rice. Overall, genomic introgression plays an indispensable
645 role throughout the entire evolutionary trajectory of rice, from initial domestication,
646 improvement, modern breeding and feralization.
647
648

**Materials and Methods**

**Genome sequencing and assembly**

To capture the genetic diversity from all ecotypes of rice, we collected 11 accessions of weedy rice from seven countries based on their phylogeny with cultivated rice and geographic positions; we also included one XI cultivar NJ11, or Nanjing 11, which is presumed to be the direct ancestor of weedy rice in the Yangtze River Basin (Qiu et al., 2020). Genomic DNA samples were extracted from young leaves of the 12 rice accessions, and their genomes were sequenced by PacBio HiFi mode according to the instructions from the manufacturer. For each accession, the sequencing depth of HiFi subreads ranged from 21.6× for accession 13-65 to 45.0× for accession YZ-2, with an average of 32.4×. Following the standard protocol, Hi-C libraries of four accessions (NJ11, CX20, YCW03 and 18XHB-83) were constructed using fresh young leaves digested with the 4-cutter restriction enzyme MboI. Hi-C libraries were sequenced on an Illumina HiSeq 4000 platform with 2×150-bp paired reads.

The genomes were first assembled using hifiasm (v0.15.1-r334, default parameters) (Cheng et al., 2020). For each accession, HiFi subreads were mapped against the corresponding assembly using minimap2 (Li, 2018), and Purge_dups was applied to purge duplicates and remove redundant sequences according to the mapping depth (Guan et al., 2020). We further used Racon (v1.4.0) to polish the assemblies with HiFi subreads for three rounds under default parameters (Vaser et al., 2017). Contigs less than 10 kb were removed from the final version. For each accession, contigs were anchored into pseudochromosomes by using a reference-guiding approach RaGOO (Alonge et al., 2019). By aligning contigs against the Nipponbare assembly using minimap2, the contigs were ordered and oriented along 12 chromosomes with no further chimeric splitting.

Previously released rice assemblies based on third-generation sequencing platforms were collected, including PacBio (Du et al., 2017; Sun et al., 2019; Wang et al., 2019; Ma et al., 2020; Xie et al., 2020; Zhou et al., 2020; Qin et al., 2021; Song et al., 2021) and Nanopore sequencing (Choi et al., 2020; Read et al., 2020; Shang et al., 2022; Zhang et al., 2022). Before adopting assemblies in the construction of the rice pangenome, we first systematically evaluated the assembly qualities and ruled out assemblies that did not meet our criteria.

**Quality assessment of rice genome assemblies**

We first assessed the quality of our 12 newly assembled rice genomes. Synteny against the reference assembly Nipponbare and gapless assembly MH63 (Song et al., 2021) confirmed their high completeness. The paired Hi-C reads of four accessions (YCW03, 18XHB-83, CX20 and NJ11) were cleaned using NGSQC toolkit (Patel and Jain, 2012) and then mapped to the corresponding assembly using Bowtie2 (v2.3.5.1) (Langmead and Salzberg, 2012). After retaining high-quality and validated paired reads (mapping quality ≥ 30, edit distance ≤ 5, number of mismatches in the alignment ≤ 3, number of gap opens ≤ 2 and number of gap extensions ≤ 2), chromosome interaction maps were plotted by AllHiC_plot (Zhang et al., 2020), and

693    they revealed high accuracy in contig ordering and orientation (Supplementary Fig.
694    1).

695    Genome assemblies of *Oryza sativa* and *Oryza rufipogon* based on third-generation
696    sequencing (through May 2022) were collected (Supplementary Table 1). DXCWR
697    (Ma et al., 2020) was excluded considering its low contig N50 of less than 200 kb.
698    IRGC109232 (Zhao et al., 2018) was removed due to the abnormal size of the
699    assembly obtained from the public database. Eleven assemblies in Zhang et al. (2022)
700    were randomly selected from 75 newly generated genomes and used in subsequent
701    quality assessment. Five indices were applied to evaluate the genome quality of all
702    rice assemblies (Supplementary Fig. 2a). Assembly continuity was evaluated by
703    contig N50 and the long-terminal repeat assembly index (LAI), which was revealed
704    by the assembly completeness of long-terminal-repeat (LTR) retrotransposons (Ou et
705    al., 2018). The LTR elements of each assembly were identified by RepeatMasker and
706    RepeatModeler (http://repeatmasker.org/). BUSCO (v4.1.2) metrics were calculated to
707    evaluate the completeness by using dataset poales_odb10 containing 4896 genes
708    (Simao et al., 2015). As expected, assemblies based on long-read sequencing
709    performed well on the above quality indices (Supplementary Fig. 2a). Hence, in
710    addition to the above evaluations at the whole-genome level, base-resolution accuracy
711    and completeness were measured by the consensus quality value (QV) and the
712    number of homozygous variants called by self-short-read mapping. Reference-free
713    QVs were calculated by Merqury (Rhie et al., 2020) and yak
714    (https://github.com/lh3/yak) by comparing *k*-mers derived from unassembled,
715    high-accuracy sequencing reads to a genome assembly. Homozygous variants (SNPs
716    and InDels) called with short reads by self-mapping are regarded as potential
717    assembly errors. Raw short-read data were first cleaned by NGSQC-toolkit and
718    mapped against the corresponding assembly by Bowtie2. Variants were detected using
719    GATK (v3.7, default parameters) (McKenna et al., 2010) and annotated by SnpEff
720    (v3.6) to profile their potential effects on the prediction of amino acid sequences and
721    further gene functions (Cingolani et al., 2012). Variants with high effects (including
722    stop loss and gain, start loss, frame-shift variant) and moderate effects (including
723    in-frame insertion/deletion, missense variant) will directly impact the reliability of
724    gene haplotypes and downstream haplotype analysis. The assembly quality of wild
725    accessions from Shang et al. (2022) and one weedy accession (YCW03) were not
726    assessed at the base level due to the unavailability of NGS data. Assemblies with low
727    base-level quality were removed for cultivated rice, mainly including Nanopore
728    sequencing-based genomes, except three accessions (DomSufid, Basmati334 and JHU)
729    from aromatic and tropical groups of GJ (*O. sativa* ssp. *japonica*), which were kept to
730    balance the sampling of each group. For wild genomes, only two accessions, W2014
731    from Ma et al. (2020) and IRGC106162 from Xie et al. (2020), sequenced using the
732    PacBio platform were available. Thus, nine assemblies in different wild groups from
733    Shang et al. (2022) sequenced by the Nanopore platform were adopted for analysis.
734

735    **Phylogenetic relationship of rice assemblies**
736    Together with the 12 new assemblies in this study, a total of 75 rice assemblies,

737    including 11 wild accessions (*O. rufipogon*), 12 weedy accessions (*O. sativa* ssp.

738    *spontanea*), 51 cultivated accessions (*O. sativa*) and an African cultivated rice

739    outgroup (*Oryza glaberrima*) accession CG14, were used in this study. We first

740    confirmed their phylogenetic relationship and assigned them to taxonomic groups

741    using whole-genome SNPs. The assemblies were aligned against the reference

742    assembly Nipponbare using nucmer implemented in the MUMmer package (v4.0.0)

743    (Marçais et al., 2018), and called SNPs were used to build their phylogeny by

744    FastTreeMP with 1000 bootstrap replicates (Price et al., 2009). According to their

745    phylogeny and prior knowledge of the 74 *Oryza sativa* and *Oryza rufipogon*

746    accessions, each accession was assigned to groups. Briefly, the subspecies GJ (*O.*

747    *sativa* ssp. *japonica*, $n = 23$ in total) includes four *aromatic* (aro), four tropical (trp),

748    two subtropical (subtrp) and 13 temperate (tmp) accessions; subspecies XI (*O. sativa*

749    ssp. *indica*, $n = 40$) includes four *aus*, three XI2, seven XI3, ten XI1A and 16 XI1B

750    accessions. The wild population (*O. rufipogon*, $n = 11$) includes two accessions from

751    Or-3, four from Or-2, three from Or-1, and two from Or-4. Genotype data for

752    approximately seven thousand rice accessions used in the principal component

753    analysis (PCA) were adopted from a previous study (Wu et al., 2022a). PCA and IBD

754    (Identity-by-descent) calculation were performed using Plink (v1.9) with a pruned

755    subset of SNPs based on linkage disequilibrium (10 SNPs in each 50-kb sliding

756    window with pairwise Pearson's correlation efficient $r^2$ less than 0.5) (Chang et al.,

757    2015).

