

1 **Title:**

2 Machine Learning Reveals Spatiotemporal Genome Evolution in Asian Rice
3 Domestication

4

5 **Authors:**

6 Hajime Ohyanagi, Kosuke Goto, Sónia Negrão, Rod A. Wing, Mark A. Tester, Kenneth L.
7 McNally, Vladimir B. Bajic, Katsuhiko Mineta, Takashi Gojobori*

8

9 **Authors' Affiliations:**

10 *King Abdullah University of Science and Technology (KAUST), Computational Bioscience Research*
11 *Center (CBRC), Biological and Environmental Sciences & Engineering Division (BESE), Thuwal,*
12 *23955-6900, Saudi Arabia*

13 Hajime Ohyanagi, Kosuke Goto & Takashi Gojobori

14

15 *School of Biology and Environmental Science, University College Dublin, Belfield, Ireland*

16 Sónia Negrão

17

18 *King Abdullah University of Science and Technology (KAUST), Biological and Environmental Sciences*
19 *& Engineering Division (BESE), Thuwal, 23955-6900, Saudi Arabia*

20 Rod A. Wing & Mark A. Tester

21

22 *King Abdullah University of Science and Technology (KAUST), Computational Bioscience Research*
23 *Center (CBRC), Computer, Electrical and Mathematical Sciences & Engineering Division (CEMSE),*
24 *Thuwal, 23955-6900, Saudi Arabia*

25 Vladimir B. Bajic & Katsuhiko Mineta

26

27 *International Rice Research Institute, Manila, Philippines*

28 Kenneth L. McNally

29

30 ***Corresponding Author:**

31 Takashi Gojobori, Distinguished Professor,

32 *King Abdullah University of Science and Technology (KAUST), Thuwal 23955-6900, Saudi Arabia*

33 *E-mail: takashi.gojobori@kaust.edu.sa, Phone: +966-12-808-2893*

34

35 **Domestication is anthropogenic evolution that fulfills mankind's critical food**
36 **demand. As such, elucidating the molecular mechanisms behind this process**
37 **promotes the development of future new crops. With the aim of understanding the**

38 whole domestication process of Asian rice and by employing the *Oryza sativa*
39 subspecies (*indica* and *japonica*) as an Asian rice domestication model, we
40 scrutinized genomic introgressions between them as traces of domestication. Here
41 we show the genome-wide introgressive region (IR) map of Asian rice, by utilizing
42 4,587 accession genotypes with a stable outgroup species, particularly at the finest
43 resolution through a machine learning-aided method. The IR map revealed that
44 14.2% of the rice genome consists of IRs, including both wide IRs (recent) and
45 narrow IRs (ancient). This introgressive landscape with their time calibration
46 indicates that introgression events happened in multiple genomic regions over
47 multiple periods. From the correspondence between our wide IRs and so-called
48 Selective Sweep Regions, we provide a definitive answer to a long-standing
49 controversy in plant science: Asian rice phylogeny appears to depend on which
50 regions and time frames are examined.

51

52 Rice is one of the most essential crops to humankind, playing a critical role in food
53 security ¹. Since it has been domesticated to fit it to humanity's needs, its genome holds
54 the secrets to ancient and modern agricultural practices, which can serve as an
55 informative reference for future breeding practices. Rice domestication history can be
56 divided into three independent episodes: *Oryza nivara* (also known as annual *O.*
57 *rufipogon* or Or-I) and *O. rufipogon* in Asia that led to Asian rice (*O. sativa* L.) ², *O.*
58 *barthii* that was domesticated by early African farmers around 3,000 years ago and led to
59 African rice (*O. glaberrima* Steud.) ³, and a New World rice domestication process by
60 Amazon farmers around 4,000 years ago that occurred in South America ⁴. In particular,

61 the Asian domesticated rice (*O. sativa*) is the most prominent species in the genus *Oryza*,
62 which has served as the major staple crop in most Asian countries for millennia.

63 Among these three domesticated rice species, Asian rice (*O. sativa*) and its origins
64 have been the most intensively studied and continue to be debated in both archeological
65 and genetic research areas ⁵⁻²⁰. In short, two conflicting domestication hypotheses have
66 been proposed: 1) a single domestication process where a single subspecies (either *indica*
67 or *japonica*) was first domesticated from a wild rice, while the other arose from a
68 hybridization with another wild rice species; and 2) independent domestication processes
69 where different species of *O. nivara* and *O. rufipogon* with distinct Asian origins gave
70 rise to different domesticated subspecies.

71 A comprehensive SNP-based genomic phylogeny (*i.e.*, a genomic phylogeny as a
72 whole) clearly shows that at least two origins of *O. sativa* subspecies exist¹⁴, *i.e.*, *O.*
73 *sativa* ssp. *indica* and *O. nivara* cluster with each other, while *O. sativa* ssp. *japonica* and
74 *O. rufipogon* make another cluster. However, this is just a subspecies phylogeny, which
75 does not reflect the domestication history. To trace back the history, plant scientists have
76 been focusing on their own self-defining genomic entities, *e.g.*, domestication-associated
77 gene regions (with flanking upstream/downstream regions), Selective Sweep Regions ¹⁴,
78 Co-located Low-Density Genomic Regions ¹⁰, transposable elements ⁶, microsatellites ¹²,
79 and so forth. In other words, there have been multiple definitions for domestication-
80 derived regions. Meanwhile, phylogenies inferred by plant scientists do not always agree
81 with one another, either supporting theories 1) or 2). In fact, the domesticated Asian rice
82 accessions have supposedly introduced agronomically advantageous traits from one
83 subspecies to another during the domestication process ^{7,9,20-22}. Therefore, their genomes

84 are presumed to be mosaics since they have been exchanging alleles over introgression
85 events throughout history. In this sense, the controversy over the origins of rice
86 domestication arose from the disagreed domestication-derived regions. Moreover, the
87 window size studied is a critical factor in the controversy. In our study, the phylogenetic
88 analysis of a domestication-associated gene with variable lengths of
89 upstream/downstream flanking regions, as also shown in the result of Choi & Purugganan
90 ⁸ showed that the gene window size profoundly affects the resultant gene phylogenies
91 (shown in **Fig. 2e, f, g, h, and i**, details will be described in **Consequence of Analysis**
92 **Window Size**).

93 Given that introgression events are representative of human intervention (*i.e.*, the
94 domestication process), our simple and robust rationale is not to focus on particular
95 genomic regions, but rather to exhaustively detect any introgressive regions (IRs)
96 between subspecies as traceable signs of domestication, employing windows with as fine
97 a resolution as possible. In keeping with this notion, we present not only gene-by-gene
98 introgressive states but also a genome-wide IR map between *O. sativa* ssp. *indica* and ssp.
99 *japonica* at the finest resolution using an efficient machine learning model, with the aim
100 of revealing the whole domestication process of Asian rice.

101

102 **Results**

103 **Invention of *Distance Difference (DD)* to Detect Introgressions**

104 To capture the entire introgressive landscape of domesticated Asian rice genomes using a
105 large-scale genotype set (**Fig. 1a and b**), we needed to overcome three major difficulties
106 described in the **Methods**. In short, i) the low density of rice genotypes, ii) over-diversity

107 within each subspecies (**Fig. 1c**), and iii) the instability of an outgroup. To overcome
108 these challenges, we employed 14x coverage genotypes supplied by the 3,000 Rice
109 Genomes Project²²⁻²⁵, a median 10th subset extraction from the comprehensive dataset,
110 and a reproductively isolated accession of *O. punctata* (BB diploid, 2n=24, with African
111 geographical origin)²⁶ as an outgroup species. For more details, see **Methods**.

112 Each domesticated subpopulation has its own particular evolutionary rate²⁷.
113 Therefore, each of *indica* and *japonica* subpopulations should show, to some extent,
114 different genetic distances to an outgroup (a wild rice accession), since they have been
115 separated from each other for a length of time (**Fig. 2a**) with the assumption that any
116 inter-subspecies cross (*i.e.*, an introgression) has not occurred. On the other hand, they
117 will show more similar genetic distances to the outgroup when an inter-subspecies cross
118 has occurred (**Fig. 2b**). In particular, subspecies in domesticated plants have been
119 artificially forced to make inter-subspecies crossings in order to introduce agronomically
120 important traits, thereby particular regions of their genomes must be strongly affected by
121 the decrease in difference of genetic distance (distance difference).

122 Even though this decrease may disturb an accurate inference of genetic phylogeny of
123 rice subspecies and wild relatives, it can be paradoxically utilized as an index of
124 introgression, *i.e.*, once a decrease is observed, it is a possible sign of an introgression
125 event. To distinguish IRs from non-IRs (**Fig. 2a** and **b**), we conceptually defined *DD*
126 (*Distance Difference* to the outgroup: A unit is number of substitutions per nucleotide
127 site) as:

$$128 \quad DD = |F84(\text{outgroup to } indica) - F84(\text{outgroup to } japonica)|$$

129 ^(*) F84 = Felsenstein84 nucleotide genetic distance²⁸

130 Here, the regions with smaller *DDs* represent IRs, while the regions with larger *DDs*
131 represent non-IRs. For more details, see **Methods**. Note that because IRs at the very early
132 stage of domestication will not show enough decrease in *DDs*, IRs of very ancient origin
133 are out of scope of this method.

134

135 **Incoherent Introgressive States of Domestication-associated Genes (D-genes)**

136 Based on the logic above, we firstly aimed to determine *DDs* of 25 manually curated
137 domestication-associated genes (D-genes, **Fig. 2c**) as indices of their introgressive states.
138 To archive the best accuracy in this limited scale analysis, we constructed 25 gene-by-
139 gene phylogenetic trees without any flanking upstream/downstream regions, and we
140 visually inspected their *DDs* thoroughly, to determine whether *indica* and *japonica* show
141 a similar genetic distance to the outgroup, or different genetic distances to the outgroup.
142 Our results show that incoherent introgressive states of D-gene regions, *i.e.* nine D-genes
143 (*Bh4*, *C1*, *GAD1*, *LABA1*, *LG1*, *Progl*, *qSW5*, *Rc*, and *sh4*) out of 25, are introgressive
144 (regardless of the direction), whereas 14 D-genes (*BADH2*, *Bph14*, *DPL2*, *Ehd1*, *Ghd7*,
145 *Gn1a*, *GS3*, *GW2*, *Phr1*, *qSH1*, *Rd*, *sd1*, *tb1*, and *waxy*) are not (**Fig. 2c** and **d**, yellow =
146 non-introgressive, red = introgressive, full size phylogenetic tree pictures with detailed
147 color system are shown in **Supplementary Fig. 1**). *Hd1* and *S5* have status-undetermined.
148 Through a statistical analysis (**Supplementary Table 2**), we found significant enrichment
149 in the introgressed proportion of D-genes to that of the control (all genes) by a G-test of
150 Goodness-of-Fit (P -value < 0.000121). However, the use of this approach with the D-
151 genes did not yield a coherent introgressive state, thus providing little insight into the
152 history of Asian rice at the present stage, emphasizing the need for a more systematic

153 approach to decipher the genome-wide status of Asian rice. For a further interpretation of
154 these results, see **Discussion**.

155

156 **Consequence of Analysis Window Size**

157 Because the introgressive states of D-genes did not give clear answer to the history of
158 Asian rice, we consequently explored the genome-wide introgressive states in a manner
159 involving significantly more computational resource costs and time.

160 Our phylogenetic analysis for one of the D-genes (*LGI*) with variable lengths of
161 flanking upstream/downstream regions (**Fig. 2e** : CDS only, **f** : +5kb-upstream/+5kb-
162 downstream, **g** : +10kb-upstream/+10kb-downstream, **h** : +20kb-upstream/+20kb-
163 downstream, and **i** : +100kb-upstream/+100kb-downstream, respectively) clearly shows
164 that region size heavily affects the resultant phylogeny. More precisely, a narrow region
165 (CDS only) showed a monophyletic topology of *LGI* between *indica* and *japonica*,
166 suggesting that it is introgressive (**Fig. 2e**), while wider region analyses resulted in a
167 polyphyletic relationship resembling non-introgressive state (**Fig. 2g, h, and i**). Full-size
168 tree pictures with a detailed color system are shown in **Supplementary Fig. 2**. Therefore,
169 we emphasize that window size is important; the window size setup in genome-wide
170 analysis is significant when we are dealing with phylogenies of domesticated Asian rice
171 at the loci-level.

172 The genome of domesticated Asian rice is polyphyletic as a whole, yet not always so
173 at the loci-level^{7,9,14,20-22}. This is in line with our inconsistent result (**Fig. 2e, f, g, h, and**
174 **i**), indicating that a narrower window setup leads to a more accurate inference of
175 phylogeny at the loci-level. Moreover, adopting a wider window size is inaccurate

176 because it does not deal with phylogenies at the loci-level^{7,9,21,22}, but rather with a whole-
177 genome phylogeny. Furthermore, our preliminary analyses with imputed 4,587 accession
178 genotypes unsuccessfully resulted in similar inconsistent phylogenetic relationships,
179 indicating that methods based on the haplotype linkages in certain-sized regions (*e.g.*,
180 wider window size; imputation) are not suitable for exploring the phylogenies at the loci-
181 level.

