

1 **Nanopore-based genome assembly and the evolutionary genomics of basmati rice**

2

3

4 Jae Young Choi<sup>1\*</sup>, Zoe N. Lye<sup>1</sup>, Simon C. Groen<sup>1</sup>, Xiaoguang Dai<sup>2</sup>, Priyesh Rughani<sup>2</sup>, Sophie  
5 Zaaiker<sup>3</sup>, Eoghan D. Harrington<sup>2</sup>, Sissel Juul<sup>2</sup> and Michael D. Purugganan<sup>1,4\*</sup>

6

7

8 <sup>1</sup>Center for Genomics and Systems Biology, Department of Biology, New York University, New  
9 York, New York, USA

10 <sup>2</sup>Oxford Nanopore Technologies, New York, New York, USA

11 <sup>3</sup>New York Genome Center, New York, New York, USA

12 <sup>4</sup>Center for Genomics and Systems Biology, NYU Abu Dhabi Research Institute, New York  
13 University Abu Dhabi, Abu Dhabi, United Arab Emirates

14

15

16 \* Corresponding authors, Email: [jyc387@nyu.edu](mailto:jyc387@nyu.edu) (JYC), [mp132@nyu.edu](mailto:mp132@nyu.edu) (MDP)

17

18

19

20

21

22

23

24 **ABSTRACT**

25 **BACKGROUND**

26         The *circum*-basmati group of cultivated Asian rice (*Oryza sativa*) contains many iconic  
27 varieties and is widespread in the Indian subcontinent. Despite its economic and cultural  
28 importance, a high-quality reference genome is currently lacking, and the group's evolutionary  
29 history is not fully resolved. To address these gaps, we used long-read nanopore sequencing and  
30 assembled the genomes of two *circum*-basmati rice varieties, Basmati 334 and Dom Sufid.

31

32 **RESULTS**

33         We generated two high-quality, chromosome-level reference genomes that represented  
34 the 12 chromosomes of *Oryza*. The assemblies showed a contig N50 of 6.32Mb and 10.53Mb for  
35 Basmati 334 and Dom Sufid, respectively. Using our highly contiguous assemblies we  
36 characterized structural variations segregating across *circum*-basmati genomes. We discovered  
37 repeat expansions not observed in japonica—the rice group most closely related to *circum*-  
38 basmati—as well as presence/absence variants of over 20Mb, one of which was a *circum*-  
39 basmati-specific deletion of a gene regulating awn length. We further detected strong evidence of  
40 admixture between the *circum*-basmati and *circum*-aus groups. This gene flow had its greatest  
41 effect on chromosome 10, causing both structural variation and single nucleotide polymorphism  
42 to deviate from genome-wide history. Lastly, population genomic analysis of 78 *circum*-basmati  
43 varieties showed three major geographically structured genetic groups: (1) Bhutan/Nepal group,  
44 (2) India/Bangladesh/Myanmar group, and (3) Iran/Pakistan group.

45

46 **CONCLUSION**

47 Availability of high-quality reference genomes from nanopore sequencing allowed  
48 functional and evolutionary genomic analyses, providing genome-wide evidence for gene flow  
49 between *circum*-aus and *circum*-basmati, the nature of *circum*-basmati structural variation, and  
50 the presence/absence of genes in this important and iconic rice variety group.

51

## 52 **KEYWORDS**

53 *Oryza sativa*, Asian rice, aromatic rice group, domestication, crop evolution, nanopore  
54 sequencing, aus, basmati, indica, japonica, admixture, awnless, *de novo* genome assembly

55

## 56 **BACKGROUND**

57 *Oryza sativa* or Asian rice is an agriculturally important crop that feeds one-half of the  
58 world's population [1], and supplies 20% of people's caloric intake (www.fao.org). Historically,  
59 *O. sativa* has been classified into two major variety groups, japonica and indica, based on  
60 morphometric differences and molecular markers [2, 3]. These variety groups can be considered  
61 as subspecies, particularly given the presence of reproductive barriers between them [4].  
62 Archaeobotanical remains suggest japonica rice was domesticated ~9,000 years ago in the  
63 Yangtze Basin of China, while indica rice originated ~4,000 years ago when domestication  
64 alleles were introduced from japonica into either *O. nivara* or a proto-indica in the Indian  
65 subcontinent [5]. More recently, two additional variety groups have been recognized that are  
66 genetically distinct from japonica and indica: the aus/*circum*-aus and aromatic/*circum*-basmati  
67 rices [6–8].

68 The rich genetic diversity of Asian rice is likely a result from a complex domestication  
69 process involving multiple wild progenitor populations and the exchange of important

70 domestication alleles between *O. sativa* variety groups through gene flow [5, 7, 9–17].  
71 Moreover, many agricultural traits within rice are variety group-specific [18–23], suggesting  
72 local adaptation to environments or cultural preferences have partially driven the diversification  
73 of rice varieties.

74 Arguably, the *circum*-basmati rice group has been the least studied among the four major  
75 variety groups, and it was only recently defined in more detail based on insights from genomic  
76 data [7]. Among its members the group boasts the iconic basmati rices (*sensu stricto*) from  
77 southern Asia and the sadri rices from Iran [6]. Many, but not all, *circum*-basmati varieties are  
78 characterized by distinct and highly desirable fragrance and texture [24]. Nearly all fragrant  
79 *circum*-basmati varieties possess a loss-of-function mutation in the *BADH2* gene that has its  
80 origins in ancestral japonica haplotypes, suggesting that an introgression between *circum*-  
81 basmati and japonica may have led to fragrant basmati rice [21, 25, 26]. Genome-wide  
82 polymorphism analysis of a smaller array of *circum*-basmati rice cultivars shows close  
83 association with japonica varieties [7, 16, 27], providing evidence that at least part of the  
84 genomic make-up of *circum*-basmati rices may indeed be traced back to japonica.

85 Whole-genome sequences are an important resource for evolutionary geneticists studying  
86 plant domestication, as well as breeders aiming to improve crop varieties. Single-molecule  
87 sequencing regularly produces sequencing reads in the range of kilobases (kb) [28]. This is  
88 particularly helpful for assembling plant genomes, which are often highly repetitive and  
89 heterozygous, and commonly underwent at least one round of polyploidization in the past [29–  
90 31]. The *Oryza sativa* genome, with a relatively modest size of ~400 Mb, was the first crop  
91 genome sequence assembled [29], and there has been much progress in generating *de novo*  
92 genome assemblies for other members of the genus *Oryza*. Currently, there are assemblies for

93 nine wild species (*Leersia perrieri* [outgroup], *O. barthii*, *O. brachyantha*, *O. glumaepatula*, *O.*  
94 *longistaminata*, *O. meridionalis*, *O. nivara*, *O. punctata*, and *O. rufipogon*) and two domesticated  
95 species (*O. glaberrima* and *O. sativa*) [32–37].

96 Within domesticated Asian rice (*O. sativa*), genome assemblies are available for cultivars  
97 in most variety groups [32, 33, 38–42]. However, several of these reference assemblies are based  
98 on short-read sequencing data and show higher levels of incompleteness compared to assemblies  
99 generated from long-read sequences [40, 41]. Nevertheless, these *de novo* genome assemblies  
100 have been critical in revealing genomic variation (*e.g.* variations in genome structure and  
101 repetitive DNA, and *de novo* species- or population-specific genes) that were otherwise missed  
102 from analyzing a single reference genome. Recently, a genome assembly based on short-read  
103 sequencing data was generated for basmati rice [42]. Not only were there missing sequences in  
104 this assembly, it was also generated from DNA of an elite basmati breeding line. Such modern  
105 cultivars are not the best foundations for domestication-related analyses due to higher levels of  
106 introgression from other rice populations during modern breeding.