758

**759    Genome annotation**

760    In total, 75 rice assemblies (including African rice CG14 as an outgroup) were

761    annotated in a unified pipeline. First, transposon elements (TEs) were identified by

762    using the Extensive *de novo* TE Annotator (EDTA) approach

763    (https://github.com/oushujun/EDTA) (Su et al., 2021). For each accession, gene

764    models were predicted on the repeat-masked genome using an approach integrating

765    *ab initio* predictions and homology-based prediction. For *ab initio* prediction,

766    Augustus (Stanke et al., 2006) and Fgenesh (Salamov and Solovyev, 2000) were

767    performed with default parameters. For homology-based prediction, previously

768    predicted protein sequences of the Nipponbare (IRGSP v1.0), gapless MH63RS3 and

769    ZS97RS3 assemblies (Song et al., 2021) and over 30 other well-annotated assemblies

770    (Qin et al., 2021) were used to search putative protein-coding gene models with

771    GMAP for each accession (Wu and Watanabe, 2005). The predictions were integrated

772    into non-redundant consensus gene models using EVidenceModeler (v1.1.1) (Haas et

773    al., 2008). Short gene models (less than 50 amino acids) and gene models with

774    homology to sequences in Repbase ($e$-value $\leq$ 1e−5, identity $\geq$ 30%, coverage $\geq$ 25%)

775    were further removed from the final annotation. The protein domains of all predicted

776    coding gene models were inferred using InterProScan (v5.24-63.0) (Zdobnov and

777    Apweiler, 2001).

778

**779    Pan-genome construction using MCL**

780    A Markov Clustering (MCL) approach OrthoFinder (v2.4.1) was applied to cluster all

781    the predicted gene models of 74 rice genomes (African rice accession CG14 was

782    excluded) with default parameters (diamond all-versus-all *e*-value < 1e-5 and inflation

783    parameter = 1.5) (Emms and Kelly, 2019). Finally, over 3.10 million predicted genes

784    were clustered into 67,080 orthogroups (OGs). Increasing the inflation parameter can

785    be used to achieve higher precision at the cost of lower recall and have a larger size of

786    OGs. Conversely, a smaller value of the inflation parameter could achieve higher

787    recall at the cost of lower precision and produce smaller sized OGs, which easily

788    clusters paralogs together (Emms and Kelly, 2019). Thus, three additional inflation

789    parameters (I = 1.8, 2.0 and 2.5) were set to profile the effects on clustering

790    (Supplementary Fig. 9). Clustered gene families were categorized into core (present in

791    all genomes, *n* = 74), soft-core (present in at least 90% of genomes, *n* = 67 to 73),

792    dispensable (present in more than one but less than 90% of genomes, *n* = 2 to 66) and

793    private (only present in one genome) on the basis of the number of rice accessions in

794    which they were identified.

795

796    **Pan-genome construction using synteny**

797    Given that the MCL approach does not efficiently and accurately distinguish paralogs

798    from orthologs with high sequence similarity, we utilized the availability of genomic

799    coordinates of gene models to build the synteny-based pan-genome. Pairwise all-to-all

800    alignments using protein sequences were performed for all 74 genomes. Aligner

801    Diamond runs much faster for large protein sequence data than BLASTP (Buchfink et

802    al., 2021). We compared the recall performance between BLASTP and Diamond in

803    detecting syntelogs between genomes (Supplementary Fig. 6). Diamond was run for

804    each pair under three modes (default, sensitive and ultra-sensitive). The alignments

805    were filtered to keep only the best hits, and DAGchainer (Haas et al., 2004) was used

806    to detect syntenic genomic regions and syntelogs (parameters: -Z 12 -D 200000 -g 1

807    -A 5). No differences were observed among the three Diamond modes in the number

808    of detected syntelogs (Supplementary Fig. 6). BLASTP was performed under default

809    parameters, and syntelogs were identified using the same pipeline. The BLASTP

810    approach searched more syntelogs, and thus, the syntelogs were used in downstream

811    pan-genome construction (Supplementary Fig. 6). Instead of a reciprocal best hit

812    search, the roles of reference and query in a pairwise BLASTP search sometimes

813    result in differences due to individual-specific tandem duplicates, and such

814    individual-specific tandem duplicates were considered as a single gene in the final

815    pan-genome. Pairwise syntelog information was merged together with Nipponbare as

816    an initial framework one by one using SynPan

817    (https://github.com/dongyawu/PangenomeEvolution). If a gene from an additional

818    genome was syntenic to a previously merged pan-genome, this gene was assigned to

819    an existing SG. If a gene from an additional genome was not syntenic to any gene in

820    the merged iterative pan-genome, a new SG was created. In total, 74 genomes were

821    merged together as a synteny-based pan-genome, including 175,528 SGs. The SGs

822    were further categorized as core, soft-core, dispensable and private following the

823    criteria used in OG categorization.

824

**Construction of the rice NLRome**

To capture the diversity of NLRs in rice, we integrated multiple software predictions and gene synteny in rice genomes to obtain a comprehensive and complete rice NLRome. The NB-ARC domain was first predicted using hmmsearch (HMMER v3.1b2) against the Pfam database (v30.0) with a threshold *e*-value less than 1e-5. The LRR domains were predicted with NLR-parser (v3.0) (Steuernagel et al., 2015) by searching for motifs 9, 11 and 19; the coil domains were predicted by searching for motifs 16 and 17; and the TIR domains were predicted by searching for motifs 13, 15 and 18. All putative NLR types from genome-wide protein sequences were also determined using RGAugury (Li et al., 2016). After the NLR genes for each genome were identified by domain prediction, the NLR genes were mapped back to the SG-based pangenome to involve NLR syntelogs lacking canonical domains and find a more comprehensive and extensive NLR inventory. NLR genes were defined to have at least one NB-ARC, TIR, or CCR (RPW8) canonical domain. LRR or CC motifs alone were not considered sufficient for NLR identification. Finally, NLRs in rice genomes were identified and categorized as CNLs (containing CC, NB-ARC, and LRR domains), CNs (containing CC and NB-ARC domains) and NLs (containing only canonical NB-ARC domains). Noncanonical architectures of some NLRs have additional integrated domains (IDs), while canonical architectures contain only NB-ARC (Pfam accession PF00931), TIR (PF01582), RPW8 (PF05659), or LRR (PF00560, PF07725, PF13306, PF13855) domains or CC motifs. For ID identification, we used the genome protein sequences as the input for InterProScan, and annotations were processed with in-house scripts to obtain ID information. There are currently several known cases in which two NLR genes are required to affect the resistance function, with one protein in the pair acting in effector recognition and the other acting in signaling activation. Since all known functional pairs are present in head-to-head arrangement in the rice genome, we identified head-to-head NLR pairs by searching for NLR genes near each other and no more than 10 kb away. Multiple sequence alignment was performed using coding sequences by MAFFT (v7.490) (Katoh & Standley, 2013), and the nucleotide diversity and Tajima's *D* value for all members in each SG were calculated using DnaSP (v6) (Rozas et al., 2017).

**Haplotype diversity and ancestral haplotype assignment**

We used the average pairwise haplotype difference to measure the haplotype diversity in one population and the divergence between two populations (Supplementary Fig. 13a). The haplotype number (defined as N100) for each SG was determined by counting the unique sequences of all predicted proteins. To exclude rare haplotypes, haplotype N90 was defined as the least haplotype number that needs to be included for covering 90% of sequences in each SG. To understand the haplotype-aware origin of genes from different rice groups, we inferred the ancestral haplotype composition for each SG by determining the dominant haplotype in each rice group and assigning ancestral haplotype IDs for all genes in a group priority-based referring strategy (Supplementary Fig. 13b). Group information is prior based on whole-genome phylogeny or population structure. We labeled five haplotype IDs (hapI to hapV) to