182

183 **Genome-wide Introgressive States Occur in Blocks**

184 We developed a machine learning classification model to distinguish the non-introgressed
185 windows (**Fig. 2a**) from introgressed windows (**Fig. 2b**) computationally. This is to
186 streamline a time-consuming visual inspection (*e.g.*, if we set 1kb windows all along the
187 rice genome (~373Mb), we would need to handle ~373,000 windows). Another merit for
188 adopting a machine learning-aided method is that it is free from null hypotheses and *P*-
189 value-dependent approach²⁹. As shown in **Methods**, we achieved 96.1% accuracy for the
190 binary classifier by the Breiman & Cutler's Random Forest Algorithm³⁰, and thus we
191 adopted it for further analyses.

192 Initially, we scanned the rice genome and developed an *indica* - *japonica* IR map at
193 100kb-resolution using a random forest classification model (for details, see **Methods**),
194 but it was blocky and the introgressive landscape was still veiled, shown in **Fig. 3a**
195 showing chromosome 1. We then increased the resolution to 20kb- (**Fig. 3b**), 10kb- (**Fig.**
196 **3c**), 5kb- (**Fig. 3d**), and finally to 1kb (**Fig. 3e**). The 1kb-resolution IR map produced a
197 sharp image that discriminate introgressive states at the gene loci-level along the entire
198 genome (IR maps for chromosome 2 to chromosome 12 are shown in **Supplementary**

199 **Fig. 3**). We identified large amounts of IR bands all along the genome (**Fig. 3e** and
200 **Supplementary Fig. 3**). Surprisingly, we determined that 14.2% of genomic contents are
201 introgressive (**Fig. 4a**). In addition, the IRs are not uniformly distributed, but rather
202 unevenly located in blocks (**Fig. 3e** and **Supplementary Fig. 3**). To be precise, there are
203 several major wide IRs in each chromosome, while thousands of narrow IRs are scattered
204 all over the genome (**Fig. 3e** and **Supplementary Fig. 3**), suggesting that there are
205 multiple genetic backgrounds behind the introgressions.

206

207 **Non-uniform Ages of Introgressions**

208 Now that we have established that a substantial amount (14.2%) of the genetic contents
209 has been exchanged between *indica* and *japonica* subpopulations, we aimed to uncover
210 what the biased introgressive pattern (**Fig. 3e** and **Supplementary Fig. 3**) means. By
211 plotting the window proportions of particular *DD*s, we observed apparent non-uniform
212 genetic backgrounds (**Fig. 4b**). We propose that these multiple genetic backgrounds
213 correspond to multiple classes of IRs, and that wide IRs and narrow IRs have different
214 *DD* values. To test our proposal, we operationally and precisely defined two IR classes
215 according to the dimensional continuity of IR windows, with wide IRs ($\geq 40\text{kb}$) and
216 narrow IRs ($=1\text{kb}$), and explored their *DD*s. The genomic positions of the wide IRs are
217 shown in **Supplementary Table 3**. The results show that wide IRs have a small *DD* of
218 5.89×10^{-6} substitutions/site, on average for all chromosomes, and narrow IRs have
219 roughly 100 times larger *DD* than wide IRs (5.84×10^{-4} substitutions/site). Non-IRs show
220 a much larger *DD* (1.71×10^{-3} substitutions/site) (**Fig. 4a** shows the average for all
221 chromosomes; results for each chromosome are shown in **Supplementary Table 4**). This

222 similar trend of *DD* can also be observed in the continuous-valued histogram (continuity
223 of IR windows; from one-IR to 15-IRs) shown in **Supplementary Fig. 4**.

224 When we roughly extrapolate the *indica-japonica* divergence time to 500,000 years
225 ago^{7,26} (**Fig. 5**, non-IRs), we can then estimate that the wide IRs are approximately 1,700
226 years old, whereas the narrow IRs are approximately 170,000 years old (**Fig. 5**). Hence,
227 we concluded that the wide IRs are relatively recently formed, while the narrow IRs have
228 existed for considerably longer time.

229

230 **Correspondence between Wide IRs and Selective Sweep Regions**

231 To gain insight into the history of the domestication of Asian rice and to address the
232 controversy on the origins of this domestication, we compare the genomic locations of
233 our IRs with those of previously reported domestication-associated genomic entities,
234 namely; Selective Sweep Regions (SSRs)¹⁴ and Co-located Low-Density Genomic
235 Regions (CLDGRs)¹⁰. We re-computed these previously described SSRs and
236 CLDGRs^{10,14} with our 4,587 rice accessions dataset (**Fig. 1a**) onto the Os-Nipponbare-
237 Reference-IRGSP-1.0 reference genome (see **Methods** for more details), as shown in
238 parallel with our IRs in **Fig. 3e, f, and g** and **Supplementary Fig. 3** (red lines: SSRs, blue
239 lines: CLDGRs). Interestingly, our results show that the SSRs correspond well with our
240 IRs, in particular with wide IRs (*i.e.*, young IRs), suggesting that the SSRs capture
241 recently happened events of introgression. However, in contrast, we observed less
242 correspondence between the CLDGRs and our wide IRs (**Fig. 3e, f and g** and
243 **Supplementary Fig. 3**), suggesting that CLDGRs do not deal with such events of
244 introgression. We discuss these patterns of correspondence further in **Discussion**.

245

246 **Discussion**

247 The genetic structure of domesticated Asian rice includes five major subpopulations ³¹. A
248 recently study shows that it can be subdivided into nine detailed subpopulations ²².
249 Ancient Chinese literature reported as early as the Han dynasty in China (100 AD) the
250 existence of two ecogeographical rice groups called ‘*Xian* (or *Hsien*)’ and ‘*Geng* (or
251 *Keng*)’, which correspond to *indica* and *japonica* subpopulations, respectively ^{32,33}. This
252 indicates that *indica* and *japonica* subpopulations have been cultivated for at least around
253 2000 years, being exposed to human intervention for a long time. For this reason, we
254 chose these two subspecies as the best model for studying the domestication of Asian rice.
255 In addition, we considered these subspecies because of the availability of high quality
256 sequenced genomes ³⁴, curated genome annotations ³⁵, more than 3,000 re-sequenced
257 closely-related accessions ²²⁻²⁵, and additional quality reference genomes (IR8 for *indica*
258 and N22 for *aus*), together with eight wild *Oryza* species ²⁶.

259 Archeological evidence indicates that Asian rice was first domesticated in the early
260 Holocene period ca. 9000 ^{5,36}, but Asian rice domestication and its origin is still a matter
261 of ongoing debate in both archeological and genetic research areas ⁵⁻²⁰. Plant scientists
262 have expected that the availability of whole-genome sequences of domesticated Asian
263 rice, its wild relatives, and ancient rice ³⁷, would provide a resolution to this long-
264 standing debate, yet the controversy is ongoing, because the genetic structure of rice
265 genomes turned out to be more complex than expected. In the two research studies of
266 evolutionary origins of domesticated Asian rice ^{10,14}, they analyzed a single dataset,
267 which included 1,529 genotypes of wild and domesticated rice ^{14,38}, leading to opposite

268 domestication scenarios. More recently, the same dataset was re-evaluated by the third
269 team, who suggested that rice originated from multiple populations of *O. rufipogon*
270 (and/or *O. nivara*): *De novo* domestication only occurred once where domestication
271 alleles were introgressed predominantly from *japonica* into *indica* subpopulations^{7,8}.

272 In this study, we explore possible events of introgression between subspecies,
273 considering them as traceable signs of domestication (**Fig. 2a** and **b**). We capture the
274 genome-wide IR map between *O. sativa* ssp. *indica* and *japonica*, with the aim of
275 encapsulating the entire history of Asian rice domestication. We exhaustively scan and
276 reveal the genome-wide introgressive landscape between *indica* and *japonica* at the finest
277 resolution using a machine learning classification model (**Fig. 3e** and **Supplementary**
278 **Fig. 3**). Our results show that a surprisingly large proportion of the rice genome (14.2%)
279 consists of wide and narrow traces of introgression between *indica* and *japonica* (**Fig. 4a**).
280 This suggests that even after the initial diversification of Asian rice roughly 500,000
281 years ago^{7,26}, *indica* and *japonica* subpopulations have been exchanging alleles between
282 each other.

283 In addition, we explore the introgressive state of 25 D-gene regions. We detected a
284 significantly large number of D-genes upon IRs, though not all of D-genes
285 (**Supplementary Table 2**), which shows that introgression was a major but non-exclusive
286 molecular mechanism for D-gene propagation. In other words, some D-genes moved
287 along the introgressive flows (regardless of the direction). Note that not all D-genes were
288 mobilized via introgression events.

289 We also observed that, in terms of *DD*, the wide IRs have emerged recently, whereas
290 the narrow IRs have existed for a much longer time (**Fig. 4a** and **Supplementary Table**

291 4). This mosaic introgressive landscape in terms of time (**Fig. 5**) clearly indicates that
292 multiple introgression events between subpopulations have taken place multiple times
293 throughout history (**Fig. 6**). In each of these events, the brand-new wide IRs would
294 comprise some beneficial alleles and many non-beneficial alleles. The beneficial alleles
295 would have been selected for and fixed in recipient subpopulations, while the non-
296 beneficial alleles would not have been fixed in the subpopulation. Thus, the genomic
297 regions with less advantageous alleles would have been replaced, eventually disappearing
298 following subsequent multiple backcrosses within the recipient subpopulation (**Fig. 6**).
299 Such genome dynamics can look like “sequentially built sandcastles” on a beach,
300 whereby newly built castles are still intact, while the older castles are already beginning
301 to crumble (**Fig. 6**). From the standpoint of our Sandcastles Model, the vast majority of
302 detected IRs correspond to non-beneficial alleles, which are mostly derived by
303 hitchhiking effects (**Fig. 6**). Extrapolating the *indica-japonica* divergence time (500,000
304 years ago corresponds to 1.71×10^{-3} substitutions/site in terms of DD)^{7,26}, we can
305 estimate that the narrow and wide IRs are approximately 170,000 and 1,700 years old,
306 respectively (**Fig. 5**). This is consistent with the Asian rice domestication timeline: It was
307 initially domesticated in the early Holocene period^{5,36} and has been maintained for at
308 least about 2,000 years^{32,33}.

309 The history, particularly the first origins of Asian rice domestication has long been a
310 subject of active discussion in plant biology⁵⁻²⁰. Studies have focused specifically on the
311 domestication-associated regions that presumably reflect the domestication process in
312 rice genomes. Those regions are typically defined by D-gene loci with flanking
313 upstream/downstream regions, SSRs, and CLDGRs. As an inevitable consequence in

314 those studies ^{10,14}, the definition of domestication-associated regions heavily affected the
315 reconstructed genetic phylogenies and the conclusions.

316 In this study, by employing highly dense SNP information and a machine learning
317 modeling approach, we elucidated a 1kb-resolution IR map and found that the young IRs
318 were well co-localized with SSRs ¹⁴, but not with CLDGRs ¹⁰. In terms of population
319 genetics, each of the IRs and SSRs were derived from a different population statistic, *i.e.*,
320 IRs were detected by a decrease in genetic distance difference to the wild relative (*DD*),
321 while SSRs were inferred by a decrease in nucleotide diversity (Π) compared to that of
322 the wild relatives. However, since gene introgressions will act in the direction of
323 decreasing Π in the domesticated population, $\Pi(\text{wild}) / \Pi(\text{domesticated})$ will have a
324 higher value, and thus the correspondence between SSRs and young IRs makes sense. In
325 terms of molecular phylogeny, the young IRs show a quite higher genetic identity
326 between *indica* and *japonica*, which could lead to monophyly (**Fig. 5**, bottom right panel).
327 On the other hand, the old IRs and non-IRs tend to represent more genetic divergence,
328 which seems to be polyphyletic (**Fig. 5**, bottom left panel and top panel). Hence the
329 discrepancy in results from the two previous studies ^{10,11,14,15} can be reasonably explained
330 by our Sandcastles Model (**Fig. 6**), *i.e.*, one study focused on the new castles (young IRs)
331 ¹⁴, while the other did not ¹⁰. We propose that focusing on certain-sized genomic regions
332 is a misleading way to understand the primal origins of domesticated life, because these
333 regions may contain recently built young IR blocks (**Fig. 6**).

334 If we pursue the very first ancestor of domesticated Asian rice, we need to eliminate
335 carefully the SSR-like entities that overlap with the young IR blocks from the analysis,
336 because they are recent and do not reflect ancient domestication history. We should

337 instead probe into other SSRs (old SSRs) and/or old IRs in the genome, which are the
338 true traces of ancient domestication history. However, since the extant domesticated
339 subspecies (*e.g.*, *indica* and *japonica*), and closely-related wild relatives (*O. nivara* and *O.*
340 *rufipogon*) as well, are not yet completely isolated reproductively^{39,40}, the subspecies
341 boundary of initial Asian rice subpopulations should be much more permeable. In that
342 sense, it may not be meaningful to explore whether the initial domesticated rice
343 individual(s) is in a single subpopulation or are in multiple subpopulations, because they
344 were too permeable in terms of conventional taxonomy. Leastwise, be that as the initial
345 domestication might happen multiple times, these were in a single population *sensu lato*.