107 Here, we report the *de novo* sequencing and assembly of the landraces (traditional  
108 varieties) Basmati 334 [21, 43, 44] and Dom Sufid [21, 24, 45, 46] using the long-read nanopore  
109 sequencing platform of Oxford Nanopore Technologies [47]. Basmati 334 is from Pakistan,  
110 evolved in a rainfed lowland environment and is known to be drought tolerant at the seedling and  
111 reproductive stages [44]. It also possesses several broad-spectrum bacterial blight resistance  
112 alleles [48, 49], making Basmati 334 desirable for breeding resilience into modern basmati  
113 cultivars [49, 50]. Dom Sufid is an Iranian sadri cultivar that, like other sadri and basmati (*sensu*  
114 *stricto*) varieties, is among the most expensive varieties currently available in the market [24]. It  
115 has desirable characteristics such as aromaticity and grain elongation during cooking, although it

116 is susceptible to disease and abiotic stress [24, 51]. Because of their special characteristics, both  
117 Basmati 334 and Dom Sufid are used in elite rice breeding programs to create high yielding and  
118 resilient aromatic rice varieties [24, 44–46, 50].

119 Based on long reads from nanopore sequencing, our genome assemblies have high  
120 quality, contiguity, and genic completeness, making them comparable in quality to assemblies  
121 associated with key rice reference genomes. We used our *circum*-basmati genome assemblies to  
122 characterize genomic variation existing within this important rice variety group, and analyze  
123 domestication-related and other evolutionary processes that shaped this variation. Our *circum*-  
124 basmati rice genome assemblies will be valuable complements to the available assemblies for  
125 other rice cultivars, unlocking important genomic variation for rice crop improvement.

126

## 127 RESULTS

128 **Nanopore sequencing of basmati and sadri rice.** Using Oxford Nanopore Technologies' long-  
129 read sequencing platform, we sequenced the genomes of the *circum*-basmati landraces Basmati  
130 334 (basmati *sensu stricto*) and Dom Sufid (sadri). We called 1,372,950 reads constituting a total  
131 of 29.2 Gb for Basmati 334 and 1,183,159 reads constituting a total of 24.2 Gb for Dom Sufid  
132 (Table 1). For both samples the median read length was > 17 kb, the read length N50 was > 33  
133 kb, and the median quality score per read was ~11.

134

135 **Table 1. Summary of nanopore sequencing read data.**

Flow-cell	Number of Reads	Median Read Length	Read Length N50	Median Quality Score (QS)	Total Bases
Basmati 334					

FAK30515	288,473	19,905	36,743	11.27	6,843,069,570
FAK30732	306,247	18,792	30,974	11.19	6,341,851,953
FAK30522	228,191	17,366	36,456	11.07	4,816,938,523
FAK27872	244,606	18,335	31,267	11.35	5,045,781,146
FAK27919	305,433	18,087	30,727	11.43	6,191,306,294
All	1,372,950	18,576	33,005	11.27	29,238,947,486

---

Dom Sufid

FAK30464	300,290	18,477	37,754	11.34	6,681,819,859
FAK30582	258,584	17,641	34,213	11.30	5,353,774,444
FAK28890	330,924	16,756	34,033	10.96	6,553,200,184
FAK30064	293,361	16,178	32,835	10.99	5,618,557,776
All	1,183,159	17,237	34,728	11.14	24,207,352,263

---

136

137

138 ***De novo* assembly of the Basmati 334 and Dom Sufid rice genomes.** Incorporating only those

139 reads that had a mean quality score of > 8 and read lengths of > 8 kb, we used a total of

140 1,076,192 reads and 902,040 reads for the Basmati 334 and Dom Sufid genome assemblies,

141 which resulted in a genome coverage of ~62× and ~51×, respectively (Table 2). We polished the

142 genome assemblies with both nanopore and short Illumina sequencing reads. The final, polished

143 genome assemblies spanned 386.5 Mb across 188 contigs for Basmati 334, and 383.6 Mb across

144 116 contigs for Dom Sufid. The genome assemblies had high contiguity, with a contig N50 of

145 6.32 Mb and 10.53 Mb for Basmati 334 and Dom Sufid, respectively. Our genome assemblies

146 recovered more than 97% of the 1,440 BUSCO [52] embryophyte gene groups, which is

147 comparable to the BUSCO statistics for the japonica Nipponbare [33] (98.4%) and indica R498  
148 reference genomes [41] (98.0%). This is an improvement from the currently available genome  
149 assembly of basmati variety GP295-1 [42], which was generated from Illumina short-read  
150 sequencing data and has a contig N50 of 44.4 kb with 50,786 assembled contigs.

151 We examined coding sequences of our *circum*-basmati genomes by conducting gene  
152 annotation using published rice gene models and the *MAKER* gene annotation pipeline [52, 53].  
153 A total of 41,270 genes were annotated for the Basmati 334 genome, and 38,329 for the Dom  
154 Sufid genome. BUSCO gene completion analysis [52] indicated that 95.4% and 93.6% of the  
155 3,278 single copy genes from the liliopsida gene dataset were found in the Basmati 334 and Dom  
156 Sufid gene annotations respectively.

157

158 **Table 2. Summary of the *circum*-basmati rice genome assemblies**

	Basmati 334	Dom Sufid
Genome Coverage	62.5×	51.4×
Number of Contigs	188	116
Total Number of Bases in Contigs	386,555,741	383,636,250
Total Number of Bases Scaffolded	386,050,525	383,245,802
Contig N50 Length	6.32 Mb	10.53 Mb
Contig L50	20	13
Total Contigs > 50 kbp	159	104
Maximum Contig Length	17.04 Mb	26.82 Mb
BUSCO Gene Completion (Assembly)	97.6%	97.0%
GC Content	43.6%	43.7%



Repeat Content	44.4%	44.2%
Number of Annotated Genes	41,270	38,329
BUSCO Gene Completion (Annotation)	95.4%	93.6%

---

159

160 **Whole genome comparison to other rice variety group genomes.** We aligned our draft  
161 genome assemblies to the japonica Nipponbare reference genome sequence [33], which  
162 represents one of the highest quality reference genome sequences (Figure 1A). Between the  
163 Nipponbare, Basmati 334 and Dom Sufid genomes, high levels of macro-synteny were evident  
164 across the japonica chromosomes. Specifically, we observed little large-scale structural variation  
165 between Basmati 334 and Dom Sufid contigs and the japonica genome. A noticeable exception  
166 was an apparent inversion in the *circum*-basmati genome assemblies at chromosome 6 between  
167 positions 12.5 Mb and 18.7 Mb (Nipponbare coordinates), corresponding to the pericentromeric  
168 region [54]. Interestingly, the same region showed an inversion between the Nipponbare and  
169 indica R498 reference genomes [41], whereas in the *circum*-aus N22 cultivar no inversions are  
170 observed (Supplemental Figure 1). While the entire region was inverted in R498, the inversion  
171 positions were disjoint in Basmati 334 and Dom Sufid, apparently occurring in multiple regions  
172 of the pericentromere. We independently verified the inversions by aligning raw nanopore  
173 sequencing reads to the Nipponbare reference genome using the long read-aware aligner *ngmlr*  
174 [55], and the structural variation detection program *sniffles* [55]. *Sniffles* detected several  
175 inversions, including a large inversion between positions 13.1 and 17.7 Mb and between 18.18  
176 and 18.23 Mb, with several smaller inversions located within the largest inversion (Supplemental  
177 Table 1).

178           Because of high macro-synteny with japonica (Figure 1A), we ordered and oriented the  
179 contigs of the Basmati 334 and Dom Sufid assemblies using a reference genome-based  
180 scaffolding approach [56]. For both Basmati 334 and Dom Sufid, over 99.9% of the assembled  
181 genomic contigs were anchored to the Nipponbare reference genome (Table 2). The scaffolded  
182 *circum*-basmati chromosomes were similar in size to those in reference genomes for cultivars in  
183 other rice variety groups (Nipponbare [33], the *circum*-aus variety N22 [37], and the indica  
184 varieties IR8 [37] and R498 [41]) that were sequenced, assembled, and scaffolded to near  
185 completion (Table 3).