869 represent the dominant haplotypes of groups tmp, XI1A, XI1B, *aus* and XI3 in order,
870 presented by red, blue, orange, yellow, and green in Fig. 4a and Supplementary Fig.
871 15, respectively. HapI was first defined as the most dominant sequence in group tmp.
872 If the dominant haplotype in the current group was defined by a former group, the
873 haplotype ID was skipped and set as missing. For example, if the dominant haplotype
874 in XI1A was the same as that in tmp (hapI), hapII was then not defined. The
875 frequency of a dominant haplotype within a group should be at least three. All other
876 rare haplotypes were compressed as hapR in gray to simplify the ancestral haplotype
877 graphs. Gene absence is represented by white blocks. Different group priority orders
878 have no influences on the calculation of haplotype diversity and divergence but only
879 change the ancestral haplotype graphs (Supplementary Fig. 13c).
880
881 **Inter-subspecies introgression blocks**
882 As observed from the ancestral haplotype landscape, haplotypes in some large
883 genomic regions were shared between XI and GJ. We used inter-population haplotype
884 divergence (HDG) to quantify haplotype sharing by calculating the average
885 differences in haplotypes from the two populations (Supplementary Fig. 13a). We
886 merged adjacent SGs whose divergence values between tmp and XI (excluding *aus*)
887 were less than 0.5 into lowly divergent blocks between subspecies as candidate
888 introgression blocks. At least 10 SGs were required within a single introgression
889 block. To examine the significance of the nonrandom clustered distribution of these
890 lowly divergent SGs in blocks, we randomly sampled the same number of SGs as
891 lowly divergent SGs observed on each chromosome and calculated the density of
892 sampled SGs in sliding windows of every 10 SGs. A total of 100 thousand random
893 samplings were replicated, and the $P$ values were determined by counting the
894 sampling times where the density of sampled SGs was higher than that observed for
895 each window. $P = 0.01$ was empirically set as a cutoff value. Assuming that the low
896 divergence in the detected blocks was caused by inter-subspecies introgression after
897 their divergence, the divergence time between GJ and XI should be younger in the
898 identified blocks than in their neighboring regions. We used the synonymous
899 substitution rate ($K_s$) to measure the relative divergence time between GJ and XI, free
900 from selection. Significantly lower $K_s$ values were observed in the majority of 18
901 large blocks (>300 kb) than in their flanking regions (Supplementary Fig. 16).
902
903 **Phylogeny of introgression blocks**
904 To investigate the origin and gene flow of the 73 candidate introgression blocks, we
905 utilized recently released genomic sequences of 184 wild accessions with high
906 sequencing depth (8× on average, much higher than <2× on average in a previously
907 used wild population)(Zheng et al., 2022; Huang et al., 2012). Raw sequencing reads
908 were first cleaned using NGSQC toolkit and mapped against the reference assembly
909 Nipponbare (IRGSP v1.0) using Bowtie2. Assemblies of two wild rice accessions
910 (W1943 and DWCWR) from group Or-3 were added. Combing the 77 assemblies and
911 184 wild genomes, high-quality SNPs were called following a previous pipeline (Qiu
912 et al., 2020). The population structure was first surveyed by PCA using Plink (v1.9)

913  (Chang et al., 2015) and the phylogeny was produced by FastTreeMP with 1000
914  bootstrap replicates based on 6.85 million high-quality SNPs (minor allele frequency
915  of 0.02 and maximum missing rate of 0.1). Four wild groups were identified: Or-1,
916  Or-2, Or-3 and Or-4. The SNPs in each introgression block were extracted and used to
917  build the phylogeny using IQ-TREE (v1.6.12) with the best substitution model
918  TIM2e+R2 determined by ModelFinder implemented in IQ-TREE (Nguyen et al.,
919  2015) and FastTreeMP with 1000 bootstrap replications, where the African cultivated
920  rice accession CG14 was set as the outgroup. To avoid over-interpretation, the
921  phylogeny in which the wild accessions were not obviously and empirically clustered
922  into four groups, as defined by the whole-genome SNPs, was discarded.
923
924  **ABBA-BABA test**
925  The availability of population-level whole-genome high-depth sequencing data of
926  wild rice from four groups (Or-1, Or-2, Or-3 and Or-4), enables us to perform
927  comprehensive statistical $D$ tests, which is widely used and robust in gene flow
928  detection (Green et al., 2010; Wu et al., 2022b). We randomly selected genomes of
929  XI1A rice accessions ($n = 100$), XI1B ($n = 100$), XI2 ($n = 80$), XI3 ($n = 100$) and *aus*
930  ($n = 60$) from 3K RGP (Wang et al., 2018). We employed $f_d$ statistic to indicate the
931  introgression from GJ to XI in sliding genomic windows (Martin et al., 2015). Under
932  a given quartet topology ((P1, P2), P3, O), positive $f_d$ statistic values indicate the
933  introgression from P3 to P2, zero represents no introgression, and negative $f_d$ statistic
934  values have no biological meaning and thus are converted to zero. We estimated the $f_d$
935  statistic values under topology ((Or-1, X), tmp, Or-4), where X is Or-2, XI1A, XI1B,
936  XI2, XI3 and *aus* in topology T1 to T6, respectively (Supplementary Fig. 18). T1 is
937  set as a background control in detecting introgression. To eliminate the influence of
938  modern breeding where inter-subspecies hybridization is frequently performed on
939  ancient introgression inference, we used genomes of only landraces in each group to
940  repeat the $f_d$ calculation, where tmp, XI1A, XI2, XI3, and *aus* included 47, 24, 21, 52
941  and 32 landrace accessions. Generally, no obvious differences are observed between
942  introgression block determination using all and landrace accessions only
943  (Supplementary Fig. 19). Python scripts are available at
944  https://github.com/simonhmartin/genomics_general. Parameters are set as "window
945  size: 20 kb, step size: 2 kb, minimum good sites per window: 50, and minimum
946  proportion of samples genotyped per site: 0.4". The final putative introgression
947  regions were determined by integrating evidences form haplotype divergence,
948  phylogeny and ABBA-BABA tests. The functional enrichment analysis was
949  performed using ShinyGO (v0.77) (Ge et al., 2020).
950
951  **Structural variations in de-domestication**
952  To gain a more detailed understanding of structural variations in de-domestication, we
953  compared the genome assemblies of weedy and cultivated rice. It has been found that
954  the sensitivity of detecting deletions is higher than that of insertions, and we adopted a
955  pairwise genome alignment strategy as mentioned previously in Jayakodi et al. (2020).
956  Each pair contains a weedy accession and a cultivar accession, which compose a

monophyly in the phylogenetic tree, and the genome assembly of cultivated rice was considered a query or reference genome. Nucmer in the MUMmer package was used to obtain the results of these two alignments (Marçais et al., 2018), and then PAVs (presence-and-absence variants, including insertions and deletions) were called using Assemblytics (v.1.2.1) (Nattestad & Schatz, 2016). The structural PAVs were classified as InDels (<50 bp) and SVs ($\geq$ 50 bp). Only deletions were kept in both alignments and converted into PAVs according to the reference genome in each pair. Genes located in or intersected with each PAV were obtained, as well as corresponding gene annotations. The variations in domestication genes and their flanking regions (15 kb for *qSH1* and 2 kb for other genes) were manually investigated, and the causative mutations during the domestication process were checked. To validate the reliability of structural variations in domestication or improvement genes (e.g., *OsC1* and *Bh4*), HiFi subreads of weedy rice were mapped against themselves and corresponding cultivated assemblies to check the local alignments using IGV (Thorvaldsdóttir et al., 2013). To infer the source of the causative mutation in *Bh4* in weedy accessions, the *Bh4* phylogeny based on SNPs was analyzed.

The 73 rice genome assemblies were aligned against the Nipponbare reference genome using the nucmer program implemented in the MUMmer package with the default parameter, and only the best position of each query on the reference was preserved. The alignments from 6.0 to 6.5 Mb on chromosome 7 in the Nipponbare assembly were extracted for visualization by synteny plots and comparison among wild, cultivated and weedy assemblies. To validate the candidate introgression event inferred from the synteny plots, SNPs around the *Rc* gene (including its flanking 2-kb regions) were extracted and used in phylogeny construction by FastTreeMP under the GTR+CAT model with 1000-times bootstrapping.

**Data and code availability**
All the PacBio HiFi subreads for 12 rice accessions and Hi-C data for four rice accessions generated in this study have been deposited in NGDC (https://ngdc.cncb.ac.cn/) under the accession code PRJCA012143. The newly generated assemblies for 12 accessions and the annotations (including GFF, CDS sequences and predicted protein sequences) for all 74 rice accessions can be found under project accession PRJCA012309 in NGDC. The raw resequencing data of previously published wild accessions can be downloaded from NCBI under accession number PRJNA657701. The VCF file of SNPs from all 75 assemblies and an additional 184 wild rice accessions with high sequencing depth is available at Zenodo (10.5281/zenodo.7196576). The gene re-annotations, NLR annotations and Pfam annotations for all 75 rice accessions are deposited at Zenodo (10.5281/zenodo.7248110). The in-house scripts used in this study have been deposited in GitHub (https://github.com/dongyawu/PangenomeEvolution).