346 In summary, we have determined that a surprisingly large proportion (14.2%) of
347 genetic contents has been exchanged between *indica* and *japonica* subpopulations. We
348 have also demonstrated that introgression events have happened in multiple genomic
349 regions over multiple periods throughout the history of domesticated Asian rice, revealing
350 the complex spatiotemporal genome dynamics in Asian rice domestication.
351 Concomitantly, we settle the major controversy in plant science between two hypotheses
352⁵⁻²⁰ using our Sandcastles Model, *i.e.*, each study was focusing on a different genomic
353 region of a different era. Moreover, because the IRs contribute to the domestication
354 process in a proactive manner, our IR map provides a unique reference for potential target
355 loci in breeding of rice. This gives insight into new breeding designs and practices based
356 on the introgressive genomic map.

357

358 **Methods summary**

359 The genotypes of domesticated and wild rice accessions were all retrieved from publicly

360 available databases. The full methods and any associated information are available in the
361 online version of the paper.

362

363 **Methods**

364 **Reference genome.** For the reference genome sequences and reference genome
365 annotations, the reference Nipponbare genome Os-Nipponbare-Reference-IRGSP-1.0 (*O.*
366 *sativa* ssp. *japonica* cv. Nipponbare) ⁴¹ ; hereinafter referred to as Nipponbare RefSeq
367 and CGSNL annotations served in RAP-DB ⁴² were employed, respectively.

368 **Domestication-associated genes (D-genes).** Based on our literature survey, we manually
369 selected and curated a total of 25 D-genes (**Fig. 2c**) for this study. The selection criteria
370 were based on agronomically beneficial effects of genes selected.

371 **Issues on rice genotypes.** In particular, we focused our analyses on two *O. sativa*
372 subspecies, ssp. *indica* and ssp. *japonica*, as an Asian rice domestication model. Despite
373 multiple studies conducted to explore the history of Asian rice introgression and
374 domestication with large-scale accessions datasets including *indica* and *japonica* ⁸
375 ^{10,14,21,22} , their genome-wide scanning procedures have been performed using relatively
376 large window size setups (5kb -100kb). The importance of window size in such analyses
377 are outlined in this study (**Fig. 2e, f, g, h, and i**) and also in Choi & Purugganan ⁸ , but due
378 to the low SNPs density (56.4% missing data rate) in the dataset ^{14,38} , the issue of
379 window size had not yet been overcome. Another problem is that each *indica* and
380 *japonica* subpopulation contains a significant amount of genetic diversity ^{14,22,31} , or
381 rather, some subspecies accessions can be intermediate accessions between the two
382 subspecies since these subpopulations are not yet completely reproductively isolated from

383 each other³⁹. In fact, both *indica* and *japonica* subpopulations show a certain degree of
384 phenotypic diversity, including some intermediate traits (**Fig. 1c**). Consequently, when
385 taking all the *indica* and *japonica* accessions into account, the conclusion may be
386 ambiguous because of the intermediate states of genetic distance. The final issue to be
387 overcome when we trace back the domestication history of Asian rice is to choose which
388 species to use as an outgroup. It is widely believed that *O. nivara* and *O. rufipogon* are
389 the immediate ancestors of ssp. *indica* and ssp. *japonica*, respectively². However, those
390 wild rice species are still able to intermate with *O. sativa*⁴⁰; thus, the genetic distance
391 between those wild rice species and *O. sativa* could be underestimated in introgressive
392 regions. Hence, those wild rice species are not always suitable for outgroup species in
393 phylogenetic analysis. Our preliminary gene-by-gene phylogenetic analyses with the
394 3,000 Rice Genomes Project²²⁻²⁵, higher coverage wilds^{26,38,43,44} and the *O. punctata*²⁶
395 datasets (**Fig. 1a**, in total 3,060 accessions) aimed to assess the suitability of *O. nivara*,
396 *O. rufipogon*, *O. glaberrima*, *O. barthii*, *O. glumaepatula* and *O. punctata* as outgroup
397 species for this study (**Supplementary Fig. 5**). Our analyses showed that in some cases
398 (e.g. *Gn1a*, *LGI1*, *Phr1*, and *qSH1*) (**Supplementary Fig. 5i, n, o and q**), a close-relatives
399 (*O. rufipogon* or *O. nivara*) can serve as an outgroup species. However, in most cases,
400 they are not suitable for an outgroup since they are not genetically isolated from
401 domesticated rice (**Supplementary Fig. 5**).

402 **Solutions on rice genotype issues.** To develop an accurate high-resolution (up to 1kb
403 window width) map of Asian rice introgression in a reasonable manner, we needed to
404 address the above-mentioned three problems: i) the low density of rice genotypes, ii)
405 over-diversity within each subspecies, and iii) the instability of outgroup. With the aim of

406 achieving good quality and quantity of rice genotypes, we collected imputation-free ~14x
407 coverage genotypes of 3,024 rice cultivars (**Fig. 1a**) from the 3,000 Rice Genomes
408 Project ²²⁻²⁵, in conjunction with other publicly available genotypes (**Fig. 1a**). We
409 appropriately converted their genomic coordinates to that of the Nipponbare RefSeq as
410 described ³⁸ when needed. We performed genomic imputation with the Beagle program ⁴⁵
411 in two batches (wild/domesticated) separately and exclusively on the 4,553 accessions
412 only for the purpose of SSRs and CLDGRs re-computation (**Fig. 1a**), but not on any
413 other accession datasets. The core dataset (**Fig. 1a**, 3,025 accessions) contained 1,712
414 *indica* and 833 *japonica* accessions with a missing genotype rate of 15.0% on average.
415 Then, to overcome the effect of intra-subspecies divergence, we dynamically picked up
416 median 10th accessions from *indica* and *japonica* window by window (see **Introgressive**
417 **Regions (IRs) detection**). Finally, to adopt an appropriate outgroup species in our study,
418 based on preliminary gene-by-gene phylogenetic analyses (**Supplementary Fig. 5**), we
419 exclusively employed the *O. punctata* (IRGC105690, BB diploid, 2n=24, geographical
420 origin: Africa) ²⁶ only, with the assumption that it has been mostly reproductively isolated
421 from *O. sativa* populations. We can ignore the underestimate effect of nucleotide distance
422 due to possible introgression events between *O. sativa* and *O. punctata* (**Supplementary**
423 **Fig. 5**).

424 **Mapping and SNPs calling.** We first quality inspected all short reads by FastQC
425 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>), and then we filtered out
426 and/or trimmed out adaptor sequences and low-quality bases using Trimmomatic ⁴⁶. After
427 those preprocessing steps, we mapped the remaining reads onto the Nipponbare RefSeq
428 using 'bwa mem' commands in BWA ⁴⁷ with default parameters, except for the proper

429 insert size limitation (-w 500 or -w 800, dictated by the data source). Repeat
430 sequences scattered within the Nipponbare RefSeq were not masked in our mapping
431 process. Next, we called variants using the GATK⁴⁸ with a conventional best practice
432 method (<https://software.broadinstitute.org/gatk/best-practices/>).

433 **Phylogenetic tree construction.** For window-base analysis, we generated each 1,000bp
434 multiple alignment. For gene-by-gene analysis, we generated a multiple alignment of
435 actual CDS for each gene (including intron regions, but not including any flanking
436 upstream/downstream regions). All nucleotide genetic distances between domesticated
437 rice windows/genes and outgroup windows/genes were estimated by PHYLIP-dnadist
438 command with default parameters (Felsenstein84 distance)²⁸. We reconstructed all
439 phylogenetic trees using the PHYLIP-neighbor command with default parameters
440 (Neighbor-Joining method)^{28,49}. Trees were drawn by FigTree software GUI
441 (<http://tree.bio.ed.ac.uk/software/figtree/>), rooted by *O. punctata* as the fixed outgroup.

442 **Invention of *Distance Difference (DD)*.** Under isolated conditions, each of *indica* and
443 *japonica* subpopulations should show different genetic distances to an outgroup (a wild
444 rice accession) to some extent, since they have been separated for a length of time in each
445 subpopulation (**Fig. 2a**). However, they will show unexpectedly similar genetic distance
446 to an outgroup when an inter-subspecies cross (*i.e.* introgression) has occurred recently
447 (**Fig. 2b**). Together with incomplete lineage sorting and other possible situations^{50,51}, this
448 is one of the reasons why a particular gene phylogeny does not always agree with the
449 (sub)species phylogeny. Here we conceptually define *DD* (genetic *Distance Difference* to
450 the outgroup) as;

451
$$DD = |F84 (\text{outgroup to } indica) - F84 (\text{outgroup to } japonica)| .$$

452 (*) F84 = Felsenstein84 nucleotide genetic distance ²⁸

453 Here, smaller *DDs* represent IRs, while larger *DDs* mean that those are non-IRs. Note
454 that IRs happened in the initial period of domestication will not show enough decrease in
455 *DD*, hence such IRs are out of scope of this method. In terms of population genetics, we
456 have multiple *indica* accessions and multiple *japonica* accessions, and each
457 subpopulation includes much genetic diversity (see **Issues on rice genotypes**). To
458 overcome the undesirable effect on intra-subspecies over-diversity in terms of nucleotide
459 distance to the outgroup, we dynamically chose the median 10th accessions from *indica*
460 (172 accessions) window by window (or gene by gene), and median 10th accessions from
461 *japonica* (84 accessions) window by window (or gene by gene), respectively. They are
462 representative subpopulations in each window (or each gene) in the sense that the most
463 mediocre members reflect the profile of population. Therefore, the actual *DD* value is not
464 computed by a single *indica* accession and a single *japonica* accession. Instead, it is
465 computed by the average form of median 10th accessions of *indica*, and by the average
466 form of median 10th accessions of *japonica*. Hence, the actual formula for *DD* is;

$$467 \quad DD = \left| \frac{\sum_{indica}^{median\ 10th} F84(outgroup\ to\ indica)}{172} - \frac{\sum_{japonica}^{median\ 10th} F84(outgroup\ to\ japonica)}{84} \right|.$$

468 (*) F84 = Felsenstein84 nucleotide genetic distance ²⁸

469 **Introgressive Regions (IRs) detection.** For the gene-by-gene analysis, we conducted
470 visual phylogeny inspection (**Fig. 2** and **Supplementary Fig. 1**). For the window-based
471 analysis, although visual inspection of each window phylogeny would give the best
472 accuracy, it is too time consuming. We thus aimed to computationally distinguish the
473 non-introgressed windows (**Fig. 2a**) from the introgressed windows (**Fig. 2b**) by the use
474 of a binary classifier through Breiman & Cutler's Random Forest Algorithm ³⁰. The

475 accuracy of the binary classifier was 96.1%, as determined by a 10-fold cross validation
476 (for more details, see **Optimization of machine learning models**). The 1kb resolution
477 machine learning classification result showed that 14.2% of the rice genome was
478 introgressive, and 50.0% was non-introgressive (was excluded 35.8% from the analysis
479 and marked as status-undetermined, for reasons outlined below) (**Fig. 4a**). In the
480 window-based analysis, we excluded windows that have less alignable length with the
481 outgroup (<5% of the window region, *i.e.* <50bp in the case of the 1kb window setup).
482 We also excluded windows with no genetic difference (*i.e.*, no SNP) from any of the
483 *indica/japonica* accessions to the outgroup at all. Those windows are shown as gray
484 windows (**Fig. 3** and **Supplementary Fig. 3**).

485 **Training of machine learning models.** For the training dataset of machine learning
486 classification models, we firstly conducted visual phylogeny inspection for randomly
487 chosen 640 1kb-windows (~0.267% of total phylogeny determined windows, see **Fig.**
488 **4a**), and we identified 114 windows as IRs and 526 windows as non-IRs. We then
489 balanced the ratio between positive cases (IRs) and negative cases (non-IRs) in 114 IRs
490 and randomly sub-sampled 114 non-IRs, respectively, and these 228 cases were finally
491 used as the actual training dataset for generating the classification models.

492 **Optimization of machine learning models.** For the features used to develop the
493 classification models, we extracted the nucleotide distance matrices for median 10th 257
494 accessions (172 *indica*, 84 *japonica*, and 1 outgroup). Since the $257^2 = 66,049$ variables
495 were too computationally demanding, we reduced the variables by equal subsampling to
496 50 accessions, retaining the original variations in each subspecies (50 *indica*, 50
497 *japonica*, and 1 outgroup). Finally, we adopted $101^2 = 10,201$ variables as the features

498 for developing the classification models. In order to find the best option for our machine
499 learning analysis, then we conducted a grid search for model parameters with a support
500 vector machine model (with non-linear Gaussian kernel) (with parameters $C = 2, 4, 8, 16,$
501 $32, 64, 128, 256, 512, 1024$; $\sigma = 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024$; 100 cases in
502 total), and a random forest model (with parameters $n\text{tree} = 16, 32, 64, 128, 256, 512,$
503 $1024, 2048, 4096, 8192$; $m\text{try} = 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024$; 100 cases in
504 total). We determined that the random forest model ($n\text{tree} = 512, m\text{try} = 256, \text{accuracy} =$
505 96.1% by 10-fold cross validation, data not shown) was the best option. We implemented
506 the support vector machine model, random forest model, and cross validation framework
507 by R language and R packages (kernlab, randomForest, and mlr) ([https://www.r-](https://www.r-project.org)
508 [project.org](https://www.r-project.org)).