186

187 **Table 3. Comparison of assembled chromosome sizes for cultivars across variety groups.**

Chromosome	Basmati 334	Dom Sufid	Nipponbare	N22	IR8	R498
1	44,411,451	44,306,286	43,270,923	44,711,178	44,746,683	44,361,539
2	35,924,761	36,365,206	35,937,250	38,372,633	37,475,564	37,764,328
3	40,305,655	38,133,813	36,413,819	36,762,248	39,065,119	39,691,490
4	34,905,232	34,714,597	35,502,694	33,558,078	35,713,470	35,849,732
5	30,669,872	31,017,353	29,958,434	28,792,057	31,269,760	31,237,231
6	29,982,228	32,412,977	31,248,787	29,772,976	32,072,649	32,465,040
7	30,410,531	29,511,326	29,697,621	29,936,233	30,380,234	30,277,827
8	29,921,941	29,962,976	28,443,022	25,527,801	30,236,384	29,952,003
9	24,050,083	23,970,096	23,012,720	22,277,206	24,243,884	24,760,661
10	25,596,481	24,989,786	23,207,287	20,972,683	25,246,678	25,582,588
11	29,979,012	29,949,236	29,021,106	29,032,419	32,337,678	31,778,392
12	29,893,278	27,912,150	27,531,856	22,563,585	25,963,606	26,601,357

---

Total	386,050,525	383,245,802	373,245,519	362,279,097	388,751,709	390,322,188
-------	-------------	-------------	-------------	-------------	-------------	-------------

---

188

189           Next, we assessed the assembly quality of the *circum*-basmati genomes by contrasting  
190 them against available *de novo*-assembled genomes within the Asian rice complex (see Materials  
191 and Method for a complete list of genomes). We generated a multi-genome alignment to the  
192 Nipponbare genome, which we chose as the reference since its assembly and gene annotation is a  
193 product of years of community-based efforts [33, 57, 58]. To infer the quality of the gene regions  
194 in each of the genome assemblies, we used the multi-genome alignment to extract the coding  
195 DNA sequence of each Nipponbare gene and its orthologous regions from each non-japonica  
196 genome. The orthologous genes were counted for missing DNA sequences (“N” sequences) and  
197 gaps to estimate the percent of Nipponbare genes covered. For all genomes the majority of  
198 Nipponbare genes had a near-zero proportion of sites that were missing in the orthologous non-  
199 Nipponbare genes (Supplemental Figure 2). The missing proportions of Nipponbare-orthologous  
200 genes within the Basmati 334 and Dom Sufid genomes were comparable to those for genomes  
201 that had higher assembly contiguity [37, 40, 41].

202           Focusing on the previously sequenced basmati GP295-1 genome [42], our newly  
203 assembled *circum*-basmati genomes had noticeably lower proportions of missing genes  
204 (Supplemental Figure 2). Furthermore, over 96% of base pairs across the Nipponbare genome  
205 were alignable against the Basmati 334 (total of 359,557,873 bp [96.33%] of Nipponbare  
206 genome) or Dom Sufid (total of 359,819,239 bp [96.40%] of Nipponbare genome) assemblies,  
207 while only 194,464,958 bp (52.1%) of the Nipponbare genome were alignable against the  
208 GP295-1 assembly.

209           We then counted the single-nucleotide and insertion/deletion (indel, up to ~60 bp)  
210 differences between the *circum*-basmati and Nipponbare assemblies to assess the overall quality  
211 of our newly assembled genomes. To avoid analyzing differences across unconstrained repeat  
212 regions, we specifically examined regions where there were 20 exact base-pair matches flanking  
213 a site that had a single nucleotide or indel difference between the *circum*-basmati and  
214 Nipponbare genomes. In the GP295-1 genome there were 334,500 (0.17%) single-nucleotide  
215 differences and 44,609 (0.023%) indels compared to the Nipponbare genome. Our newly  
216 assembled genomes had similar proportions of single-nucleotide differences with the Nipponbare  
217 genome, where the Basmati 334 genome had 780,735 (0.22%) differences and the Dom Sufid  
218 genome had 731,426 (0.20%). For indels the Basmati 334 genome had comparable proportions  
219 of differences with 104,282 (0.029%) variants, but the Dom Sufid genome had higher  
220 proportions with 222,813 (0.062%) variants. In sum, our draft *circum*-basmati genomes had high  
221 contiguity and completeness as evidenced by assembly to the chromosome level, and comparison  
222 to the Nipponbare genome. In addition, our genome assemblies were comparable to the Illumina  
223 sequence-generated GP295-1 genome for the proportion of genomic differences with the  
224 Nipponbare genome, suggesting they had high quality and accuracy as well.

225           Our *circum*-basmati genome assemblies should also be of sufficiently high quality for  
226 detailed gene-level analysis. For instance, a hallmark of many *circum*-basmati rices is  
227 aromaticity, and a previous study had determined that Dom Sufid, but not Basmati 334, is a  
228 fragrant variety [21]. We examined the two genomes to verify the presence or absence of the  
229 mutations associated with fragrance. There are multiple different loss-of-function mutations in  
230 the *BADH2* gene that cause rice varieties to be fragrant [21, 25, 26], but the majority of fragrant  
231 rices carry a deletion of 8 nucleotides at position chr8:20,382,861-20,382,868 of the Nipponbare

232 genome assembly (version Os-Nipponbare-Reference-IRGSP-1.0). Using the genome alignment,  
233 we extracted the *BADH2* sequence region to compare the gene sequence of the non-fragrant  
234 Nipponbare to that of Basmati 334 and Dom Sufid. Consistent with previous observations [21],  
235 we found that the genome of the non-fragrant Basmati 334 did not carry the deletion and  
236 contained the wild-type *BADH2* haplotype observed in Nipponbare. The genome of the fragrant  
237 Dom Sufid, on the other hand, carried the 8-bp deletion, as well as the 3 single-nucleotide  
238 polymorphisms flanking the deletion. This illustrates that the Basmati 334 and Dom Sufid  
239 genomes are accurate enough for gene-level analysis.

240

241 ***Circum-basmati* gene analysis.** Our annotation identified ~40,000 coding sequences in the  
242 *circum-basmati* assemblies. We examined population frequencies of the annotated gene models  
243 across a *circum-basmati* population dataset to filter out mis-annotated gene models or genes at  
244 very low frequency in a population. We obtained Illumina sequencing reads from varieties  
245 included in the 3K Rice Genome Project [7] and sequenced additional varieties to analyze a total  
246 of 78 *circum-basmati* cultivars (see Supplemental Table 2 for a list of varieties). The Illumina  
247 sequencing reads were aligned to the *circum-basmati* genomes, and if the average coverage of a  
248 genic region was  $< 0.05\times$  for an individual this gene was called as a deletion in that variety.  
249 Since we used a low threshold for calling a deletion, the genome-wide sequencing coverage of a  
250 variety did not influence the number of gene deletions detected (Supplemental Figure 3). Results  
251 showed that gene deletions were indeed rare across the *circum-basmati* population (Figure 2A),  
252 consistent with their probable deleterious nature. We found that 31,565 genes (76.5%) in  
253 Basmati 334 and 29,832 genes (77.8%) in the Dom Sufid genomes did not have a deletion across  
254 the population (see Supplemental Table 3 for a list of genes).

255           There were 517 gene models from Basmati 334 and 431 gene models from Dom Sufid  
256 that had a deletion frequency of  $\geq 0.3$  (see Supplemental Table 4 for a list of genes). These gene  
257 models with high deletion frequencies were not considered further in this analysis. The rest were  
258 compared against the *circum*-aus N22, indica R498, and japonica Nipponbare gene models to  
259 determine their orthogroup status (Figure 2B; see Supplemental Table 5 for a list of genes and  
260 their orthogroup status), which are sets of genes that are orthologs and recent paralogs of each  
261 other [59].

262           The most frequent orthogroup class observed was for groups in which every rice variety  
263 group has at least one gene member. There were 13,894 orthogroups within this class, consisting  
264 of 17,361 genes from N22, 18,302 genes from Basmati 334, 17,936 genes from Dom Sufid,  
265 17,553 genes from R498, and 18,351 genes from Nipponbare. This orthogroup class likely  
266 represents the set of core genes of *O. sativa* [42]. The second-highest orthogroup class observed  
267 was for groups with genes that were uniquely found in both *circum*-basmati genomes (3,802  
268 orthogroups). These genes represent those restricted to the *circum*-basmati group.