## Reference

Alonge, M., Soyk, S., Ramakrishnan, S., Wang, X., Goodwin, S., Sedlazeck, F.J., Lippman, Z.B., and Schatz, M.C. (2019). RaGOO: fast and accurate reference-guided scaffolding of draft genomes. Genome Biol 20.

Buchfink, B., Reuter, K., and Drost, H. (2021). Sensitive protein alignments at tree-of-life scale using diamond. Nat Methods 18:366-368.

Carpentier, M.C., Manfroi, E., Wei, FJ., Wu, H.P., Lasserre E., Llauro C., Debladis E., Akakpo R., Hsing Y.I. and Panaud O. (2019) Retrotranspositional landscape of Asian rice revealed by 3000 genomes. Nat Commun 10, 24.

Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015) Second-generation PLINK: rising to the challenge of larger and richer datasets. Gigascience 4:7.

Chen, E., Huang, X., Tian, Z., Wing, R.A., and Han, B. (2019). The genomics of *Oryza* species provides insights into rice domestication and heterosis. Annu Rev Plant Biol 70:639-665.

Cheng, H., Concepcion, G.T., Feng, X., Zhang, H., and Li, H. (2021). Haplotype-resolved *de novo* assembly using phased assembly graphs with hifiasm. Nat Methods 18:170-175.

Choi, J.Y., Lye, Z.N., Groen, S.C., Dai, X., Rughani, P., Zaaijer, S., Harrington, E.D., Juul, S., and Purugganan, M.D. (2020). Nanopore sequencing-based genome assembly and evolutionary genomics of circum-basmati rice. Genome Biol 21.

Choi, J.Y., Platts, A.E., Fuller, D.Q., Hsing, Y., Wing, R.A., and Purugganan, M.D. (2017). The rice paradox: multiple origins but single domestication in Asian rice. Mol Biol Evol 34(4):969-979.

Choi J.Y., and Purugganan M.D. (2018) Multiple origin but single domestication led to *Oryza sativa*. G3 8(3):797-803.

Cingolani, P., Platts, A., Wang, L.L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X., and Ruden, D.M. (2014). A program for annotating and predicting the effects of single nucleotide polymorphisms, snpeff. Fly 6:80-92.

Civáň, P., Craig, H., Cox, C.J., and Brown, T.A. (2015). Three geographically separate domestications of Asian rice. Nat Plants 1:15164.

Civáň, P., and Brown, T.A. (2017). Origin of rice (*Oryza sativa* L.) domestication genes. Genet Resour Crop Evol. 64(6):1125-1132.

Civáň, P., and Brown, T.A. (2018). Misconceptions regarding the role of introgression in the origin of *Oryza sativa* subsp. *indica*. Front Plant Sci. 9:1750.

Du, H., Yu, Y., Ma, Y., Gao, Q., Cao, Y., Chen, Z., Ma, B., Qi, M., Li, Y., and Zhao, X., et al. (2017). Sequencing and *de novo* assembly of a near complete *indica* rice genome. Nat Commun 8.

Edelman, N.B., Frandsen, P.B., Miyagi, M., Clavijo, B., Davey, J., Dikow, R.B., Garcia-Accinelli, G., Van Belleghem, S.M., Patterson, N., and Neafsey, D.E., et al. (2019). Genomic architecture and introgression shape a butterfly radiation. Science 366:594-599.

Emms, D.M., and Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for

1045        comparative genomics. Genome Biol 20(1):238.

1046 Fujino, K., Sekiguchi, H., Matsuda, Y., Sugimoto, K., Ono, K., and Yano, M. (2008).
1047        Molecular identification of a major quantitative trait locus, *qltg3-1*, controlling
1048        low-temperature germinability in rice. Proc Natl Acad Sci U S A
1049        105:12623-12628.

1050 Ge, S.X., Jung. D., and Yao, R. (2020). ShinyGO: a graphical gene-set enrichment
1051        tool for animals and plants. Bioinformatics 36(8):2628-2629.

1052 Green, R.E., Krause, J., Briggs, A.W., Maricic, T., Stenzel, U., Kircher, M., Patterson,
1053        N., Li, H., Zhai, W., Fritz, M.H., et al. (2010). A draft sequence of the Neandertal
1054        genome. Science 328(5979):710-722.

1055 Gross, B.L., and Zhao, Z. (2014). Archaeological and genetic insights into the origins
1056        of domesticated rice. Proc Natl Acad Sci U S A 111:6190-6197.

1057 Guan, D., McCarthy, S.A., Wood, J., Howe, K., Wang, Y., and Durbin, R. (2020).
1058        Identifying and removing haplotypic duplication in primary genome assemblies.
1059        Bioinformatics 36:2896-2898.

1060 Gutaker, R.M., Groen, S.C., Bellis, E.S., Choi, J.Y., Pires, I.S., Bocinsky, R.K.,
1061        Slayton, E.R., Wilkins, O., Castillo, C.C., Negrão, S., et al. (2020). Genomic
1062        history and ecology of the geographic spread of rice. Nat Plants 6(5):492-502.

1063 Haas, B.J., Delcher, A.L., Wortman, J.R., and Salzberg, S.L. (2004). DAGchainer: a
1064        tool for mining segmental genome duplications and synteny. Bioinformatics
1065        20(18):3643-3646.

1066 Haas, B.J., Salzberg, S.L., Zhu, W., Pertea, M., Allen, J.E., Orvis, J., White, O., Buell,
1067        C.R., and Wortman, J.R. (2008). Automated eukaryotic gene structure annotation
1068        using EVidenceModeler and the program to assemble spliced alignments.
1069        Genome Biol 9(1):R7.

1070 He, L., Li, M., Qiu, Z., Chen, D., Zhang, G., Wang, X., Chen, G., Hu, J., Gao, Z., and
1071        Dong, G., et al. (2020). Primary leaf-type ferredoxin 1 participates in
1072        photosynthetic electron transport and carbon assimilation in rice. The Plant
1073        Journal 104:44-58.

1074 He S., Tan L., Hu Z., Chen G., Wang G., and Hu T. (2012). Molecular characterization
1075        and functional analysis by heterologous expression in *E. coli* under diverse
1076        abiotic stresses for *OsLEA5*, the atypical hydrophobic LEA protein from *Oryza*
1077        *sativa* L. Mol Genet Genomics 287(1):39-54.

1078 Hua, L., Wang, D.R., Tan, L., Fu, Y., Liu, F., Xiao, L., Zhu, Z., Fu, Q., Sun, X., and
1079        Gu, P., et al. (2015). *LABA1*, a domestication gene associated with long, barbed
1080        awns in wild rice. The Plant Cell 27:1875-1888.

1081 Huang, X., Kurata, N., Wei, X., Wang, Z., Wang, A., Zhao, Q., Zhao, Y., Liu, K., Lu,
1082        H., and Li, W., et al. (2012). A map of rice genome variation reveals the origin of
1083        cultivated rice. Nature 490:497-501.

1084 Izumi, M., Hidema, J., Wada, S., Kondo, E., Kurusu, T., Kuchitsu, K., Makino, A.,
1085        and Ishida, H. (2015). Establishment of monitoring methods for autophagy in
1086        rice reveals autophagic recycling of chloroplasts and root plastids during energy
1087        limitation. Plant Physiol 167:1307-1320.

1088 Jayakodi, M., Padmarasu, S., Haberer, G., Bonthala, V.S., Gundlach, H., Monat, C.,

1089  Lux, T., Kamal, N., Lang, D., and Himmelbach, A., et al. (2020) The barley
1090  pan-genome reveals the hidden legacy of mutation breeding. Nature
1091  588(7837):284-289.

1092  Jiao, Y., Wang, Y., Xue, D., Wang, J., Yan, M., Liu, G., Dong, G., Zeng, D., Lu, Z.,
1093  and Zhu, X., et al. (2010). Regulation of *OsSPL14* by OsmiR156 defines ideal
1094  plant architecture in rice. Nat Genet 42:541-544.

1095  Jin, J., Huang, W., Gao, J., Yang, J., Shi, M., Zhu, M., Luo, D., and Lin, H. (2008).
1096  Genetic control of rice plant architecture under domestication. Nat Genet
1097  40:1365-1369.

1098  Jones, P., Binns, D., Chang, H.Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H.,
1099  Maslen, J., Mitchell, A., and Nuka, G., et al. (2014). Interproscan 5:
1100  genome-scale protein function classification. Bioinformatics 30:1236-1240.

1101  Kato, Y., Konishi, M., Shigyo, M., Yoneyama, T., and Yanagisawa, S. (2010).
1102  Characterization of plant eukaryotic translation initiation factor 6 (*eif6*) genes:
1103  the essential role in embryogenesis and their differential expression in
1104  *Arabidopsis* and rice. Biochem Bioph Res Co 397:673-678.