509 **Verification of the machine learning model.** To verify the effectiveness of our random
510 forest classifier, we drew an identical conclusion by adopting another statistical
511 classification method as shown below. Assuming that the median 10th subset data are not
512 normally distributed, we tested whether the difference between F84 (outgroup to *indica*)
513 and F84 (outgroup to *japonica*) is statistically significant or not, using the non-parametric
514 statistical test method (Mann-Whitney U test, $P\text{-value} < 10^{-7}$), window by window. When
515 the null hypothesis is rejected, the window will be non-introgressive (**Fig. 2a**,
516 significantly different). Otherwise (*i.e.*, not significantly different), it is considered a
517 candidate for introgression (**Fig. 2b**). As noted above, although the P -value threshold is
518 quite conservative ($P\text{-value} < 10^{-7}$), 54.8% of the rice genome (similarly to random forest
519 model at 50.0%) was still determined as significant (*i.e.*, non-introgressive). We
520 determined that genomic locations were introgressive similarly to the random forest

521 model (data not shown), and our conclusion was identical to that of the random forest
522 model. Even if we adopted a more aggressive P -value < 0.05 , the significant (*i.e.*, non-
523 introgressive) window percentages were still quite similar (56.4%), the genomic locations
524 as introgressive were still similar to those of the random forest model (data not shown)
525 and again we reached identical conclusions, thus demonstrating the robustness of our
526 random forest model. Moreover, manual phylogeny curation of 25 gene-by-gene results
527 was well in line with the window-based results of random forest (**Fig. 3** and
528 **Supplementary Fig. 3**), reconfirming the accuracy of our random forest model.

529 **Enrichment test for D-genes on IRs.** We tested whether the 25 D-genes (**Fig. 2c**) are
530 statistically significantly enriched (or depleted) on IRs or not. A G-test of Goodness-of-
531 Fit showed statistically significant enrichment on the proportion of introgressive D-genes
532 (9 genes) against non-introgressive D-genes (14 genes) (**Supplementary Table 2**) (2 D-
533 genes (*Hdl* and *S5*) showed undetermined phylogeny). For the control (all genes, *i.e.*,
534 expected proportion), we computationally determined each gene's IRs concordance when
535 the entire gene locus was inclusively contained in any continuous IRs of 1kb resolution
536 (introgressive = 3,498 genes: 9.24%; non-introgressive = 34,350 genes: 90.8%). The G-
537 test was conducted with the following R script:

```
538 > observed      = c(9,14)
539 > expected.prop = c(0.0924, 0.908)
540 > degrees = 1
541 > expected.count = sum(observed)*expected.prop
542 > G = 2 * sum(observed * log(observed / expected.count))
543 > G
544 [1] 14.78253
545 > pchisq(G,df=degrees,lower.tail=FALSE)
546 [1] 0.0001206482
547 > q()
548
```

549 **Re-computation of Selective Sweep Regions (SSRs) and Co-located Low-Density**
550 **Genomic Regions (CLDGRs).** For the already reported domestication-associated

551 genomic entities (Selective Sweep Regions (SSRs)¹⁴ and Co-located Low-Density
552 Genomic Regions (CLDGRs)¹⁰), we re-computed their SSRs and CLDGRs using our
553 4,587 accessions dataset (**Fig. 1a**) on the Nipponbare RefSeq, and we conducted
554 independent permutation tests to determine the appropriate $\Pi(\text{wild}) / \Pi(\text{domesticated})$
555 threshold. In **Fig. 3e** and **Supplementary Fig. 3**, re-computed SSRs and CLDGRs were
556 shown as red lines and blue lines, respectively. The re-computation procedures are
557 summarized in **Supplementary Fig. 6** and **7**.

558 **Data availability.** All the intermediate and final analysis results in this study are
559 available from the corresponding author upon request.

560

561 **D-genes' References (will be imported to Fig. 2c):**

*BADH2*⁵²

*Bh4*⁵³

*Bph14*⁵⁴

*C1*⁵⁵

*DPL2*⁵⁶

*Ehd1*⁵⁷

*GAD1*⁵⁸

*Ghd7*⁵⁹

*Gn1a*⁶⁰

*GS3*⁶¹

*GW2*⁶²

*Hd1*⁶³

*LABA1*⁶⁴

*LG1*⁶⁵

*Phr1*⁶⁶

*Prog1*⁶⁷

*qSH1*⁶⁸

*qSW5*⁶⁹

*Rc*⁷⁰

*Rd*⁷¹

*S5*⁷²

*sd1*⁷³
*sh4*⁷⁴
*tb1*⁷⁵
*waxy*⁷⁶

562

563 **References:**

- 564 1. FAO. FAO Statistical Yearbook Part3 : Feeding the world. (2013).
565 2. Kumagai, M., Tanaka, T., Ohyanagi, H., Hsing, Y.C. & Itoh, T. Genome Sequence of Oryza
566 Species. in *Rice Genomics, Genetics and Breeding* (eds. Sasaki, T. & Ashikari, M.) 1-20
567 (Springer, 2018).
568 3. Wang, M. *et al.* The genome sequence of African rice (*Oryza glaberrima*) and evidence
569 for independent domestication. *Nat Genet* **46**, 982-8 (2014).
570 4. Hilbert, L. *et al.* Evidence for mid-Holocene rice domestication in the Americas. *Nat Ecol*
571 *Evol* **1**, 1693-1698 (2017).
572 5. Callaway, E. Domestication: The birth of rice. *Nature* **514**, S58-9 (2014).
573 6. Carpentier, M.C. *et al.* Retrotranspositional landscape of Asian rice revealed by 3000
574 genomes. *Nat Commun* **10**, 24 (2019).
575 7. Choi, J.Y. *et al.* The Rice Paradox: Multiple Origins but Single Domestication in Asian
576 Rice. *Mol Biol Evol* **34**, 969-979 (2017).
577 8. Choi, J.Y. & Purugganan, M.D. Multiple Origin but Single Domestication Led to *Oryza*
578 *sativa*. *G3 (Bethesda)* **8**, 797-803 (2018).
579 9. Civan, P. & Brown, T.A. Role of genetic introgression during the evolution of cultivated
580 rice (*Oryza sativa* L.). *BMC Evol Biol* **18**, 57 (2018).
581 10. Civan, P., Craig, H., Cox, C.J. & Brown, T.A. Three geographically separate domestications
582 of Asian rice. *Nat Plants* **1**, 15164 (2015).
583 11. Civan, P., Craig, H., Cox, C.J. & Brown, T.A. Multiple domestications of Asian rice. *Nat*
584 *Plants* **2**, 16037 (2016).
585 12. Gao, L.Z. & Innan, H. Nonindependent domestication of the two rice subspecies, *Oryza*
586 *sativa* ssp. *indica* and ssp. *japonica*, demonstrated by multilocus microsatellites.
587 *Genetics* **179**, 965-76 (2008).
588 13. Gross, B.L. & Zhao, Z.J. Archaeological and genetic insights into the origins of
589 domesticated rice. *Proceedings of the National Academy of Sciences of the United States*
590 *of America* **111**, 6190-6197 (2014).
591 14. Huang, X. *et al.* A map of rice genome variation reveals the origin of cultivated rice.
592 *Nature* **490**, 497-501 (2012).
593 15. Huang, X.H. & Han, B. Rice domestication occurred through single origin and multiple
594 introgressions. *Nature Plants* **2**(2016).
595 16. Londo, J.P., Chiang, Y.C., Hung, K.H., Chiang, T.Y. & Schaal, B.A. Phylogeography of Asian
596 wild rice, *Oryza rufipogon*, reveals multiple independent domestications of cultivated
597 rice, *Oryza sativa*. *Proc Natl Acad Sci U S A* **103**, 9578-83 (2006).
598 17. Molina, J. *et al.* Molecular evidence for a single evolutionary origin of domesticated rice.
599 *Proc Natl Acad Sci U S A* **108**, 8351-6 (2011).
600 18. Sang, T. & Ge, S. Genetics and phylogenetics of rice domestication. *Curr Opin Genet Dev*
601 **17**, 533-8 (2007).

- 602 19. Vitte, C., Ishii, T., Lamy, F., Brar, D. & Panaud, O. Genomic paleontology provides
603 evidence for two distinct origins of Asian rice (*Oryza sativa* L.). *Mol Genet Genomics* **272**,
604 504-11 (2004).
- 605 20. Yang, C.C. *et al.* Independent domestication of Asian rice followed by gene flow from
606 japonica to indica. *Mol Biol Evol* **29**, 1471-9 (2012).
- 607 21. Santos, J.D. *et al.* Fine scale genomic signals of admixture and alien introgression among
608 Asian rice landraces. *Genome Biol Evol* (2019).
- 609 22. Wang, W. *et al.* Genomic variation in 3,010 diverse accessions of Asian cultivated rice.
610 *Nature* **557**, 43-49 (2018).
- 611 23. Alexandrov, N. *et al.* SNP-Seek database of SNPs derived from 3000 rice genomes.
612 *Nucleic Acids Res* **43**, D1023-7 (2015).
- 613 24. Li, J.Y., Wang, J. & Zeigler, R.S. The 3,000 rice genomes project: new opportunities and
614 challenges for future rice research. *Gigascience* **3**, 8 (2014).
- 615 25. Mansueto, L. *et al.* Rice SNP-seek database update: new SNPs, indels, and queries.
616 *Nucleic Acids Res* **45**, D1075-D1081 (2017).
- 617 26. Stein, J.C. *et al.* Genomes of 13 domesticated and wild rice relatives highlight genetic
618 conservation, turnover and innovation across the genus *Oryza*. *Nat Genet* **50**, 285-296
619 (2018).
- 620 27. Sun, X., Jia, Q., Guo, Y., Zheng, X. & Liang, K. Whole-genome analysis revealed the
621 positively selected genes during the differentiation of indica and temperate japonica
622 rice. *PLoS One* **10**, e0119239 (2015).
- 623 28. Felsenstein, J. PHYLIP - Phylogeny Inference Package (Version 3.2). *Cladistics* **5**, 164-6
624 (1989).
- 625 29. Johnson, D.H. The insignificance of statistical significance testing. *Journal of Wildlife*
626 *Management* **63**, 763-772 (1999).
- 627 30. Breiman, L. Random forests. *Machine Learning* **45**, 5-32 (2001).
- 628 31. Garris, A.J., Tai, T.H., Coburn, J., Kresovich, S. & McCouch, S. Genetic structure and
629 diversity in *Oryza sativa* L. *Genetics* **169**, 1631-8 (2005).
- 630 32. Oka, H.I. *Origin of Cultivated Rice*, (Elsevier Science, Tokyo, 1988).
- 631 33. Ting, Y. Chronological studies of the cultivation and the distribution of rice varieties,
632 Keng and Sen. *Sun Yatsen University Agronomy Bulletin* **6**, 1-32 (1949).
- 633 34. International Rice Genome Sequencing, P. The map-based sequence of the rice genome.
634 *Nature* **436**, 793-800 (2005).
- 635 35. Ohyanagi, H. *et al.* The Rice Annotation Project Database (RAP-DB): hub for *Oryza sativa*
636 ssp. japonica genome information. *Nucleic Acids Res* **34**, D741-4 (2006).
- 637 36. Zuo, X.X. *et al.* Dating rice remains through phytolith carbon-14 study reveals
638 domestication at the beginning of the Holocene. *Proceedings of the National Academy*
639 *of Sciences of the United States of America* **114**, 6486-6491 (2017).
- 640 37. Kumagai, M. *et al.* Rice Varieties in Archaic East Asia: Reduction of Its Diversity from Past
641 to Present Times. *Mol Biol Evol* **33**, 2496-505 (2016).
- 642 38. Ohyanagi, H. *et al.* OryzaGenome: Genome Diversity Database of Wild *Oryza* Species.
643 *Plant Cell Physiol* **57**, e1 (2016).
- 644 39. Guo, J. *et al.* Overcoming inter-subspecific hybrid sterility in rice by developing indica-
645 compatible japonica lines. *Sci Rep* **6**, 26878 (2016).
- 646 40. Wang, H.R., Vieira, F.G., Crawford, J.E., Chu, C.C. & Nielsen, R. Asian wild rice is a hybrid
647 swarm with extensive gene flow and feralization from domesticated rice. *Genome*
648 *Research* **27**, 1029-1038 (2017).