269           In comparison to genes in other rice variety groups, the *circum*-basmati genes shared the  
270 highest number of orthogroups with *circum*-aus (2,648 orthogroups), followed by japonica (1,378  
271 orthogroups), while sharing the lowest number of orthogroups with indica (663 orthogroups). In  
272 fact, genes from indica variety R498 had the lowest number assigned to an orthogroup (Figure  
273 2B inset table), suggesting this genome had more unique genes, *i.e.* without orthologs/paralogs to  
274 genes in other rice variety groups.

275  
276       **Genome-wide presence/absence variation within the *circum*-basmati genomes.** Our  
277 assembled *circum*-basmati genomes were >10 Mb longer than the Nipponbare genome, but

278 individual chromosomes showed different relative lengths (Table 3) suggesting a considerable  
279 number of presence/absence variants (PAVs) between the genomes. We examined the PAVs  
280 between the *circum*-basmati and Nipponbare genomes using two different computational  
281 packages: (i) *sniffles*, which uses raw nanopore reads aligned to a reference genome to call  
282 PAVs, and (ii) *assemblytics* [60], which aligns the genome assemblies to each other and calls  
283 PAVs. The results showed that, while the total number of PAVs called by *sniffles* and  
284 *assemblytics* were similar, only ~36% of PAVs had overlapping positions (Table 4). In addition,  
285 the combined total size of PAVs was larger for predictions made by *sniffles* compared to those  
286 by *assemblytics*. For subsequent analysis we focused on PAVs that were called by both methods.  
287

288 **Table 4. Comparison of presence/absence variation called by two different computational**  
289 **packages.**

	<i>sniffles</i>	<i>assemblytics</i>	overlap
Basmati 334			
Deletion Counts	11,989	11,247	4051
Deleted Basepairs	43,768,763	29,048,238	19,328,854
Insertion Counts	11,447	12,161	3734
Inserted Basepairs	19,650,518	14,498,550	5,783,551
Dom Sufid			
Deletion Counts	9901	10,115	3649
Deleted Basepairs	36,600,114	26,128,143	17,274,967
Insertion Counts	9834	11,134	3340
Inserted Basepairs	16,527,995	12,902,410	5,160,503

290  
291           The distribution of PAV sizes indicated that large PAVs were rare across the *circum-*  
292 basmati genomes, while PAVs < 500 bps in size were the most common (Figure 3A). Within  
293 smaller-sized PAVs those in the 200-500 bp size range showed a peak in abundance. A closer  
294 examination revealed that sequence positions of more than 75% of these 200-500 bp-sized PAVs  
295 overlapped with transposable element coordinates in the *circum-basmati* genomes (Supplemental  
296 Table 6). A previous study based on short-read Illumina sequencing data reported a similar  
297 enrichment of short repetitive elements such as the long terminal repeats (LTRs) of  
298 retrotransposons, *Tc1/mariner* elements, and *mPing* elements among PAVs in this size range  
299 [61].

300           PAVs shorter than 200 bps also overlapped with repetitive sequence positions in the  
301 *circum-basmati* genomes, but the relative abundance of each repeat type differed among  
302 insertion and deletion variants. Insertions in the Basmati 334 and Dom Sufid genomes had a  
303 higher relative abundance of simple sequence repeats (*i.e.* microsatellites) compared to deletions  
304 (Supplemental Table 6). These inserted simple sequence repeats were highly enriched for (AT)<sub>n</sub>  
305 dinucleotide repeats, which in Basmati 334 accounted for 66,624 bps out of a total of 72,436 bps  
306 (92.0%) of simple sequence repeats, and for Dom Sufid 56,032 bps out of a total of 63,127 bps  
307 (88.8%).

308           Between the Basmati 334 and Dom Sufid genomes, ~45% of PAVs had overlapping  
309 genome coordinates (Figure 3B) suggesting that variety-specific insertion and deletion  
310 polymorphisms were common. We plotted PAVs for each of our *circum-basmati* genomes to  
311 visualize their distribution (Figure 3C). Chromosome-specific differences in the distribution of  
312 PAVs were seen for each *circum-basmati* genome: in Basmati 334, for example, chromosome 1



313 had the lowest density of PAVs, while in Dom Sufid this was the case for chromosome 2  
314 (Supplemental Figure 5). On the other hand, both genomes showed significantly higher densities  
315 of PAVs on chromosome 10 (Tukey's range test  $P < 0.05$ ). This suggested that, compared to  
316 Nipponbare, chromosome 10 was the most differentiated in terms of insertion and deletion  
317 variations in both of our *circum-basmati* genomes.

318

319 **Evolution of *circum-basmati* rice involved group-specific gene deletions.** The proportion of  
320 repeat sequences found within the larger-sized PAVs (*i.e.* those  $> 2$  kb) was high, where between  
321 84% and 98% of large PAVs contained transposable element-related sequences (Supplemental  
322 Table 6). Regardless, these larger PAVs also involved loss or gain of coding sequences. For  
323 instance, gene ontology analysis of domesticated rice gene orthogroups showed enrichment for  
324 genes related to electron transporter activity among both *circum-basmati*-specific gene losses and  
325 gains (see Supplemental Table 7 for gene ontology results for *circum-basmati*-specific gene  
326 losses and Supplemental Table 8 for gene ontology results for *circum-basmati*-specific gene  
327 gains).

328 Many of these genic PAVs could have been important during the rice domestication  
329 process [11]. Gene deletions, in particular, are more likely to have a functional consequence than  
330 single-nucleotide polymorphisms or short indels and may underlie drastic phenotypic variation.  
331 In the context of crop domestication and diversification this could have led to desirable  
332 phenotypes in human-created agricultural environments. For instance, several domestication  
333 phenotypes in rice are known to be caused by gene deletions [35, 62–66].

334 There were 873 gene orthogroups for which neither of the *circum-basmati* genomes had a  
335 gene member, but for which genomes for all three other rice variety groups (N22, Nipponbare,

336 and R498) had at least one gene member. Among these, there were 545 orthogroups for which  
337 N22, Nipponbare, and R498 each had a single-copy gene member, suggesting that the deletion of  
338 these genes in both the Basmati 334 and Dom Sufid genomes could have had a major effect in  
339 *circum*-basmati. We aligned Illumina sequencing data from our *circum*-basmati population  
340 dataset to the japonica Nipponbare genome, and calculated deletion frequencies of Nipponbare  
341 genes that belonged to the 545 orthogroups (see Supplemental Table 9 for gene deletion  
342 frequencies in the *circum*-basmati population for the Nipponbare genes that are missing in  
343 Basmati 334 and Dom Sufid). The vast majority of these Nipponbare genes (509 orthogroups or  
344 93.4%) were entirely absent in the *circum*-basmati population, further indicating that these were  
345 *circum*-basmati-specific gene deletions fixed within this variety group.

346         One of the genes specifically deleted in *circum*-basmati rice varieties was *Awn3-1*  
347 (Os03g0418600), which was identified in a previous study as associated with altered awn length  
348 in japonica rice [67]. Reduced awn length is an important domestication trait that was selected  
349 for ease of harvesting and storing rice seeds [68]. This gene was missing in both *circum*-basmati  
350 genomes and no region could be aligned to the Nipponbare *Awn3-1* genic region (Figure 2C).  
351 Instead of the *Awn3-1* coding sequence, this genomic region contained an excess of transposable  
352 element sequences, suggesting an accumulation of repetitive DNA may have been involved in  
353 this gene's deletion. The flanking arms upstream and downstream of Os03g0418600 were  
354 annotated in both *circum*-basmati genomes and were syntenic to the regions in both Nipponbare  
355 and N22. These flanking arms, however, were also accumulating transposable element  
356 sequences, indicating that this entire genomic region may be degenerating in both *circum*-  
357 basmati rice genomes.