1105  Katoh, K., and Standley, D.M. (2013) MAFFT multiple sequence alignment software
1106  version 7: improvements in performance and usability. Mol Biol Evol
1107  30(4):772-780.

1108  Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with bowtie 2.
1109  Nat Methods 9:357-359.

1110  Li, C., Zhou, A., and Sang, T. (2006). Rice domestication by reducing shattering.
1111  Science 311:1932-1936.

1112  Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences.
1113  Bioinformatics 34:3094-3100.

1114  Li, L., Li, Y., Jia, Y., Caicedo, A.L., and Olsen, K.M. (2017). Signatures of adaptation
1115  in the weedy rice genome. Nat Genet 49:811-814.

1116  Li, P., Quan, X., Jia, G., Xiao, J., Cloutier, S., and You, F.M. (2016). RGAugury: a
1117  pipeline for genome-wide prediction of resistance gene analogs (RGAs) in plants.
1118  BMC Genomics 17(1):852

1119  Liu, J., Chen, J., Zheng, X., Wu, F., Lin, Q., Heng, Y., Tian, P., Cheng, Z., Yu, X., and
1120  Zhou, K., et al. (2017). *GW5* acts in the brassinosteroid signalling pathway to
1121  regulate grain width and weight in rice. Nat Plants 3.

1122  Long, Y., Zhao, L., Niu, B., Su, J., Wu, H., Chen, Y., Zhang, Q., Guo, J., Zhuang, C.,
1123  and Mei, M., et al. (2008). Hybrid male sterility in rice controlled by interaction
1124  between divergent alleles of two adjacent genes. Proc Natl Acad Sci U S A
1125  105:18871-18876.

1126  Luo, J., Liu, H., Zhou, T., Gu, B., Huang, X., Shangguan, Y., Zhu, J., Li, Y., Zhao, Y.,
1127  and Wang, Y., et al. (2013). *An-1* encodes a basic helix-loop-helix protein that
1128  regulates awn development, grain size, and grain number in rice. The Plant Cell
1129  25:3360-3376.

1130  Ma, X., Fan, J., Wu, Y., Zhao, S., Zheng, X., Sun, C., and Tan, L. (2020).
1131  Whole-genome *de novo* assemblies reveal extensive structural variations and
1132  dynamic organelle-to-nucleus DNA transfers in African and Asian rice. The Plant

1133    Journal 104:596-612.

1134    Marçais, G., Delcher, A.L., Phillippy, A.M., Coston, R., Salzberg, S.L., and Zimin, A.
1135        (2018). Mummer4: a fast and versatile genome alignment system. PloS Comput
1136        Biol 14:e1005944.

1137    Martin, S.H., Davey, J.W., and Jiggins, C.D. (2015). Evaluating the use of
1138        ABBA-BABA statistics to locate introgressed loci. Mol Biol Evol. 32(1):244-57.

1139    McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A.,
1140        Garimella, K., Altshuler, D., Gabriel, S., and Daly, M., et al. (2010). The genome
1141        analysis toolkit: a mapreduce framework for analyzing next-generation DNA
1142        sequencing data. Genome Res 20:1297-1303.

1143    Molina, J., Sikora, M., Garud, N., Flowers, J.M., Rubinstein, S., Reynolds, A., Huang,
1144        P., Jackson, S., Schaal, B.A., and Bustamante, C.D., et al. (2011). Molecular
1145        evidence for a single evolutionary origin of domesticated rice. Proc Natl Acad
1146        Sci U S A 108:8351-8356.

1147    Nattestad, M., and Schatz, M.C. (2016) Assemblytics: a web analytics tool for the
1148        detection of variants from an assembly. Bioinformatics. 32(19):3021-3023.

1149    Nguyen, L.T., Schmidt, H.A., von Haeseler, A., and Minh, B.Q. (2015). IQ-TREE: a
1150        fast and effective stochastic algorithm for estimating maximum-likelihood
1151        phylogenies. Mol Biol Evol 32(1):268-274.

1152    Ou, S., and Jiang, N. (2018). LTR_retriever: a highly accurate and sensitive program
1153        for identification of long terminal repeat retrotransposons. Plant Physiol
1154        176:1410-1422.

1155    Patel, R.K., and Jain, M. (2012). NGS QC toolkit: a toolkit for quality control of next
1156        generation sequencing data. PloS One 7:e30619.

1157    Price, M.N., Dehal, P.S., and Arkin, A.P. (2009). FastTree: computing large minimum
1158        evolution trees with profiles instead of a distance matrix. Mol Biol Evol
1159        26:1641-1650.

1160    Qin, P., Lu, H., Du, H., Wang, H., Chen, W., Chen, Z., He, Q., Ou, S., Zhang, H., and
1161        Li, X., et al. (2021). Pan-genome analysis of 33 genetically diverse rice
1162        accessions reveals hidden genomic variations. Cell 184:3542-3558.

1163    Qiu, J., Jia, L., Wu, D., Weng, X., Chen, L., Sun, J., Chen, M., Mao, L., Jiang, B., and
1164        Ye, C., et al. (2020). Diverse genetic mechanisms underlie worldwide convergent
1165        rice feralization. Genome Biol 21.

1166    Qiu, J., Zhou, Y., Mao, L., Ye, C., Wang, W., Zhang, J., Yu, Y., Fu, F., Wang, Y., and
1167        Qian, F., et al. (2017). Genomic variation associated with local adaptation of
1168        weedy rice during de-domestication. Nat Commun 8.

1169    Read, A.C., Moscou, M.J., Zimin, A.V., Pertea, G., Meyer, R.S., Purugganan, M.D.,
1170        Leach, J.E., Triplett, L.R., Salzberg, S.L., and Bogdanove, AJ. (2020) Genome
1171        assembly and characterization of a complex zfBED-NLR gene-containing
1172        disease resistance locus in Carolina Gold Select rice with Nanopore sequencing.
1173        PLoS Genet 16(1):e1008571.

1174    Rhie, A., Walenz, B.P., Koren, S., and Phillippy, A.M. (2020). Merqury: reference-free
1175        quality, completeness, and phasing assessment for genome assemblies. Genome
1176        Biol 21(1):245.

1177  Rozas, J., Ferrer-Mata, A., Sánchez-DelBarrio, J.C., Guirao-Rico, S., Librado, P.,
1178      Ramos-Onsins, S.E., and Sánchez-Gracia, A. (2017). DnaSP 6: DNA sequence
1179      polymorphism analysis of large data sets. Mol Biol Evol 34(12):3299-3302.
1180  Salamov, A.A., and Solovyev, V.V. (2000) *Ab initio* gene finding in *Drosophila*
1181      genomic DNA. Genome Res 10(4):516-522.
1182  Shang, L., Li, X., He, H., Yuan, Q., Song, Y., Wei, Z., Lin, H., Hu, M., Zhao, F., and
1183      Zhang, C., et al. (2022). A super pan-genomic landscape of rice. Cell Res
1184      32:878-896.
1185  Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., and Zdobnov, E.M.
1186      (2015). BUSCO: assessing genome assembly and annotation completeness with
1187      single-copy orthologs. Bioinformatics 31:3210-3212.
1188  Song, B., Chuah, T., Tam, S.M., and Olsen, K.M. (2014). Malaysian weedy rice
1189      shows its true stripes: wild *Oryza* and elite rice cultivars shape agricultural weed
1190      evolution in southeast Asia. Mol Ecol 23:5003-5017.
1191  Song, J., Xie, W., Wang, S., Guo, Y., Koo, D., Kudrna, D., Gong, C., Huang, Y., Feng,
1192      J., and Zhang, W., et al. (2021). Two gap-free reference genomes and a global
1193      view of the centromere architecture in rice. Mol Plant 14:1757-1767.
1194  Stanke, M., Keller, O., Gunduz, I., Hayes, A., Waack, S., and Morgenstern, B. (2006)
1195      AUGUSTUS: *ab initio* prediction of alternative transcripts. Nucleic Acids Res
1196      34:W435-9.
1197  Steuernagel, B., Jupe, F., Witek, K., Jones, J.D., and Wulff, B.B. (2015). NLR-parser:
1198      rapid annotation of plant NLR complements. Bioinformatics 31(10):1665-1667.
1199  Su, W., Ou S., Hufford, M.B., and Peterson, T. (2021) A tutorial of EDTA: Extensive
1200      De Novo TE Annotator. Methods Mol Biol 2250:55-67.
1201  Sun, J., Ma, D., Tang, L., Zhao, M., Zhang, G., Wang, W., Song, J., Li, X., Liu, Z.,
1202      and Zhang, W., et al. (2019). Population genomic analysis and *de novo* assembly
1203      reveal the origin of weedy rice as an evolutionary game. Mol Plant 12:632-647.
1204  Sun, J., Zhang, G., Cui, Z., Cui Z., Kong X., Yu X., Gui R., Han Y., Li Z., and Lang
1205      H., et al. (2022) Regain flood adaptation in rice through a 14-3-3 protein
1206      OsGF14h. Nat Commun 13, 5664 (2022).
1207  Tanabe, S., Ashikari, M., Fujioka, S., Takatsuto, S., Yoshida, S., Yano, M., Yoshimura,
1208      A., Kitano, H., Matsuoka, M., and Fujisawa, Y, et al. (2005) A novel cytochrome
1209      P450 is implicated in brassinosteroid biosynthesis via the characterization of a
1210      rice dwarf mutant, *dwarf11*, with reduced seed length. Plant Cell 17(3):776-790.
1211  Thorvaldsdóttir, H., Robinson, J.T., and Mesirov, J.P. (2013) Integrative Genomics
1212      Viewer (IGV): high-performance genomics data visualization and exploration.
1213      Brief Bioinform 14(2):178-192.
1214  Tseng, I., Hong, C., Yu, S., and Ho, T.D. (2013). Abscisic acid- and stress-induced
1215      highly proline-rich glycoproteins regulate root growth in rice. Plant Physiol
1216      163:118-134.
1217  Vaser, R., Sović, I., Nagarajan, N., and Šikić, M. (2017). Fast and accurate *de novo*
1218      genome assembly from long uncorrected reads. Genome Res 27:737-746.
1219  Wang, H., Vieira, F.G., Crawford, J.E., Chu, C., and Nielsen, R. (2017). Asian wild
1220      rice is a hybrid swarm with extensive gene flow and feralization from