- 649 41. Kawahara, Y. *et al.* Improvement of the *Oryza sativa* Nipponbare reference genome
650 using next generation sequence and optical map data. *Rice (N Y)* **6**, 4 (2013).
- 651 42. Sakai, H. *et al.* Rice Annotation Project Database (RAP-DB): an integrative and interactive
652 database for rice genomics. *Plant Cell Physiol* **54**, e6 (2013).
- 653 43. Xu, X. *et al.* Resequencing 50 accessions of cultivated and wild rice yields markers for
654 identifying agronomically important genes. *Nat Biotechnol* **30**, 105-11 (2011).
- 655 44. Zhao, Q. *et al.* Pan-genome analysis highlights the extent of genomic variation in
656 cultivated and wild rice. *Nat Genet* **50**, 278-284 (2018).
- 657 45. Browning, S.R. & Browning, B.L. Rapid and accurate haplotype phasing and missing-data
658 inference for whole-genome association studies by use of localized haplotype clustering.
659 *American Journal of Human Genetics* **81**, 1084-1097 (2007).
- 660 46. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
661 sequence data. *Bioinformatics* **30**, 2114-20 (2014).
- 662 47. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
663 transform. *Bioinformatics* **25**, 1754-60 (2009).
- 664 48. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing
665 next-generation DNA sequencing data. *Genome Res* **20**, 1297-303 (2010).
- 666 49. Saitou, N. & Nei, M. The neighbor-joining method: a new method for reconstructing
667 phylogenetic trees. *Mol Biol Evol* **4**, 406-25 (1987).
- 668 50. Pamilo, P. & Nei, M. Relationships between Gene Trees and Species Trees. *Molecular*
669 *Biology and Evolution* **5**, 568-583 (1988).
- 670 51. Yang, C.C., Sakai, H., Numa, H. & Itoh, T. Gene tree discordance of wild and cultivated
671 Asian rice deciphered by genome-wide sequence comparison. *Gene* **477**, 53-60 (2011).
- 672 52. Kovach, M.J., Calingacion, M.N., Fitzgerald, M.A. & McCouch, S.R. The origin and
673 evolution of fragrance in rice (*Oryza sativa* L.). *Proc Natl Acad Sci U S A* **106**, 14444-9
674 (2009).
- 675 53. Zhu, B.F. *et al.* Genetic control of a transition from black to straw-white seed hull in rice
676 domestication. *Plant Physiol* **155**, 1301-11 (2011).
- 677 54. Du, B. *et al.* Identification and characterization of Bph14, a gene conferring resistance to
678 brown planthopper in rice. *Proc Natl Acad Sci U S A* **106**, 22163-8 (2009).
- 679 55. Saitoh, K., Onishi, K., Mikami, I., Thidar, K. & Sano, Y. Allelic diversification at the C
680 (*Osc1*) locus of wild and cultivated rice: Nucleotide changes associated with
681 phenotypes. *Genetics* **168**, 997-1007 (2004).
- 682 56. Mizuta, Y., Harushima, Y. & Kurata, N. Rice pollen hybrid incompatibility caused by
683 reciprocal gene loss of duplicated genes. *Proc Natl Acad Sci U S A* **107**, 20417-22 (2010).
- 684 57. Doi, K. *et al.* Ehd1, a B-type response regulator in rice, confers short-day promotion of
685 flowering and controls FT-like gene expression independently of Hd1. *Genes Dev* **18**,
686 926-36 (2004).
- 687 58. Jin, J. *et al.* GAD1 Encodes a Secreted Peptide That Regulates Grain Number, Grain
688 Length, and Awn Development in Rice Domestication. *Plant Cell* **28**, 2453-2463 (2016).
- 689 59. Xue, W.Y. *et al.* Natural variation in Ghd7 is an important regulator of heading date and
690 yield potential in rice. *Nature Genetics* **40**, 761-767 (2008).
- 691 60. Ashikari, M. *et al.* Cytokinin oxidase regulates rice grain production. *Science* **309**, 741-5
692 (2005).
- 693 61. Fan, C. *et al.* GS3, a major QTL for grain length and weight and minor QTL for grain width
694 and thickness in rice, encodes a putative transmembrane protein. *Theor Appl Genet* **112**,
695 1164-71 (2006).

- 696 62. Song, X.J., Huang, W., Shi, M., Zhu, M.Z. & Lin, H.X. A QTL for rice grain width and weight
697 encodes a previously unknown RING-type E3 ubiquitin ligase. *Nat Genet* **39**, 623-30
698 (2007).
- 699 63. Yano, M. *et al.* Hd1, a major photoperiod sensitivity quantitative trait locus in rice, is
700 closely related to the arabidopsis flowering time gene CONSTANS. *Plant Cell* **12**, 2473-
701 2483 (2000).
- 702 64. Hua, L. *et al.* LABA1, a Domestication Gene Associated with Long, Barbed Awns in Wild
703 Rice. *Plant Cell* **27**, 1875-1888 (2015).
- 704 65. Zhu, Z.F. *et al.* Genetic control of inflorescence architecture during rice domestication.
705 *Nature Communications* **4**(2013).
- 706 66. Yu, Y.C. *et al.* Independent Losses of Function in a Polyphenol Oxidase in Rice:
707 Differentiation in Grain Discoloration between Subspecies and the Role of Positive
708 Selection under Domestication. *Plant Cell* **20**, 2946-2959 (2008).
- 709 67. Tan, L. *et al.* Control of a key transition from prostrate to erect growth in rice
710 domestication. *Nat Genet* **40**, 1360-4 (2008).
- 711 68. Konishi, S. *et al.* An SNP caused loss of seed shattering during rice domestication.
712 *Science* **312**, 1392-6 (2006).
- 713 69. Shomura, A. *et al.* Deletion in a gene associated with grain size increased yields during
714 rice domestication. *Nat Genet* **40**, 1023-8 (2008).
- 715 70. Sweeney, M.T., Thomson, M.J., Pfeil, B.E. & McCouch, S. Caught red-handed: Rc encodes
716 a basic helix-loop-helix protein conditioning red pericarp in rice. *Plant Cell* **18**, 283-94
717 (2006).
- 718 71. Furukawa, T. *et al.* The Rc and Rd genes are involved in proanthocyanidin synthesis in
719 rice pericarp. *Plant J* **49**, 91-102 (2007).
- 720 72. Du, H., Ouyang, Y., Zhang, C. & Zhang, Q. Complex evolution of S5, a major reproductive
721 barrier regulator, in the cultivated rice *Oryza sativa* and its wild relatives. *New Phytol*
722 **191**, 275-87 (2011).
- 723 73. Asano, K. *et al.* Artificial selection for a green revolution gene during japonica rice
724 domestication. *Proc Natl Acad Sci U S A* **108**, 11034-9 (2011).
- 725 74. Li, C., Zhou, A. & Sang, T. Rice domestication by reducing shattering. *Science* **311**, 1936-9
726 (2006).
- 727 75. Takeda, T. *et al.* The OsTB1 gene negatively regulates lateral branching in rice. *Plant J* **33**,
728 513-20 (2003).
- 729 76. Olsen, K.M. *et al.* Selection under domestication: Evidence for a sweep in the rice Waxy
730 genomic region. *Genetics* **173**, 975-983 (2006).

731

732 **Acknowledgements**

733 The research reported in this publication was supported through funding from King
734 Abdullah University of Science and Technology (KAUST), under award numbers
735 BAS/1/1059-01-01 (to T.G.), BAS/1/1606-01-01 (to V.B.B.), FCC/1/1976-03-01 (to T.G.)
736 and FCC/1/1976-20-01 (to T.G.).

737

738 **Author contributions**

739 H.O. designed the study, performed the bioinformatics and statistical analysis, and wrote
740 the manuscript. K.G. performed the bioinformatics analysis. S.N. wrote the manuscript
741 and contributed to insightful discussions. R.A.W., M.A.T., K.M. and V.B.B. edited the
742 manuscript and contributed to insightful discussions. K.L.M. provided easy access to the
743 genotypes and phenotypes of 3,000 Rice Genomes Project and contributed to insightful
744 discussions. T.G. designed the study and wrote the manuscript. All the authors discussed
745 the results and commented on the manuscript.

746

747 **Competing interests**

748 The authors declare no competing interests.

749

750 **Corresponding author**

751 Correspondence to Takashi Gojobori: takashi.gojobori@kaust.edu.sa

752

753 **Figure Legends**

754 **Fig. 1.** Passport data of domesticated and wild Asian rice accessions in this study (**a**, in
755 total 4,587 accessions. for more details in higher coverage wilds, see **Supplementary**
756 **Table 1**), and geographical origin of accessions in 3,000 Rice Genomes Project (**b**, 3,024
757 accessions). A typical phenotypic diversity within subspecies (**c**, grain length over grain
758 width in *O. sativa* ssp. *indica* (n=1269, green) and *japonica* (n=533, blue)).

759 **Fig. 2.** Schematic view of underestimate on genetic *Distance Difference* (**a** and **b**), and
760 phylogenetic analysis of manually curated D-genes (25 genes) and their determined
761 introgressive states (**c** and **d**). Under isolated conditions, each of *indica* and *japonica*
762 subpopulation shall show different genetic distance to the outgroup (a wild rice
763 accession) to some extent, since they have been isolated from each other for a length of
764 time (**a**), whereas they will show unexpectedly similar genetic distance to the outgroup
765 when they made an inter-subspecies crossing (*i.e.* introgression) recently (**b**). Manually
766 curated D-genes (25 genes) and their determined introgressive state (**c**). Reconstructed
767 phylogenetic trees of 25 D-genes (**d**), green nodes : *indica*, blue nodes : *japonica*. Non-
768 introgressive genes were shown in yellow background. Introgressive genes were shown
769 in red background. Genes of undetermined phylogeny were shown in gray background.
770 Phylogenetic trees for one of the D-genes (*LGI*) with variable length of flanking
771 upstream/downstream regions (**e** : CDS only, **f** : +5kb-upstream/+5kb-downstream, **g** :
772 +10kb-upstream/+10kb-downstream, **h** : +20kb-upstream/+20kb-downstream, and **i** :
773 +100kb-upstream/+100kb-downstream, respectively). Full size tree pictures with detailed
774 color system are shown in **Supplementary Fig. 1** and **Supplementary Fig. 2**.

775 **Fig. 3.** 100kb- (**a**), 20kb- (**b**), 10kb- (**c**), 5kb- (**d**), and 1kb-resolution (**e**) IR maps
776 (showing chromosome 1 only). The chromosome coordinate was shown in bp on the left

777 side of horizontal chromosomal rectangles, lined in every 2,500,000 bp. Introgressive
778 windows were shown in red. Non-introgressive windows were shown in yellow.
779 Windows of undetermined phylogeny were shown in gray. Each green rectangle stands
780 for a D-gene region. The 1kb-resolution windows (**e**) were shown in parallel with SSRs
781 (red lines) and CDRGs (blue lines). Magnified views for two regions in chr01 (**f**) and
782 chr04 (**g**) were exemplified as well.

783 **Fig. 4.** Numerical distribution of *DD* (*Distance Difference*). The *DD* statistics according
784 to dimensional continuity of all 1kb windows (a, average of all 12 chromosomes) and the
785 window proportion histogram of particular *DD*s (b, x-axis : *DD* in logarithmic scale, y-
786 axis : frequency of windows). *DD* is defined as below:

$$787 \quad DD = | F84 (\text{outgroup to } indica) - F84 (\text{outgroup to } japonica) |$$

788 ^(*) F84 = Felsenstein 84 nucleotide genetic distance

789 For more details of the formula, see **Methods**.

790 **Fig. 5.** Conceptual diagram of estimated introgression ages. The magnitudes of *DD*s
791 (*Distance Differences*, red scales) were overdrawn.

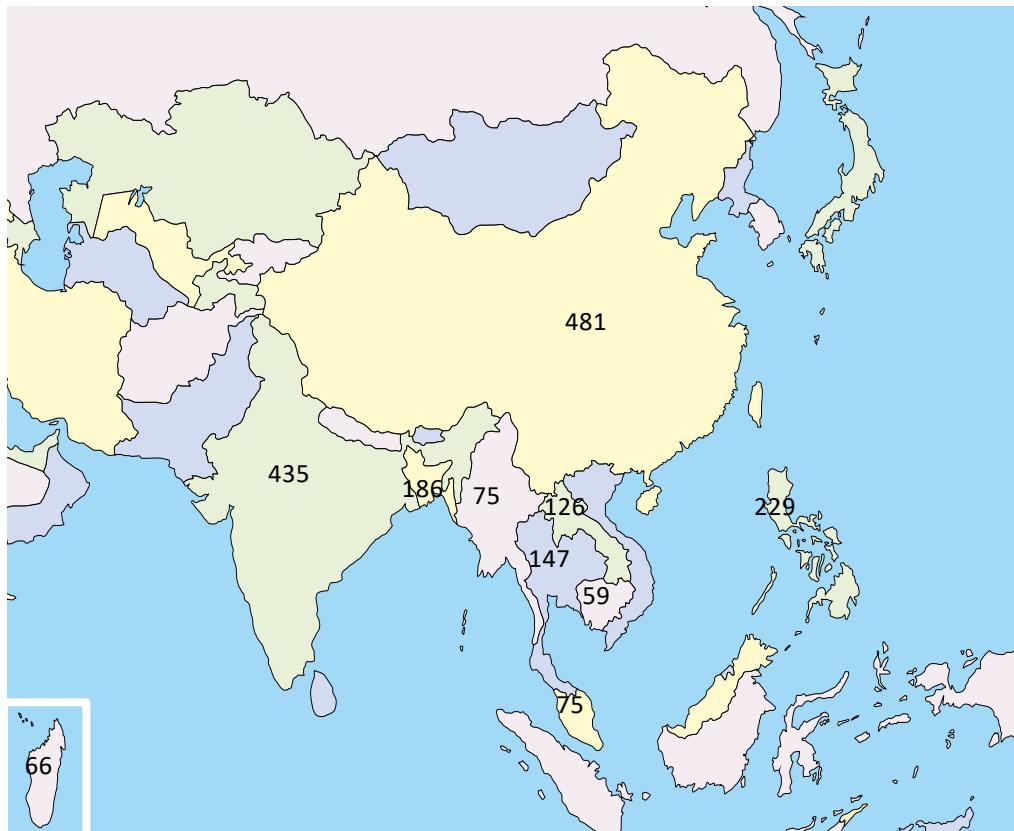
792 **Fig. 6.** The Sandcastles Model in domestication, a case scenario with three independent
793 introgression events. Each * (asterisk) stands for an agronomically beneficial allele.