358

359 **Repetitive DNA and retrotransposon dynamics in the *circum-basmati* genomes.** Repetitive  
360 DNA makes up more than 44% of the Basmati 334 and Dom Sufid genome assemblies (Table 2).  
361 Consistent with genomes of other plant species [69] the repetitive DNA was largely composed of  
362 Class I retrotransposons, followed by Class II DNA transposons (Figure 4A). In total, 171.1 Mb  
363 were annotated as repetitive for Basmati 334, and 169.5 Mb for Dom Sufid. The amount of  
364 repetitive DNA in the *circum-basmati* genomes was higher than in the Nipponbare (160.6 Mb)  
365 and N22 genomes (152.1 Mb), but lower than in the indica R498 (175.9 Mb) and IR8 (176.0 Mb)  
366 genomes. These differences in the total amount of repetitive DNA were similar to overall  
367 genome assembly size differences (Table 3), indicating that variation in repeat DNA  
368 accumulation is largely driving genome size differences in rice [70].

369 We focused our attention on retrotransposons, which made up the majority of the rice  
370 repetitive DNA landscape (Figure 4A). Using *LTRharvest* [71, 72], we identified and *de novo*-  
371 annotated LTR retrotransposons in the *circum-basmati* genomes. *LTRharvest* annotated 5,170  
372 and 5,150 candidate LTR retrotransposons in Basmati 334 and Dom Sufid, respectively  
373 (Supplemental Tables 10 and 11). Of these, 4,180 retrotransposons (80.9% of all candidate LTR  
374 retrotransposons) in Basmati 334 and 4,228 (82.1%) in Dom Sufid were classified as LTR  
375 retrotransposons by *RepeatMasker's RepeatClassifier* tool (<http://www.repeatmasker.org>). Most  
376 LTR retrotransposons were from the *gypsy* and *copia* superfamilies [73, 74], which made up  
377 77.1% (3,225 *gypsy* elements) and 21.9% (915 *copia* elements) of LTR retrotransposons in the  
378 Basmati 334 genome, and 76.4% (3,231 *gypsy* elements) and 22.8% (962 *copia* elements) of  
379 LTR retrotransposons in the Dom Sufid genome, respectively. Comparison of LTR  
380 retrotransposon content among reference genomes from different rice variety groups  
381 (Supplemental Figure 4) revealed that genomes assembled to near completion (*i.e.* Nipponbare,

382 N22, Basmati 334, Dom Sufid, and indica varieties IR8 and R498, as well as MH63 and ZS97  
383 [40]) had higher numbers of annotated retrotransposons than genomes generated from short-read  
384 sequencing data (GP295-1, *circum*-aus varieties DJ123 [38] and Kasalath [39], and indica variety  
385 IR64 [38]), suggesting genome assemblies from short-read sequencing data may be missing  
386 certain repetitive DNA regions.

387         Due to the proliferation mechanism of LTR transposons, the DNA divergence of an LTR  
388 sequence can be used to approximate the insertion time for an LTR retrotransposon [75].  
389 Compared to other rice reference genomes, the insertion times for the Basmati 334 and Dom  
390 Sufid LTR retrotransposons were most similar to those observed for elements in the *circum*-aus  
391 N22 genome (Supplemental Figure 4). Within our *circum*-basmati assemblies, the *gypsy*  
392 superfamily elements had a younger average insertion time (~2.2 million years ago) than  
393 elements of the *copia* superfamily (~2.7 million years ago; Figure 4B).

394         Concentrating on *gypsy* and *copia* elements with the *rve* (integrase; Pfam ID: PF00665)  
395 gene, we examined the evolutionary dynamics of these LTR retrotransposons by reconstructing  
396 their phylogenetic relationships across reference genomes for the four domesticated rice variety  
397 groups (N22, Basmati 334, Dom Sufid, R498, IR8, and Nipponbare), and the two wild rice  
398 species (*O. nivara* and *O. rufipogon*; Fig 3C). The retrotransposons grouped into distinct  
399 phylogenetic clades, which likely reflect repeats belonging to the same family or subfamily [76].  
400 The majority of phylogenetic clades displayed short external and long internal branches,  
401 consistent with rapid recent bursts of transposition observed across various rice LTR  
402 retrotransposon families [77].

403         The *gypsy* and *copia* superfamilies each contained a clade in which the majority of  
404 elements originated within *O. sativa*, present only among the four domesticated rice variety

405 groups (Figure 4C, single star; see Supplemental Tables 12 and 13 for their genome coordinates).  
406 Elements in the *gypsy* superfamily phylogenetic clade had sequence similarity (963 out of the  
407 1,837 retrotransposons) to elements of the *hopi* family [78], while elements in the *copia*  
408 superfamily phylogenetic clade had sequence similarity (88 out of the 264) to elements in the  
409 *osr4* family [79]. Elements of the *hopi* family are found in high copy number in genomes of  
410 domesticated rice varieties [80] and this amplification has happened recently [81].

411         Several retrotransposon clades were restricted to certain rice variety groups. The *gypsy*  
412 superfamily harbored a phylogenetic clade whose elements were only present in genomes of  
413 *circum*-aus, *circum*-basmati, and indica varieties (Figure 4C, double star; see Supplemental  
414 Table 14 for their genome coordinates), while we observed a clade comprised mostly of *circum*-  
415 basmati-specific elements within the *copia* superfamily (Figure 4C, triple star; see Supplemental  
416 Table 15 for their genome coordinates). Only a few members of the *gypsy*-like clade had  
417 sequence similarity (7 out of 478) to elements of the *rire3* [82] and *m215* [83] families.  
418 Members of both families are known to be present in high copy numbers in genomes of  
419 domesticated rice varieties, but their abundance differs between the japonica and indica variety  
420 groups [80], suggesting a *rire3*- or *m215*-like element expansion in the *circum*-aus, *circum*-  
421 basmati, and indica genomes. A majority of the *circum*-basmati-specific *copia*-like elements had  
422 sequence similarity (109 out of 113) to members of the *houba* family [78], which are found in  
423 high copy numbers in certain individuals, but in lower frequency across the rice population [80].  
424 This suggests the *houba* family might have undergone a recent expansion specifically within the  
425 *circum*-basmati genomes.

426

427 **Phylogenomic analysis on the origins of *circum*-basmati rice.** We estimated the phylogenetic  
428 relationships within and between variety groups of domesticated Asian rice. Our maximum-  
429 likelihood phylogenetic tree, based on four-fold degenerate sites from the Nipponbare coding  
430 sequences (Figure 5A), showed that each cultivar was monophyletic with respect to its variety  
431 group of origin. In addition, the *circum*-basmati group was sister to japonica rice, while the  
432 *circum*-aus group was sister to indica. Consistent with previous observations, the wild rices *O.*  
433 *nivara* and *O. rufipogon* were sister to the *circum*-aus and japonica rices, respectively [14].  
434 While this suggests that each domesticated rice variety group may have had independent wild  
435 progenitors of origin, it should be noted that recent hybridization between wild and domesticated  
436 rice [84, 85] could lead to similar phylogenetic relationships.

437 To further investigate phylogenetic relationships between *circum*-basmati and japonica,  
438 we examined phylogenetic topologies of each gene involving the trio Basmati 334, Nipponbare,  
439 and *O. rufipogon*. For each gene we tested which of three possible topologies for a rooted three-  
440 species tree - *i.e.* [(P1, P2), P3], O, where O is outgroup *O. barthii* and P1, P2, and P3 are  
441 Basmati 334 (or Dom Sufid), Nipponbare, and *O. rufipogon*, respectively - were found in highest  
442 proportion. For the trio involving Basmati 334, Nipponbare, and *O. rufipogon* there were 7,581  
443 genes (or 32.6%), and for the trio involving Dom Sufid, Nipponbare, and *O. rufipogon* there  
444 were 7,690 genes (or 33.1%), that significantly rejected one topology over the other two using an  
445 Approximately Unbiased (AU) topology test [86]. In both trios, the majority of those genes  
446 supported a topology that grouped *circum*-basmati and Nipponbare as sister to each other (Figure  
447 5B; 3,881 [or 51.2%] and 4,407 [or 57.3%] genes for Basmati 334 and Dom Sufid, respectively).  
448 A lower number of genes (3,018 [or 39.8%] and 2,508 [or 32.6%] genes for Basmati 334 and

449 Dom Sufid, respectively) supported the topology that placed Nipponbare and *O. rufipogon*  
450 together.

451 The topology test result suggested that [(japonica, *circum-basmati*), *O. rufipogon*] was  
452 the true species topology, while the topology [(japonica, *O. rufipogon*), *circum-basmati*]  
453 represented possible evidence of admixture (although it could also arise from incomplete lineage  
454 sorting). To test for introgression, we employed D-statistics from the ABBA-BABA test [87, 88].  
455 The D-statistics for the topology [(japonica, *circum-basmati*), *O. rufipogon*] were significantly  
456 negative - Figure 5C left panel; z-score = -14.60 and  $D \pm s.e = -0.28 \pm 0.019$  for topology  
457 [(Nipponbare, Basmati 334), *O. rufipogon*], and z-score = -9.09 and  $D = -0.20 \pm 0.022$  for  
458 topology [(Nipponbare, Dom Sufid), *O. rufipogon*] - suggesting significant evidence of  
459 admixture between japonica and *O. rufipogon*.