1221      domesticated rice. Genome Res. 27(6):1029-1038.

1222 Wang, J., Deng, Q., Li, Y., Yu, Y., Liu, X., Han, Y., Luo, X., Wu, X., Ju, L., and Sun, J.,
1223      et al. (2020). Transcription factors *Rc* and *OsVP1* coordinately regulate
1224      preharvest sprouting tolerance in red pericarp rice. J Agr Food Chem
1225      68:14748-14757.

1226 Wang, L., Zhao, L., Zhang, X., Zhang, Q., Jia, Y., Wang, G., Li, S., Tian, D., Li, W.,
1227      and Yang, S. (2019). Large-scale identification and functional analysis of NLR
1228      genes in blast resistance in the Tetep rice genome sequence. Proc Natl Acad Sci
1229      U S A 116:18479-18487.

1230 Wang, M., Li, W., Fang, C., Xu, F., Liu, Y., Wang, Z., Yang, R., Zhang, M., Liu, S.,
1231      and Lu, S., et al. (2018a). Parallel selection on a dormancy gene during
1232      domestication of crops from multiple families. Nat Genet 50:1435-1441.

1233 Wang, W., Mauleon, R., Hu, Z., Chebotarov, D., Tai, S., Wu, Z., Li, M., Zheng, T.,
1234      Fuentes, R.R., Zhang, F., et al. (2018b). Genomic variation in 3,010 diverse
1235      accessions of Asian cultivated rice. Nature 557:43-49.

1236 Wang, Z., Wang, W., Xie, X., Wang, Y., Yang, Z., Peng, H., Xin, M., Yao, Y., Hu, Z.,
1237      and Liu, J., et al. (2022). Dispersed emergence and protracted domestication of
1238      polyploid wheat uncovered by mosaic ancestral haploblock inference. Nat
1239      Commun 13.

1240 Wang, Z., Wei, K., Xiong, M., Wang, J.D., Zhang, C.Q., Fan, X.L., Huang, L.C., Zhao,
1241      D.S., Liu, Q.Q., and Li, Q.F. (2021) Glucan, Water-Dikinase 1 (*GWD1*), an ideal
1242      biotechnological target for potential improving yield and quality in rice. Plant
1243      Biotechnol J 19(12):2606-2618.

1244 Wu, D., Lao, S., and Fan, L. (2021). De-domestication: an extension of crop evolution.
1245      Trends Plant Sci 26:560-574.

1246 Wu, D., Qiu, J., Sun, J., Song, B., Olsen, K.M., and Fan, L. (2022a). Weedy rice, a
1247      hidden gold mine in the paddy field. Mol Plant 15:566-568.

1248 Wu, D., Shen E., Jiang, B., Feng, Y., Tang, W., Lao, S., Jia, L., Lin, H.Y., Xie, L., and
1249      Weng, X., et al. (2022b) Genomic insights into the evolution of *Echinochloa*
1250      species as weed and orphan crop. Nat Commun 13(1):689.

1251 Wu, T.D., and Watanabe, C.K. (2005) GMAP: a genomic mapping and alignment
1252      program for mRNA and EST sequences. Bioinformatics 21(9):1859-75.

1253 Xie, X., Du, H., Tang, H., Tang, J., Tan, X., Liu, W., Li, T., Lin, Z., Liang, C., and Liu,
1254      Y. (2021). A chromosome-level genome assembly of the wild rice *Oryza*
1255      *rufipogon* facilitates tracing the origins of Asian cultivated rice. Science China
1256      Life Sciences 64:282-293.

1257 Zhan, C., Lei, L., Liu, Z., Zhou, S., Yang, C., Zhu, X., Guo, H., Zhang, F., Peng, M.,
1258      and Zhang, M., et al. (2020). Selection of a subspecies-specific diterpene gene
1259      cluster implicated in rice disease resistance. Nat Plants 6:1447-1454.

1260 Zhang, F., Wang, C., Li, M., Cui, Y., Shi, Y., Wu, Z., Hu, Z., Wang, W., Xu, J., and Li,
1261      Z. (2021). The landscape of gene-cds-haplotype diversity in rice: properties,
1262      population organization, footprints of domestication and breeding, and
1263      implications for genetic improvement. Mol Plant 14:787-804.

1264 Zhang, F., Xue, H., Dong, X., Li, M., Zheng, X., Li, Z., Xu, J., Wang, W., and Wei, C.

(2022). Long-read sequencing of 111 rice genomes reveals significantly larger pan-genomes. Genome Res 32(5):853-863.

Zhang, X., Zhang, S., Zhao, Q., Ming, R., and Tang, H. (2019). Assembly of allele-aware, chromosomal-scale autopolyploid genomes based on Hi-C data. Nat Plants 5:833-845.

Zhao, Q., Feng, Q., Lu, H., Li, Y., Wang, A., Tian, Q., Zhan, Q., Lu, Y., Zhang, L., Huang, T., et al. (2018). Pan-genome analysis highlights the extent of genomic variation in cultivated and wild rice. Nat Genet 50:278-284.

Zheng, X., Pang, H., Wang, J., Yao, X., Song, Y., Li, F., Lou, D., Ge, J., Zhao, Z, and Qiao, W., et al. (2022). Genomic signatures of domestication and adaptation during geographical expansions of rice cultivation. Plant Biotechnol J 20(1):16-18.

Zhou, Y., Chebotarov, D., Kudrna, D., Llaca, V., Lee, S., Rajasekar, S., Mohammed, N., Al-Bader, N., Sobel-Sorenson, C., and Parakkal, P., et al. (2020). A platinum standard pan-genome resource that represents the population structure of Asian rice. Sci Data 7(1):113.

Zhu, B., Si, L., Wang, Z., Jingjie Zhu, Y.Z., Shangguan, Y., Lu, D., Fan, D., Li, C., Lin, H., and Qian, Q., et al. (2011). Genetic control of a transition from black to straw-white seed hull in rice domestication. Plant Physiol 155:1301-1311.

Zdobnov, E.M., and Apweiler, R. (2001). InterProScan-an integration platform for the signature-recognition methods in InterPro. Bioinformatics 17(9):847-848.

## Author contributions

L.F. and Q.Q. conceived and supervised the study. D.W., L.J. and C.D. assembled and annotated the genomes. D.W. and L.J. collected previously released assemblies of rice genomes. Y.H. and L.X. evaluated the quality of rice assemblies. D.W., L.X. and Y.H. constructed the rice syntelog-based pan-genome and analyzed the ancestral haplotypes. L.X. carried out the analysis of NLR genes. Y.S. and L.X. performed the analysis of structural variations. L.F., C.Y. and Q.Q. discussed the results. D.W., Y.S. and L.X. wrote the manuscript and L.F. and C.Y. revised it. All authors discussed the results and commented on the manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

The supplementary material is available online.