794

a

	3000 Rice Genomes Project	RiceHap3	OryzaGenome	Rice3000+RiceHap3+OryzaGenome	Higher coverage wilds (AA)	<i>Oryza punctata</i> (BB, diploid)	Grand Total
reference	The 3000 rice genomes project 2014 Alexandrov et al. 2015 Mansueto et al. 2017 Wang et al. 2018	Huang et al. 2012	Ohyanagi et al. 2016	(This study)	Xu et al., 2012 Ohyanagi et al. 2016 Zhao et al. 2018 Stein et al. 2018	Stein et al. 2018	
reference genome	Os-Nipponbare-Reference-IRGSP-1.0	IRGSP-build4.0	Os-Nipponbare-Reference-IRGSP-1.0	Os-Nipponbare-Reference-IRGSP-1.0	Os-Nipponbare-Reference-IRGSP-1.0 (This study)	Os-Nipponbare-Reference-IRGSP-1.0 (This study)	
# of accessions cultivated	3,024	3,024	1,529	446	4,553	35	4,587
B#	246 (3KRice 2014 TableS1B)	-	-	-	246 (3KRice 2014 TableS1B)	-	-
CX#	312 (3KRice 2014 TableS1B)	-	-	-	312 (3KRice 2014 TableS1B)	-	-
IRIS_313-#	2466 (3KRice 2014 TableS1A)	-	-	-	2466 (3KRice 2014 TableS1A)	-	-
HP#	-	621 (Huang et al. 2012 TableS7)	-	-	621 (Huang et al. 2012 TableS7)	-	-
GP#	-	462 (Huang et al. 2012 TableS7)	-	-	462 (Huang et al. 2012 TableS7)	-	-
close-wild (<i>nivara</i> & <i>rufipogon</i>)	-	446 (Huang et al. 2012 TableS2)	446 (Ohyanagi et al. 2016 sup.data)	446 (Ohyanagi et al. 2016 sup.data)	-	32	-
distant-wild	-	-	-	-	-	3	1
Coverage (against Nipponbare)	High (14x in average)	Low (1x or 2x)	Low (2x)	High + Low (imputed)	High (12x each, at least)	High (140x)	
Is employed in preliminary outgroup assessment?	Yes	No	No	(No)	Yes	Yes	3,060 (Outgroup assessment)
Is employed in main analysis (introgression detection)?	Yes	No	No	(No)	No	Yes	3,025 (Main analysis)
Is employed in SSRs & CLDGRs recomputation?	(Yes)	(Yes)	(Yes)	Yes	No	No	4,553 (SSRs & CLDGRs recomputation)

b

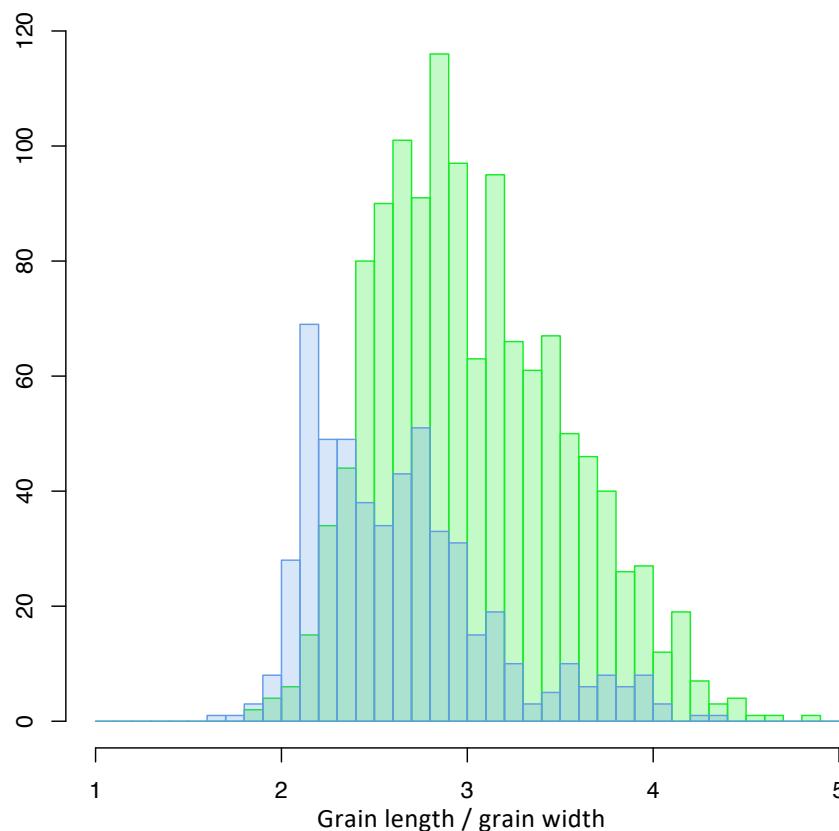


Origin of country	Number of accessions
China	481
India	435
Philippines	229
Bangladesh	186
Thailand	147
Laos	126
Myanmar	75
Malaysia	75
Madagascar	66
Cambodia	59
(Other countries)	374
(Origin unknown)	771

(In total 89 countries)

c

Frequency



The shortest grains

KHAW KAR 13::IRGC 36711-1
(*japonica*, 6.3 / 3.8 = 1.66)

MUTTU SAMBA::IRGC 36333-1
(*indica*, 5.7 / 3.0 = 1.90)

The longest grains

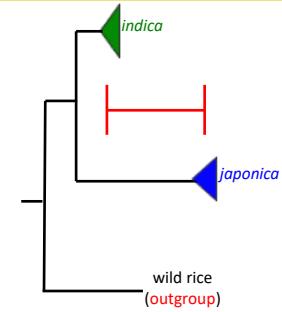
FORTUNA COLORADO::IRGC 703-1
(*japonica*, 10.4 / 2.4 = 4.33)

MAVOLATSIKA::IRGC 83137-1
(*indica*, 9.7 / 2.0 = 4.85)

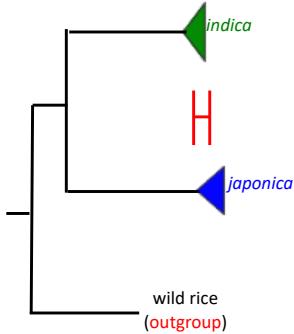
Grain length / grain width

Fig. 1. Passport data of domesticated and wild Asian rice accessions in this study (**a**, in total 4,587 accessions. for more details in higher coverage wilds, see **Supplementary Table 1**), and geographical origin of accessions in 3,000 Rice Genomes Project (**b**, 3,024 accessions). A typical phenotypic diversity within subspecies (**c**, grain length over grain width in *O. sativa* ssp. *indica* (n=1269, green) and *japonica* (n=533, blue)) .