460 Our initial topology test suggested that the trio involving Dom Sufid, Nipponbare, and *O.*  
461 *rufipogon* had a higher proportion of genes supporting the [(*circum-basmati*, japonica), *O.*  
462 *rufipogon*] topology compared to the trio involving Basmati 334, Nipponbare, and *O. rufipogon*  
463 (Figure 5B). This suggested within-population variation in the amount of japonica or *O.*  
464 *rufipogon* ancestry across the *circum-basmati* genomes due to differences in gene flow. We  
465 conducted ABBA-BABA tests involving the topology [(Basmati 334, Dom Sufid), Nipponbare  
466 or *O. rufipogon*] to examine the differences in introgression between the *circum-basmati* and  
467 japonica or *O. rufipogon* genomes. The results showed significantly positive D-statistics for the  
468 topology [(Basmati 334, Dom Sufid), Nipponbare] (Figure 5C left panel; z-score = 8.42 and  $D =$   
469  $0.27 \pm 0.032$ ), indicating that Dom Sufid shared more alleles with japonica than Basmati 334 did  
470 due to a history of more admixture with japonica. The D-statistics involving the topology  
471 [(Basmati 334, Dom Sufid), *O. rufipogon*] were also significantly positive (Figure 5C left panel;

472 z-score = 5.57 and  $D = 0.21 \pm 0.038$ ). While this suggests admixture between Dom Sufid and *O.*  
473 *rufipogon*, it may also be an artifact due to the significant admixture between japonica and *O.*  
474 *rufipogon*.

475  
476 **Signatures of admixture between *circum-basmati* and *circum-aus* rice genomes.** Due to  
477 extensive admixture between rice variety group genomes [14] we examined whether the basmati  
478 genome was also influenced by gene flow with other divergent rice variety groups (*i.e.* *circum-*  
479 *aus* or *indica* rices). A topology test was conducted for a rooted, three-population species tree.  
480 For the trio involving Basmati 334, *circum-aus* variety N22, and *indica* variety R498 there were  
481 7,859 genes (or 35.3%), and for the trio involving Dom Sufid, N22, and R498 there were 8,109  
482 genes (or 37.8%), that significantly rejected one topology over the other two after an AU test. In  
483 both trios, more than half of the genes supported the topology grouping *circum-aus* and *indica* as  
484 sisters (Figure 5D). In addition, more genes supported the topology grouping *circum-aus* and  
485 *circum-basmati* as sisters than the topology grouping *indica* and *circum-basmati* as sisters. This  
486 suggested that the *circum-aus* variety group might have contributed a larger proportion of genes  
487 to *circum-basmati* through gene flow than the *indica* variety group did.

488 To test for evidence of admixture, we conducted ABBA-BABA tests involving trios of  
489 the *circum-basmati*, N22, and R498 genomes. Results showed significant evidence of gene flow  
490 between *circum-aus* and both *circum-basmati* genomes - Figure 5C, right panel; z-score = 5.70  
491 and  $D = 0.082 \pm 0.014$  for topology [(R498, N22), Basmati 334]; and z-score = 8.44 and  $D =$   
492  $0.11 \pm 0.013$  for topology [(R498, N22), Dom Sufid]. To test whether there was variability in the  
493 *circum-aus* or *indica* ancestry in each of the *circum-basmati* genomes, we conducted ABBA-  
494 BABA tests for the topology [(Basmati 334, Dom Sufid), N22 or R498]. Neither of the ABBA-



495 BABA tests involving the topology [(Basmati 334, Dom Sufid), N22] (Figure 5C, right panel; z-  
496 score = 1.20 and  $D = 0.025 \pm 0.021$ ) or the topology [(Basmati 334, Dom Sufid), R498] (Figure  
497 5C, right panel; z-score = -2.24 and  $D = -0.06 \pm 0.026$ ) was significant, suggesting the amount of  
498 admixture from *circum*-aus to each of the two *circum*-basmati genomes was similar.

499 In sum, the phylogenomic analysis indicated that *circum*-basmati and japonica share the  
500 most recent common ancestor, while *circum*-aus has admixed with *circum*-basmati during its  
501 evolutionary history (Figure 5F). We then examined whether admixture from *circum*-aus had  
502 affected each of the *circum*-basmati chromosomes to a similar degree. For both *circum*-basmati  
503 genomes most chromosomes had D-statistics that were not different from the genome-wide D-  
504 statistics value or from zero (Figure 5E). Exceptions were chromosomes 10 and 11, where the  
505 bootstrap D-statistics were significantly higher than the genome-wide estimate.

506

507 **Population analysis on the origin of *circum*-basmati rice.** Since our analysis was based on  
508 single representative genomes from each rice variety group, we compared the results of our  
509 phylogenomic analyses to population genomic patterns in an expanded set of rice varieties from  
510 different groups. We obtained high coverage (>14×) genomic re-sequencing data (generated with  
511 Illumina short-read sequencing) from landrace varieties in the 3K Rice Genome Project [7] and  
512 from *circum*-basmati rice landraces we re-sequenced. In total, we analyzed 24 *circum*-aus, 18  
513 *circum*-basmati, and 37 tropical japonica landraces (see Supplemental Table 16 for variety  
514 names). The raw Illumina sequencing reads were aligned to the scaffolded Basmati 334 genome  
515 and computationally genotyped. A total of 4,594,290 polymorphic sites were called across the  
516 three rice variety groups and used for further analysis.

517 To quantify relationships between *circum-aus*, *circum-basmati*, and *japonica*, we  
518 conducted a topology-weighting analysis [89]. For three populations there are three possible  
519 topologies and we conducted localized sliding window analysis to quantify the number of unique  
520 sub-trees that supported each tree topology. Consistent with the phylogenomic analysis results,  
521 the topology weight was largest for the topology that grouped *japonica* and *circum-basmati* as  
522 sisters (Figure 6A; topology weight = 0.481 with 95% confidence interval [0.479-0.483]). The  
523 topology that grouped *circum-aus* and *circum-basmati* together as sisters weighed significantly  
524 more (topology weight = 0.318 with 95% confidence interval [0.316-0.320]) than the topology  
525 that grouped *japonica* and *circum-aus* as sisters (topology weight = 0.201 with 95% confidence  
526 interval [0.199-0.203]). This was consistent with the admixture results from the comparative  
527 phylogenomic analysis, which detected evidence of gene flow between *circum-aus* and *circum-*  
528 *basmati*.

529 We then examined topology weights for each individual chromosome, since the ABBA-  
530 BABA tests using the genome assemblies had detected variation in *circum-aus* ancestry between  
531 different chromosomes (Figure 5E). The results showed that for most of the chromosomes the  
532 topology [(*japonica*, *circum-basmati*), *circum-aus*] always weighed more than the remaining two  
533 topologies. An exception was observed for chromosome 10 where the topology weight grouping  
534 *circum-aus* and *circum-basmati* as sisters was significantly higher (topology weight = 0.433 with  
535 95% confidence interval [0.424-0.442]) than the weight for the genome-wide topology that  
536 grouped *japonica* and *circum-basmati* as sisters (topology weight = 0.320 with 95% confidence  
537 interval [0.312-0.328]). This change in predominant topology was still observed when the  
538 weights were calculated across wider local windows (Supplemental Figure 6). Another exception  
539 could be seen for chromosome 6 where the genome-wide topology [(*japonica*, *circum-basmati*),

540 *circum-aus*] (topology weight = 0.367 with 95% confidence interval [0.359-0.374) and the  
541 admixture topology [*circum-aus*, *circum-basmati*), japonica] (topology weight = 0.355 with 95%  
542 confidence interval [0.349-0.362]) had almost equal weights. In larger window sizes the weight  
543 of the admixed topology was slightly higher than that of the genome-wide topology  
544 (Supplemental Figure 6).