**Figure legends**

**Figure 1. Quality assessment of rice genome assemblies and their relationship based on phylogeny and synteny.** (**a**) Consensus quality values and number of homozygous SNPs based on self-mapping for genomes in four pan-genome projects. (**b**) The relationship among rice genomes measured by genetic distance (IBD, identity-by-descent) based on SNPs and synteny based on gene orders.

**Figure 2. Syntelog-based pan-genome of rice and differential gene PAVs in subspecies.** (**a**) Number of syntelog groups (SGs) represented in all 74 rice genomes versus the number of genomes. The subspecies pan-genome compositions for GJ and XI were extracted from the whole rice pan-genome. The numbers of core, soft-core, dispensable, private and absent SGs were counted. (**b**) Frequency differences in gene presence between GJ and XI along chromosome 7. Each dot represents one SG. The gene PAV profiling of the antimicrobial diterpenoid biosynthetic gene cluster DGC7 and tandem duplicates *RePRP2.1* and *RePRP2.2*, suppressors of root cell expansion, is shown.

**Figure 3. Syntelog-based ancestral haplotypes suggest widespread genomic introgression in rice evolution.** (**a**) Ancestral haplotype landscape of SGs on chromosome 4. An SNP-based phylogeny of rice groups is illustrated on the left. For each SG, seven blocks with different colors represent different haplotypes of predicted protein sequences. Candidate introgression regions are numbered by genomic length and marked by gray blocks. The centromere position information of chromosome 4 is obtained from the Rice Genome Annotation Project and indicated by a red triangle. Functionally important genes are annotated within each block, and domestication genes are highlighted in red. (**b**) Haplotype divergence between XI and GJ for all SGs on chromosome 4. (**c**) Significance test on the non-random distribution of clustered SGs in introgression blocks by sampling 100,000 replicates. The horizontal red dashed line represents a *P* value of 0.01. (**d**) Phylogeny analysis of four large introgression blocks (length greater than 300 kb) indicates a single origin of domestication alleles from the Or-3 wild group. Four wild groups are highlighted in different colors. (e) Final putative introgression blocks from proto-GJ to XI, combing the evidence from haplotype divergence, phylogeny and ABBA-BABA tests. (**f**) Complex genomic contributions from wild groups to the emergence of *aus* group revealed by the phylogenetic trees in 64 introgression blocks.

**Figure 4. Structural variations in rice de-domestication.** (**a**) Dot plots comparing all 74 assemblies against the de-domestication genomic island (from 6.0 to 6.5 Mb) on chromosome 7 of Nipponbare. The predominant translocations in wild and weedy accessions are highlighted by red boxes. The pink seed icons behind the accession IDs represent red or brown pericarp. (**b**) Phylogeny revealed by SNPs in the *Rc* region indicates introgression from Or-1 and Or-3 to XI and GJ, respectively. Labels with green, red, yellow and blue represent accessions from Or-4, Or-3, Or-2 and Or-1,

1355 respectively. Blue dots represent weedy accessions. The numbers on each branch
1356 indicate bootstrap values of less than 90%, based on 1000 replicates. (**c**) Structural
1357 variations in *OsC1* between the weedy and cultivated GJ accessions. Black rectangles
1358 represent exons. (**d**) Structural variations in *Bh4* between weedy and cultivated
1359 accessions and their seed appearances.
1360
1361 **Figure 5. A brief schematic illustration of Asian rice evolution.** The evolutionary
1362 scenario highlights that complex introgression events have contributed indispensably
1363 to rice domestication and de-domestication.

**Supplementary Information**

**Supplementary Table 1.** Meta information of 75 rice genome assemblies used in this study.

**Supplementary Table 2.** Haplotype divergence between XI and GJ for domestication and improvement genes

**Supplementary Table 3.** Putative introgression blocks identified by haplotype divergence between GJ and XI and validation by phylogeny and ABBA-BABA tests

**Supplementary Table 4.** Gene functional enrichment in final putative introgression blocks

**Supplementary Table 5.** Cloned genes in the putative introgression blocks

**Supplementary Table 6.** Structural variations between each pair of weedy and cultivated accessions in agronomy-related genes

**Supplementary Fig. 1** Hi-C interaction heatmaps for each chromosome from four rice accessions (a, YCW03; b, 18XHB-83; c, CX20; d, NJ11).

**Supplementary Fig. 2** Statistical information of the rice genomes used in this study. (a) Assembly (size), annotation (gene number, TE size and proportion) and quality assessment (contig N50, BUSCO and LAI) of rice genomes. (b) Differences in genomic features between subspecies GJ and XI. In the boxplots, the horizontal line shows the median value, and the whiskers show the 25% and 75% quartile values of each genomic feature. *P* values were calculated by Student's *t* test.

**Supplementary Fig. 3** Dot plots of newly generated rice assemblies in this study against the reference assembly Nipponbare (IRGSP) and gapless assembly MH63RS3.

**Supplementary Fig. 4** Assembly quality assessment in base accuracy. (a) QVs for rice genomes in four rice pan-genome projects using yak (https://github.com/lh3/yak). (b) The number and annotation to SNPs and InDels for each genome by mapping NGS short reads against their own assembly.

**Supplementary Fig. 5** PCA reveals the representativeness and diversity of genome assemblies used in this study. The first two principle components are shown. Filled circles indicate the assemblies used in this study. Green, blue and red represent wild, cultivated and weedy assemblies, respectively.

**Supplementary Fig. 6** Performance of BLASTP and Diamond (under different modes) in syntelog identification. (a) The BLASTP approach identifies more syntelogs than Diamond. (b) Syntelog number between accession 02428 and other accessions identified using BLASTP and Diamond.

**Supplementary Fig. 7** Pairwise synteny reveals evolutionary signatures in groups and individuals. (a) Syntelog numbers between rice groups. (b) Syntelog numbers between each accession and other accessions from the GJ and XI (including *aus*) subspecies. In the boxplots, the horizontal line shows the median value, and the whiskers show the 25% and 75% quartile values of syntelog numbers.

**Supplementary Fig. 8** Comparison of MCL and synteny-based clustering. (a) A brief scheme illustrating MCL ortholog clustering and syntelog clustering. (b) Group size

1408 comparison using synteny-based clustering (SynPan) and MCL clustering
1409 (OrthoFinder).

1410 **Supplementary Fig. 9** Benchmarking analysis on the influences of inflation
1411 parameters in the MCL clustering in rice genomes. (a) The OG numbers shared by
1412 different sizes of rice genomes under different inflation parameters. (b) Gene counts
1413 in OGs shared by different sizes of genomes. (c) Average gene number per OG under
1414 different inflation parameters.

1415 **Supplementary Fig. 10** Composition and features of the rice syntelog-based
1416 pangenome. (a) Pan-gene composition (core, soft-core, dispensable and private) of the
1417 rice pan-genome. (b) Proportions of domain-annotated genes in four categories of the
1418 rice pan-genome. (c) Percentage of SGs (blue) and genes (red) with domain
1419 gain-and-loss.

1420 **Supplementary Fig. 11** Comparison of NLR genes in rice genomes from different
1421 ecotypes and subspecies. (a) Distribution of different NLR genes in wild, cultivated
1422 and weedy accessions. (b) Distribution of different NLR genes in wild rice, XI and GJ.
1423 (c) Clustered NLRs in wild rice, XI and GJ. $P$ values are calculated using Wilcoxon
1424 test. (d) Relationship between genome assembly completeness (as indicated by LAI)
1425 and NLR size. In the boxplots, the horizontal line shows the median value, and the
1426 whiskers show the 25% and 75% quartile values of NLR sizes.

1427 **Supplementary Fig. 12** Genomic features of NLRs in rice genomes. (a) Dynamics of
1428 domain architectures in the rice NLRome. CNL, NL and null (no canonical
1429 architectures) are defined as three NLR architecture types in rice. Multiple types are
1430 observed for most SGs. (b) Nucleotide diversity and selection of NLRs from core,
1431 soft-core and dispensable SGs. In the boxplots, the horizontal line shows the median
1432 value, and the whiskers show the 25% and 75% quartile values of Pi and Tajima's $D$.
1433 Significance $P$ values are performed using Wilcoxon test. *, $P < 0.05$.