a Non-Introgressive



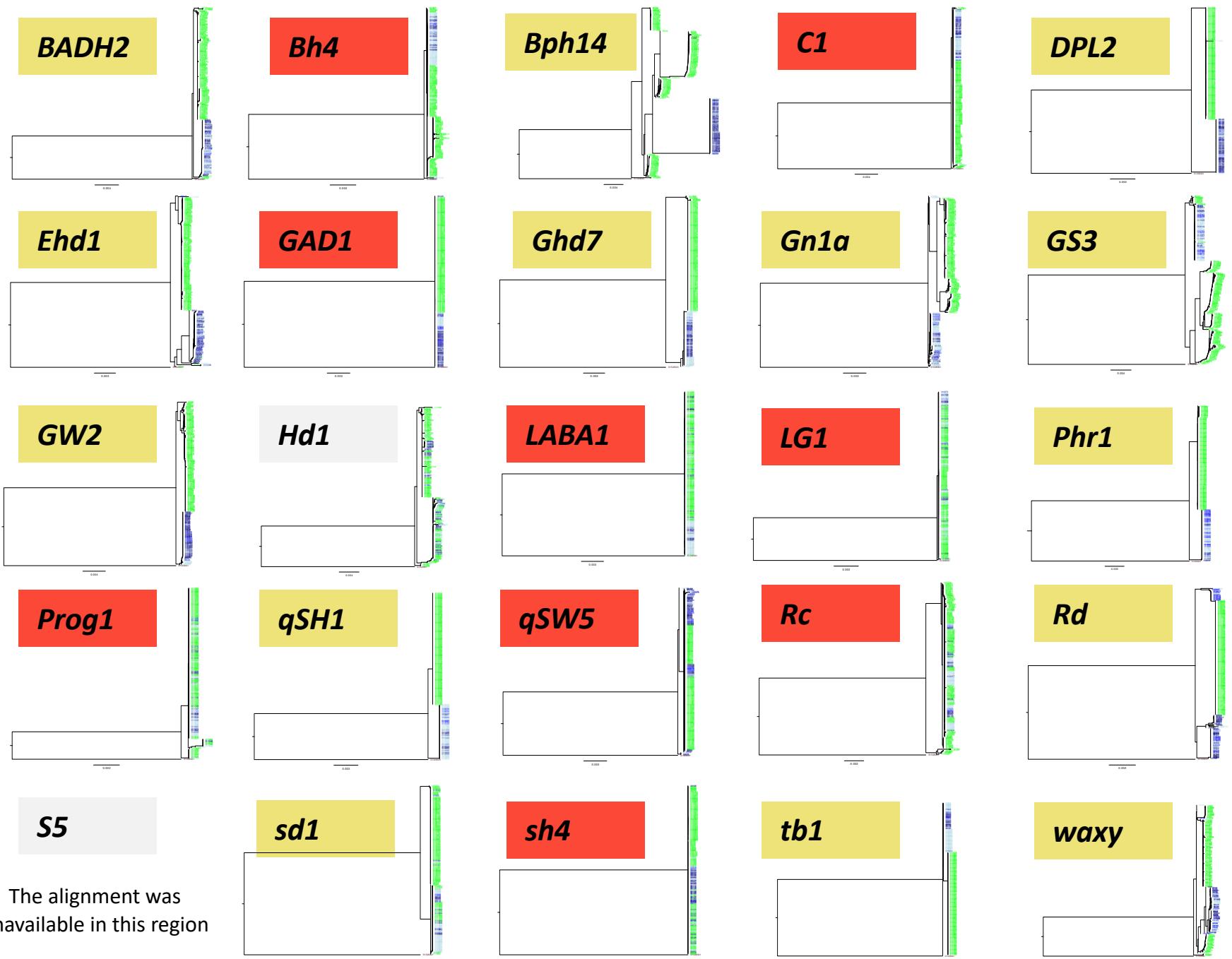
b Introgressive



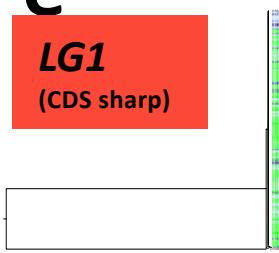
c

GeneSymbol	Description	Reference	LocusID	Location	Introgressive state (by visual inspection)
<i>BADH2</i>	Fragrance	52	Os08g0424500	chr08:20379823..20385975 (+ strand)	Non-introgressive
<i>Bh4</i>	Change hull color	53	Os04g0460200	chr04:22969845..22971859 (+ strand)	Introgressive
<i>Bph14</i>	Brown planthopper resistance	54	Os03g0848700	chr03:35693286..35699010 (- strand)	Non-introgressive
<i>C1</i>	Leaf sheath color and apiculus color	55	Os06g0205100	chr06:5315163..5316640 (+ strand)	Introgressive
<i>DPL2</i>	Hybrid incompatibility	56	Os06g0184100	chr06:4201250..4202851 (- strand)	Non-introgressive
<i>Ehd1</i>	Early heading date	57	Os10g0463400	chr10:17076098..17081344 (- strand)	Non-introgressive
<i>GAD1</i>	Grain number, length and awn development	58	Os08g0485500	chr08:23998787..24000176 (+ strand)	Introgressive
<i>Ghd7</i>	Heading date and yield potential	59	Os07g0261200	chr07:9152377..9155030 (- strand)	Non-introgressive
<i>Gn1a</i>	Grain number	60	Os01g0197700	chr01:5270449..5275585 (- strand)	Non-introgressive
<i>GS3</i>	Increase grain length	61	Os03g0407400	chr03:16729501..16735109 (- strand)	Non-introgressive
<i>GW2</i>	Grain width and weight	62	Os02g0244100	chr02:8115223..8121651 (+ strand)	Non-introgressive
<i>Hd1</i>	Heading date	63	Os06g0275000	chr06:9336376..9338569 (+ strand)	(undetermined)
<i>LABA1</i>	Long and barned awns	64	Os04g0518800	chr04:25959399..25963504 (+ strand)	Introgressive
<i>LG1</i>	Inflorescence architecture	65	Os04g0656500	chr04:33488722..33492700 (+ strand)	Introgressive
<i>Phr1</i>	Change hull color	66	Os03g0329900	chr03:12126320..12131242 (+ strand)	Non-introgressive
<i>Prog1</i>	Tiller erectness	67	Os07g0153600	chr07:2839194..2840089 (- strand)	Introgressive
<i>qSH1</i>	Seed shattering	68	Os01g0848400	chr01:36445456..36449951 (- strand)	Non-introgressive
<i>qSW5</i>	Increase grain width	69	Os05g0187500	chr05:5365122..5366701 (+ strand)	Introgressive
<i>Rc</i>	Change pericarp color	70	Os07g0211500	chr07:6062889..6069304 (+ strand)	Introgressive
<i>Rd</i>	Change pericarp color	71	Os01g0633500	chr01:25382714..25384678 (+ strand)	Non-introgressive
<i>S5</i>	Hybrid sterility	72	Os06g0213100	chr06:5759685..5761518 (+ strand)	(undetermined)
<i>sd1</i>	Reduce tiller length	73	Os01g0883800	chr01:38382385..38385469 (+ strand)	Non-introgressive
<i>sh4</i>	Seed shattering	74	Os04g0670900	chr04:34231186..34233221 (- strand)	Introgressive
<i>tb1</i>	Teosinte branched	75	Os03g0706500	chr03:28428504..28430438 (+ strand)	Non-introgressive
<i>waxy</i>	Amylose content	76	Os06g0133000	chr06:1765622..1770653 (+ strand)	Non-introgressive

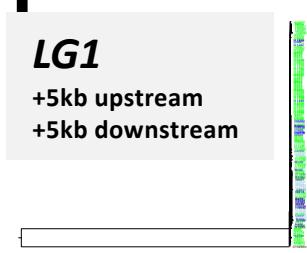
d



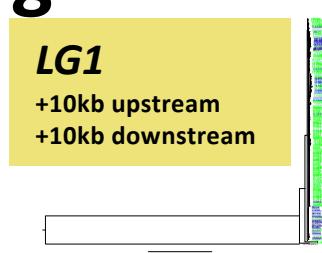
e



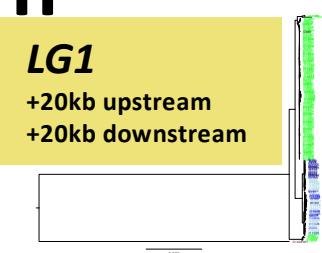
f



g



h



i

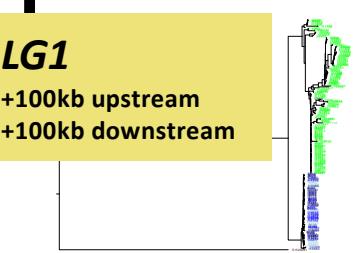


Fig. 2. Schematic view of underestimate on genetic *Distance Difference* (**a** and **b**), and phylogenetic analysis of manually curated D-genes (25 genes) and their determined introgressive states (**c** and **d**). Under isolated conditions, each of *indica* and *japonica* subpopulation shall show different genetic distance to the outgroup (a wild rice accession) to some extent, since they have been isolated from each other for a length of time (**a**), whereas they will show unexpectedly similar genetic distance to the outgroup when they made an inter-subspecies crossing (*i.e.* introgression) recently (**b**). Manually curated D-genes (25 genes) and their determined introgressive state (**c**). Reconstructed phylogenetic trees of 25 D-genes (**d**), green nodes : *indica*, blue nodes : *japonica*. Non-introgressive genes were shown in yellow background. Introgressive genes were shown in red background. Genes of undetermined phylogeny were shown in gray background. Phylogenetic trees for one of the D-genes (*LG1*) with variable length of flanking upstream/downstream regions (**e** : CDS only, **f** : +5kb-upstream/+5kb-downstream, **g** : +10kb-upstream/+10kb-downstream, **h** : +20kb-upstream/+20kb-downstream, and **i** : +100kb-upstream/+100kb-downstream, respectively). Full size tree pictures with detailed color system are shown in **Supplementary Fig. 1** and **Supplementary Fig. 2**.

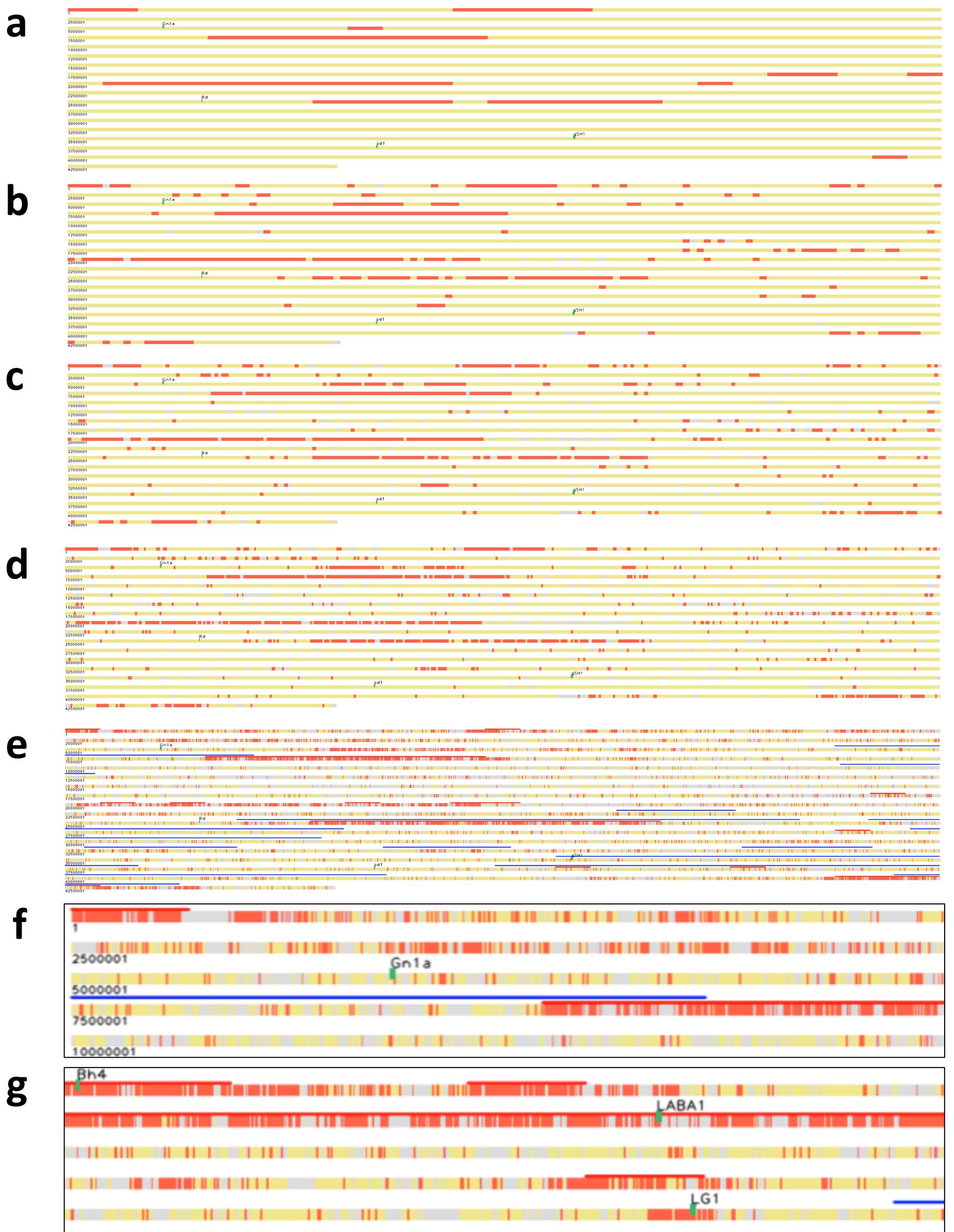


Fig. 3. 100kb- (a), 20kb- (b), 10kb- (c), 5kb- (d), and 1kb-resolution (e) IR maps (showing chromosome 1 only). The chromosome coordinate was shown in bp on the left side of horizontal chromosomal rectangles, lined in every 2,500,000 bp. Introgressive windows were shown in red. Non-introgressive windows were shown in yellow. Windows of undetermined phylogeny were shown in gray. Each green rectangle stands for a D-gene region. The 1kb-resolution windows (e) were shown in parallel with SSRs (red lines) and CDRGs (blue lines). Magnified views for two regions in chr01 (f) and chr04 (g) were exemplified as well.

a

all chromosomes

	counts	counts (%)	outgroup to <i>indica</i> (F84 distance)	outgroup to <i>japonica</i> (F84 distance)	<i>DD</i>
overall windows	373,204	100			
phylogeny N.D. windows	133,623	35.8			
phylogeny determined windows	239,581	64.2	0.055106967	0.053881707	1.23E-03
non-introgressive windows	186,567	50.0	0.055653817	0.053942041	1.71E-03
introgressive windows (all)	53,014	14.2	0.05318249	0.05366938	4.87E-04
introgressive windows (narrow = 1)	18,814	5.04	0.052480345	0.053064024	5.84E-04
introgressive windows (wide >= 40)	334	0.0895	0.055056613	0.055050718	5.89E-06