545 To estimate the evolutionary/domestication scenario that might explain the observed  
546 relationships between the *circum-aus*, *circum-basmati*, and japonica groups, we used the  
547 diffusion-based approach of the program  $\delta a d i$  [90] and fitted specific demographic models to the  
548 observed allele frequency spectra for the three rice variety groups. Because all three rice groups  
549 have evidence of admixture with each other [7, 9, 14, 16] we examined 13 demographic  
550 scenarios involving symmetric, asymmetric, and “no migration” models between variety groups,  
551 with and without recent population size changes (Supplemental Figure 7). To minimize the effect  
552 of genetic linkage on the demography estimation, polymorphic sites were randomly pruned in  
553 200 kb windows, resulting in 1,918 segregating sites. The best-fitting demographic scenario was  
554 one that modeled a period of lineage splitting and isolation, while gene flow only occurred after  
555 formation of the three populations and at a later time (Figure 6C; visualizations of the 2D site  
556 frequency spectrum and model fit can be seen in Supplemental Figure 8). This best-fitting model  
557 was one of the lesser-parameterized models we tested, and the difference in Akaike Information  
558 Criterion ( $\Delta A I C$ ) with the model with the second-highest likelihood was 25.46 (see  
559 Supplemental Table 17 for parameter estimates and maximum likelihood estimates for each  
560 demographic model).

561

562 **Genetic structure within the *circum-basmati* group.** We used the *circum-basmati* population  
563 genomic data for the 78 varieties aligned to the scaffolded Basmati 334 genome, and called the  
564 polymorphic sites segregating within this variety group. After filtering, a total of 4,430,322 SNPs  
565 across the *circum-basmati* dataset remained, which were used to examine population genetic  
566 relationships within *circum-basmati*.

567 We conducted principal component analysis (PCA) using the polymorphism data and  
568 color-coded each *circum-basmati* rice variety according to its country of origin (Figure 7A). The  
569 PCA suggested that *circum-basmati* rices could be divided into three major groups with clear  
570 geographic associations: (Group 1) a largely Bhutan/Nepal-based group, (Group 2) an  
571 India/Bangladesh/Myanmar-based group, and (Group 3) an Iran/Pakistan-based group. The rice  
572 varieties that could not be grouped occupied an ambiguous space across the principal  
573 components, suggesting these might represent admixed rice varieties.

574 To obtain better insight into the ancestry of each rice variety, we used *fastSTRUCTURE*  
575 [91] and varied assumed ancestral population (K) from 2 to 5 groups so the ancestry proportion  
576 of each rice variety could be estimated (Figure 7B). At K=2, the India/Bangladesh/Myanmar and  
577 Iran/Pakistan rice groups were shown to have distinct ancestral components, while the  
578 Bhutan/Nepal group was largely an admixture of the other two groups. At K=3, the grouping  
579 status designated from the PCA was largely concordant with the ancestral components. At K=4,  
580 most India/Bangladesh/Myanmar rices had a single ancestral component, but Iran/Pakistan rices  
581 had two ancestral components that were shared with several Bhutan/Nepal landraces.

582 Furthermore, several of the cultivars from the latter group seemed to form an admixed group  
583 with India/Bangladesh/Myanmar varieties. In fact, when a phylogenetic tree was reconstructed  
584 using the polymorphic sites, varieties within the India/Bangladesh/Myanmar and Iran/Pakistan

585 groups formed a monophyletic clade with each other. On the other hand, Bhutan/Nepal varieties  
586 formed a paraphyletic group where several clustered with the Iran/Pakistan varieties  
587 (Supplemental Figure 9).

588 In summary, the *circum*-basmati rices have evolved across a geographic gradient with at  
589 least three genetic groups (Figure 7C). These existed as distinct ancestral groups that later  
590 admixed to form several other *circum*-basmati varieties. Group 1 and Group 3 rices in particular  
591 may have experienced greater admixture, while the Group 2 landraces remained genetically more  
592 isolated from other *circum*-basmati subpopulations. We also found differences in agronomic  
593 traits associated with our designated groups (Figure 7D). The grain length to width ratio, which  
594 is a highly prized trait in certain *circum*-basmati rices [24], was significantly larger in Group 3  
595 Iran/Pakistan varieties. The thousand-kernel weights, on the other hand, were highest for Group  
596 2 India/Bangladesh/Myanmar varieties and were significantly higher than those for the  
597 ungrouped and Group 1 Bhutan/Nepal varieties.

598

## 599 **DISCUSSION**

600 Nanopore sequencing is becoming an increasingly popular approach to sequence and  
601 assemble the often large and complex genomes of plants [92–94]. Here, using long-read  
602 sequences generated with Oxford Nanopore Technologies' sequencing platform, we assembled  
603 genomes of two *circum*-basmati rice cultivars, with quality metrics that were comparable to other  
604 rice variety group reference genome assemblies [37, 40, 41]. With modest genome coverage, we  
605 were able to develop reference genome assemblies that represented a significant improvement  
606 over a previous *circum*-basmati reference genome sequence, which had been assembled with a >  
607 3-fold higher genome coverage than ours, but from short-read sequences [42]. With additional

608 short-read sequencing reads, we were able to correct errors from the nanopore sequencing reads,  
609 resulting in two high-quality *circum*-basmati genome assemblies.

610 Even with long-read sequence data, developing good plant reference genome sequences  
611 still requires additional technologies such as optical mapping or Hi-C sequencing for improving  
612 assembly contiguity [95–98], which can be error prone as well [56]. Our assemblies were also  
613 fragmented into multiple contigs, but sizes of these contigs were sufficiently large that we could  
614 use reference genome sequences from another rice variety group to anchor the majority of  
615 contigs and scaffold them to higher-order chromosome-level assemblies. Hence, with a highly  
616 contiguous draft genome assembly, reference genome-based scaffolding can be a cost-efficient  
617 and powerful method of generating chromosome-level assemblies.

618 Repetitive DNA constitutes large proportions of plant genomes [99], and there is an  
619 advantage to using long-read sequences for genome assembly as it enables better annotation of  
620 transposable elements. Many transposable element insertions have evolutionarily deleterious  
621 consequences in the rice genome [54, 100, 101], but some insertions could have beneficial  
622 effects on the host [102]. Using our genome assembly, we have identified retrotransposon  
623 families that have expanded specifically within *circum*-basmati genomes. While more study will  
624 be necessary to understand the functional effects of these insertions, long-read sequences have  
625 greatly improved the assembly and identification of repeat types.

626 Due to a lack of archaeobotanical data, the origins of *circum*-basmati rice have remained  
627 elusive. Studies of this variety group's origins have primarily focused on genetic differences that  
628 exist between *circum*-basmati and other Asian rice variety groups [6, 7]. Recently, a study  
629 suggested that *circum*-basmati rice (called 'aromatic' in that study) was a product of  
630 hybridization between the *circum*-aus and japonica rice variety groups [17]. This inference was

631 based on observations of phylogenetic relationships across genomic regions that showed  
632 evidence of domestication-related selective sweeps. These regions mostly grouped *circum-*  
633 *basmati* with *japonica* or *circum-aus*. In addition, chloroplast haplotype analysis indicated that  
634 most *circum-basmati* varieties carried a chloroplast derived from a wild rice most closely related  
635 to *circum-aus* landraces [103]. Our evolutionary analysis of *circum-basmati* rice genomes  
636 generally supported this view. Although our results suggest that *circum-basmati* had its origins  
637 primarily in *japonica*, we also find significant evidence of gene flow originating from *circum-*  
638 *aus*, which we detected both in comparative genomic and population genomic analyses.  
639 Demographic modeling indicated a period of isolation among *circum-aus*, *circum-basmati*, and  
640 *japonica*, with gene flow occurring only after lineage splitting of each group. Here, our model is  
641 consistent with the current view that gene flow is a key evolutionary process associated with the  
642 diversification of rice [10, 12–14, 16, 104, 105].