1434 **Supplementary Fig. 13** Haplotype analysis on rice syntelogs. (a) Definition of
1435 haplotype diversity and divergence. Haplotype diversity and divergence represent
1436 average haplotype differences among sequences in a syntelog group within one group
1437 and among two groups, respectively, where Xi and Xj are the presence count of
1438 haplotype i and haplotype j in group X, and $\sum$X is the total sequence count within a
1439 syntelog group. (b) Brief scheme illustrating the assignment and visualization of
1440 ancestral haplotypes for each genome. A group priority-based referring strategy is
1441 used to assign ancestral haplotypes. Group information is prior based on
1442 whole-genome phylogeny or population structure. The most dominant sequence in
1443 syntelogs from Group 1 is set as hapI. If the most dominant sequence from Group 2 is
1444 not HapI, then define HapII, otherwise HapII is skipped. If the most dominant
1445 sequence from Group 3 is neither HapI nor HapII, then define HapIII, otherwise
1446 HapIII is skipped. By analogy, dominant haplotypes are determined and colored for
1447 each syntelog group. Rare haplotypes are named as HapR colored by dark gray. In this
1448 study, the group priority order is set as tmp > XI1A > XI1B > *aus* > XI3. Different
1449 orders have no influences on the calculation of haplotype diversity and divergence. (c)
1450 mosaic graphs of ancestral haplotypes in chromosome 1 across 74 rice genomes with
1451 different priority orders, tmp > XI1A > XI1B > *aus* > XI3, XI1B > XI1A > tmp >

1452  *aus* > XI3 and XI1A > XI1B > tmp > *aus* > XI3. For each window, the same color

1453  indicates the same haplotype and dark and light gray indicates rare haplotypes and

1454  syntelog absence.

1455  **Supplementary Fig. 14** Genetic diversity measured by full-length gene, coding

1456  region, and predicted protein sequences in rice genomes. (a) Comparison of diversity

1457  measured by haplotype N100, N90 and haplotype diversity using full-length

1458  nucleotide sequences, coding sequences and predicted protein sequences. In the

1459  boxplots, the horizontal line shows the median value, and the whiskers show the 25%

1460  and 75% quartile values of each diversity indice. (b) Haplotype diversity in GJ and XI,

1461  compared to all rice genomes.

1462  **Supplementary Fig. 15** Ancestral haplotype landscape on chromosome 1 to

1463  chromosome 12. Putative introgression blocks are shown in gray rectangles and

1464  numbered. Haplotype divergence and *P* values (scaled by -log10) of non-random

1465  distribution significance tests are shown along each chromosome. Red dashed lines

1466  represent thresholds to determine introgression blocks.

1467  **Supplementary Fig. 16** Synonymous substitution rates ($K_s$) of genes in putative

1468  introgression blocks and their neighboring left and right regions. Three replicates (a, b,

1469  and c) between the GJ and XI genomes were performed. In the boxplots, the

1470  horizontal line shows the median value, and the whiskers show the 25% and 75%

1471  quartile values of $K_s$. *P* values are calculated using Wilcoxon test.

1472  **Supplementary Fig. 17** Population structure of wild rice accessions used in this study.

1473  (a) Phylogenetic tree of *Oryza rufipogon* and *Oryza sativa*. Four wild groups (Or-1,

1474  Or-2, Or-3 and Or-4) are indicated in different colors. (b) PCA plots of the first three

1475  principle components, where "W", "J" and "I" represent wild, GJ and XI accessions.

1476  (c) Geographical sources of wild accessions in different groups used in this study.

1477  **Supplementary Fig. 18** Introgression $f_d$ distributions of ABBA-BABA test on

1478  chromosome 1 to chromosome 12 in topology T1 to T6, where P1 is Or-1 (*n* = 37), P3

1479  is tmp (GJ, *n* = 100), O/outgroup is Or-4 (*n* = 25), and P2 was set as Or-2 (*n* = 42),

1480  XI1A (*n* = 100), XI1B (*n* = 100), XI2 (*n* = 80), XI3 (*n* = 100) and *aus* (*n*=60),

1481  respectively. T1 was set as a background control in introgression detection. Genomic

1482  positions of putative introgression regions are indicated by gray rectangles and

1483  detailed coordinates are provided in Supplementary Table 3. Blocks larger than 300kb

1484  are highlighted in red and blocks not supported by $f_d$ are underlined.

1485  **Supplementary Fig. 19** Comparison of $f_d$ using all genomes and landraces only under

1486  different topologies on chromosome 1. (a) $f_d$ distribution along chromosome 1. Group

1487  tmp, XI1A, XI2, XI3, and *aus* include 47, 24, 21, 52 and 32 landrace accessions,

1488  respectively. (b) Comparison of $f_d$ on chromosome 1 in T2 *vs* T12, T4 *vs* T14, and T5

1489  *vs* T15.

1490  **Supplementary Fig. 20** Genome synteny between each weedy rice assembly and the

1491  corresponding phylogenetically closest cultivated rice assembly.

1492  **Supplementary Fig. 21** Summary of structural variations between weedy and

1493  cultivated rice genomes.

1494  **Supplementary Fig. 22** Structural variations (>50 bp) between assemblies of cultivar

1495  accession NJ11 and weedy rice accession CX20 on 12 chromosomes. The largest six

1496   SVs (numbered from 1 to 6) are zoomed in and annotated, including four insertions (1,

1497   2, 5 and 6) and two deletions (3 and 4) in CX20.

1498   **Supplementary Fig. 23** The Integrative Genomics Viewer (IGV) snapshots show the

1499   structural variations in *OsC1* and *Bh4* between weedy and cultivated rice.

1500   **Supplementary Fig. 24** Phylogeny of *Bh4* in rice and morphology of rice seed hulls.

1501   Bootstrap values less than 0.90 are indicated on branches.

# Figure 1



**a**

Rice pan-genome projects
- 1: this study (*n* = 10)
- 2: Zhang et al., 2022, Genome Res. (*n* = 11)
- 3: Qin et al., 2021, Cell (*n* = 30)
- 4: Zhou et al., 2020, Sci. Data (*n* = 10)

QV (consensus quality value)

$\log_{10}$(number of homozygous SNPs)

**b**

Or-4
Or-3
aromatic
trp
subtrp
tmp
Or-2
Or-1
aus
XI2
XI3
XI1A
XI1B

GJ ssp. *japonica*
XI ssp. *indica*

outgroup
wild (*O. rufipogon*)
cultivated (*O. sativa*)
weedy (*O. sativa* ssp. *spontanea*)

Syntelog Proportion

0.7  0.8  0.9

**c**

Pearson's correlation = 0.886
*P* value < 2.2e-16

Synteny affinity

Genetic affinity (IBD)

- GJ *vs* GJ
- XI *vs* XI
- aus *vs* aus
- GJ *vs* aus
- XI *vs* aus
- GJ *vs* XI
- wild vs all

# Figure 2

# Figure 3



**a** Chromosome 4    Haplotype    hapI  hapII  hapIII  hapIV  hapV  hapR  Absence

4#5    centromere position    4#2    *An-1*    4#8  4#9  4#6  4#4  4#7    4#13    4#1    4#12  4#11  4#10    4#3

**b** Divergence
*An-1*    *LABA1*    *sh4*
*D11*
*Bh4*    *OsGA2ox6*  *OsAMT1.1*
*OsNAC2*    *OsARF11*  *RLI1*  *PAO5*

**c** -log₁₀(*P* value)

**d**
Chr04#01 (length:1753kb)    Chr04#02 (length:607kb)

Or-4
Or-3
Or-2
Or-1

Chr04#03 (length:389kb)    Chr04#04 (length:371kb)

**e** Block size  <300kb  >300kb
chromosome

**f**
Or-4        *aus*        Or-4
Or-2        Or-4        Or-2
Or-1        Or-2        *aus*
Or-3        Or-1        Or-1
XI          Or-3        Or-3
GJ          XI          XI
*aus*        GJ          GJ
**32/64**    **8/64**    **24/64**

# Figure 4



a

Ecotype ● wild ● cultivated ● weedy

b

Or-4
Or-3
Or-2
Or-1

● Weedy rice

c

18WR-118 chr06.732
Nipponbare chr06.742        +983bp
13-65 chr06.729
YCW03 chr06.734    +3bp        +2bp
Nipponbare chr06.742
WR04-6 chr06.739    +3bp

*OsC1*

d

Nipponbare chr04.2141
WR04-6 chr04.2093        +22bp
N22 Chr4.2027
PI653439 Chr4.2151    +70bp    +22bp

*Bh4*

Nipponbare    WR04-6    N22    PI653439

# Figure 5