b

Frequency

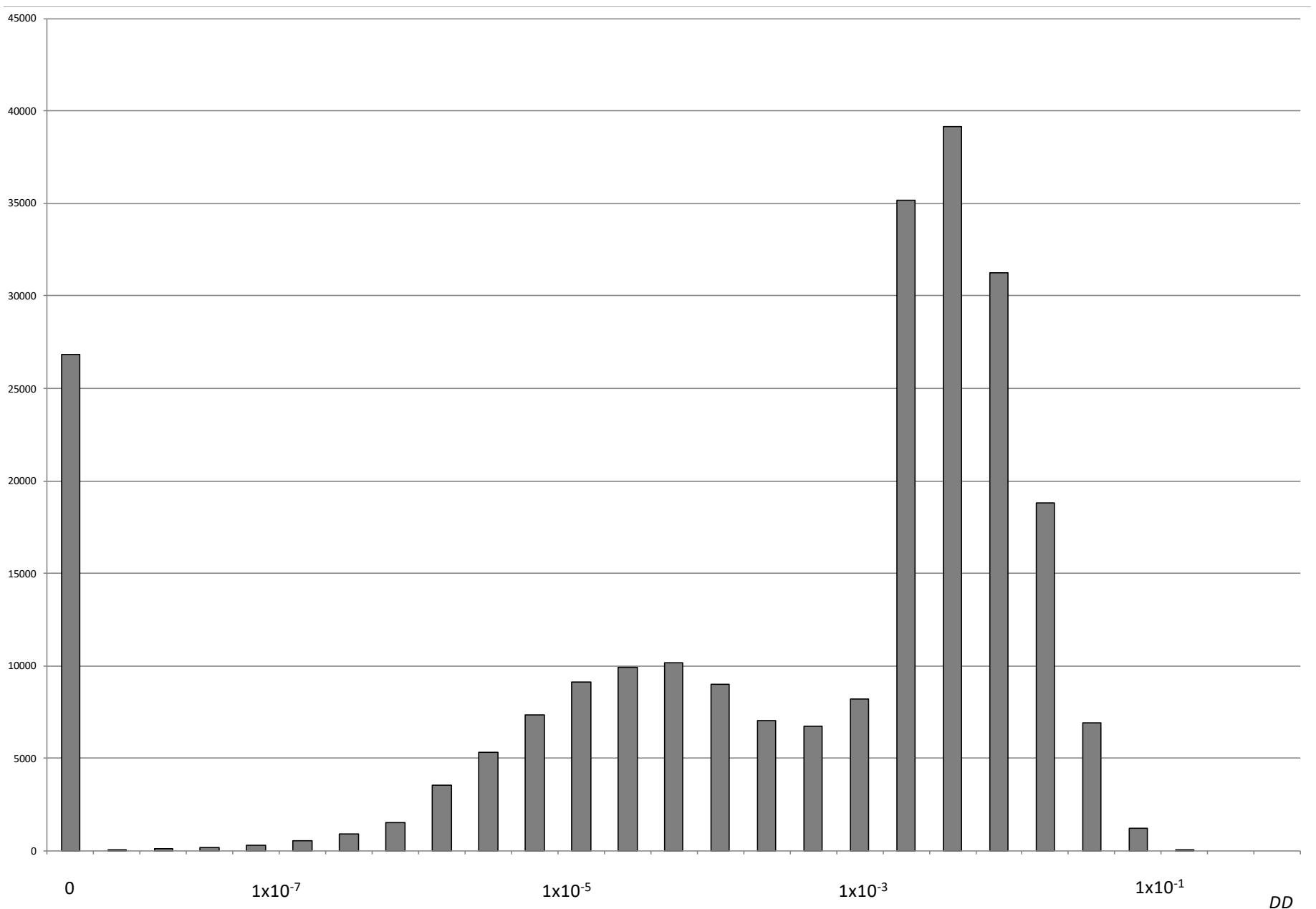


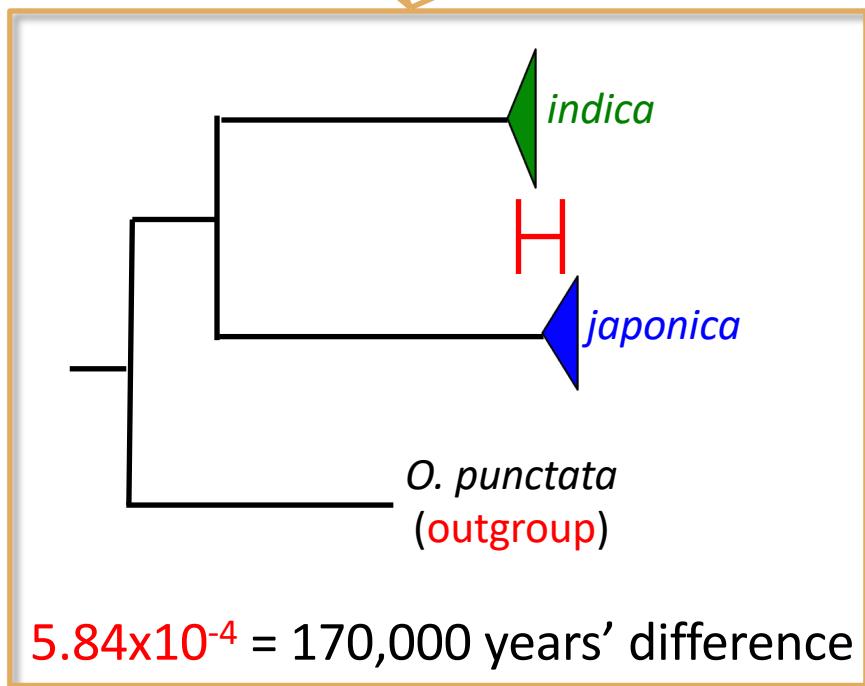
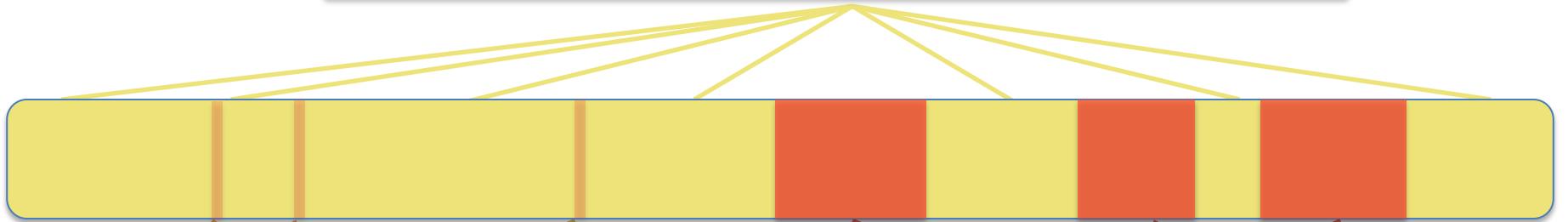
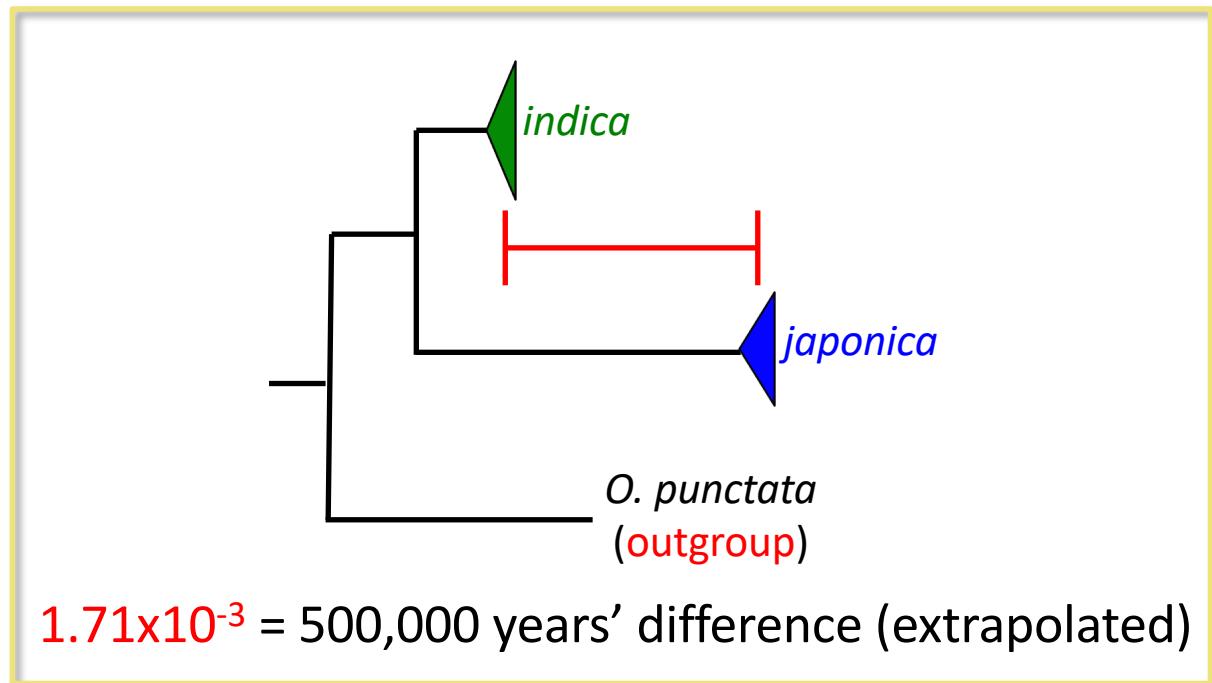
Fig. 4. Numerical distribution of *DD* (*Distance Difference*). The *DD* statistics according to dimensional continuity of all 1kb windows (**a**, average of all 12 chromosomes) and the window proportion histogram of particular *DD*s (**b**, x-axis : *DD* in logarithmic scale, y-axis : frequency of windows). *DD* is defined as below:

$$DD = | \text{F84 (outgroup to } indica) - \text{F84 (outgroup to } japonica) |$$

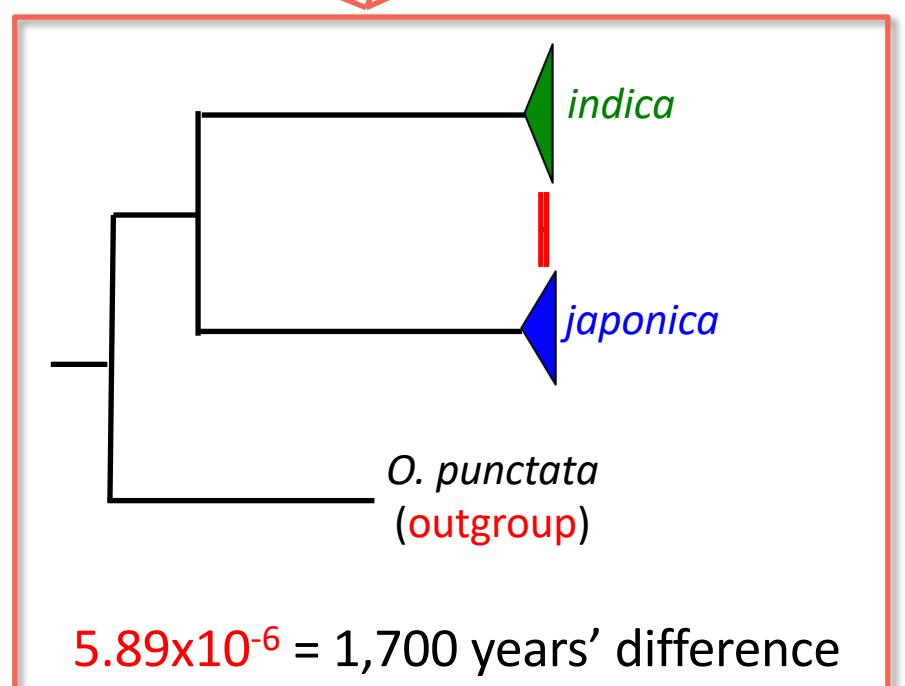
(*) F84 = Felsenstein 84 nucleotide genetic distance

For more details of the formula, see **Methods**.

Non-IRs



Narrow IRs



Wide IRs

Fig. 5. Conceptual diagram of estimated introgression ages. The magnitudes of *DDs* (Distance Differences, red scales) were overdrawn.

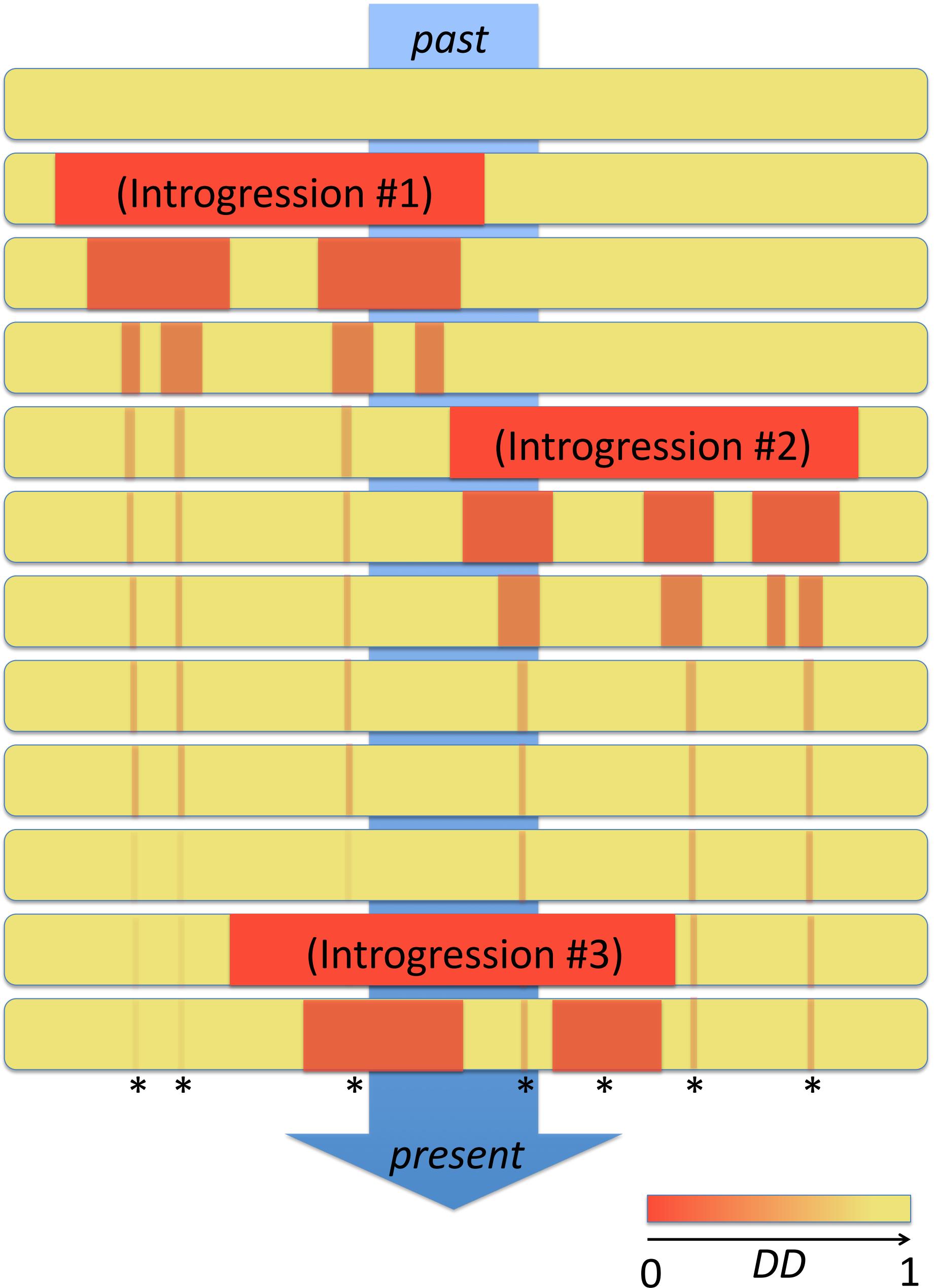


Fig. 6. The Sandcastles Model in domestication, a case scenario with three independent introgression events. Each * (asterisk) stands for an agronomically beneficial allele.