643         Interestingly, we found that chromosome 10 of *circum-basmati* had an evolutionary  
644 history that differed significantly from that of other chromosomes. Specifically, compared to  
645 *japonica*, this chromosome had the highest proportion of presence/absence variation, and shared  
646 more alleles with *circum-aus*. Based on this result, we hypothesize that this is largely due to  
647 higher levels of introgression from *circum-aus* into chromosome 10 compared to other  
648 chromosomes. Such a deviation of evolutionary patterns on a single chromosome has been  
649 observed in the *Aquilegia* genus [106], but to our knowledge has not been observed elsewhere.  
650 Why this occurred is unclear at present, but it may be that selection has driven a higher  
651 proportion of *circum-aus* alleles into chromosome 10. Future work will be necessary to clarify  
652 the consequence of this higher level of admixture on chromosome 10.

653           Very little is known about population genomic diversity within *circum*-basmati. Our  
654 analysis suggests the existence of at least three genetic groups within this variety group, and  
655 these groups showed geographic structuring. Several varieties from Group 1 (Bhutan/Nepal) and  
656 Group 3 (Iran/Pakistan) had population genomic signatures consistent with an admixed  
657 population, while Group 2 (India/Bangladesh/Myanmar) was genetically more distinct from the  
658 other two subpopulations. In addition, the geographic location of the India/Bangladesh/Myanmar  
659 group largely overlaps the region where *circum*-aus varieties were historically grown [107, 108].  
660 Given the extensive history of admixture that *circum*-basmati rices have with *circum*-aus, the  
661 India/Bangladesh/Myanmar group may have been influenced particularly strongly by gene flow  
662 from *circum*-aus. How these three genetic subpopulations were established may require a deeper  
663 sampling with in-depth analysis, but the geographically structured genomic variation shows that  
664 the diversity of *circum*-basmati has clearly been underappreciated. In addition, the Basmati 334  
665 and Dom Sufid varieties, for which we generated genome assemblies in this study, both belong to  
666 the Iran/Pakistan genetic group. Thus, our study still leaves a gap in our knowledge of genomic  
667 variation in the Bhutan/Nepal and India/Bangladesh/Myanmar genetic groups, and varieties in  
668 these groups would be obvious next targets for generating additional genome assemblies.

669

## 670 **CONCLUSIONS**

671           In conclusion, our study shows that generating high-quality plant genome assemblies is  
672 feasible with relatively modest amounts of resources and data. Using nanopore sequencing, we  
673 were able to produce contiguous, chromosome-level genome assemblies for cultivars in a rice  
674 variety group that contains economically and culturally important varieties. Our reference  
675 genome sequences have the potential to be important genomic resources for identifying single



676 nucleotide polymorphisms and larger structural variations that are unique to *circum*-basmati rice.  
677 Analyzing *de novo* genome assemblies for a larger sample of Asian rice will be important for  
678 uncovering and studying hidden population genomic variation too complex to study with only  
679 short-read sequencing technology.

680

## 681 **MATERIALS AND METHODS**

682 **Plant material.** Basmati 334 (IRGC 27819; GeneSys passport:

683 <https://purl.org/germplasm/id/23601903-f8c3-4642-a7fc-516a5bc154f7>) is a basmati (*sensu*

684 *stricto*) landrace from Pakistan and was originally donated to the International Rice Research

685 Institute (IRRI) by the Agricultural Research Council (ARC) in Karachi (donor accession ID:

686 PAK. SR. NO. 39). Dom Sufid (IRGC 117265; GeneSys passport:

687 <https://purl.org/germplasm/id/fb861458-09de-46c4-b9ca-f5c439822919>) is a sadri landrace from

688 Iran. Seeds from accessions IRGC 27819 and IRGC 117265 were obtained from the IRRI seed

689 bank, surface-sterilized with bleach, and germinated in the dark on a wet paper towel for four

690 days. Seedlings were transplanted individually in pots containing continuously wet soil in a

691 greenhouse at New York University's Center for Genomics and Systems Biology and cultivated

692 under a 12h day-12h night photoperiod at 30°C. Plants were kept in the dark in a growth cabinet

693 under the same climatic conditions for four days prior to tissue harvesting. Continuous darkness

694 induced chloroplast degradation, which diminishes the amount of chloroplast DNA that would

695 otherwise end up in the DNA extracted from the leaves.

696

697 **DNA extractions.** Thirty-six 100-mg samples (3.6 g total) of leaf tissue from a total of 10 one-

698 month-old plants were flash-frozen at harvest for each accession and stored at -80°C. DNA

699 extractions were performed by isolating the cell nuclei and gently lysing the nuclei to extract  
700 intact DNA molecules [109]. Yields ranged between 140ng/ul and 150ng/ul.

701

702 **Library preparation and nanopore sequencing.** Genomic DNA was visualized on an agarose  
703 gel to determine shearing. DNA was size-selected using BluePippin BLF7510 cassette (Sage  
704 Science) and high-pass mode (>20 kb) and prepared using Oxford Nanopore Technologies'  
705 standard ligation sequencing kit SQK-LSK109. FLO-MIN106 (R9.4) flowcells were used for  
706 sequencing on the GridION X5 platform.

707

708 **Library preparation and Illumina sequencing.** Extracted genomic DNA was prepared for  
709 short-read sequencing using the Illumina Nextera DNA Library Preparation Kit. Sequencing was  
710 done on the Illumina HiSeq 2500 – HighOutput Mode v3 with 2×100 bp read configuration, at  
711 the New York University Genomics Core Facility.

712

713 **Genome assembly, polishing, and scaffolding.** After completion of sequencing, the raw signal  
714 intensity data was used for base calling using *flip flop* (version 2.3.5) from Oxford Nanopore  
715 Technologies. Reads with a mean qscore (quality) greater than 8 and a read length greater than 8  
716 kb were used, and trimmed for adaptor sequences using *Porechop*  
717 (<https://github.com/rrwick/Porechop>). Raw nanopore sequencing reads were corrected using the  
718 program *Canu* [110], and then assembled with the genome assembler *Flye* [111].

719 The initial draft assemblies were polished for three rounds using the raw nanopore reads  
720 with *Racon* ver. 1.2.1 [112], and one round with *Medaka*  
721 (<https://github.com/nanoporetech/medaka>) from Oxford Nanopore Technologies. Afterwards,

722 reads from Illumina sequencing were used by *bwa-mem* [113] to align to the draft genome  
723 assemblies. The alignment files were then used by *Pilon* ver. 1.22 [114] for three rounds of  
724 polishing.

725 Contigs were scaffolded using a reference genome-guided scaffolding approach  
726 implemented in *RaGOO* [56]. Using the Nipponbare genome as a reference, we aligned the  
727 *circum-basmati* genomes using *Minimap2* [115]. *RaGOO* was then used to order the assembly  
728 contigs. Space between contigs was artificially filled in with 100 ‘N’ blocks.

729 Genome assembly statistics were calculated using the *bbmap stats.sh* script from the  
730 *BBTools* suite (<https://jgi.doe.gov/data-and-tools/bbtools/>). Completeness of the genome  
731 assemblies was evaluated using *BUSCO* ver. 2.0 [116]. Synteny between the *circum-basmati*  
732 genomes and the Nipponbare genome was visualized using *D-GENIES* [117]. Genome-wide  
733 dotplot from *D-GENIES* indicated the initial genome assembly of Dom Sufid had an evidence of  
734 a large chromosomal fusion between the ends of chromosome 4 and 10. Closer examination of  
735 this contig (named contig\_28 of Dom Sufid) showed the break point overlapped the telomeric  
736 repeat sequence, indicating there had been a misassembly between the ends of chromosome 4  
737 and 10. Hence, contig\_28 was broken up into two so that each contig represented the respective  
738 chromosome of origin, and were then subsequently scaffolded using *RaGOO*.

739 Inversions that were observed in the dot plot were computationally verified  
740 independently using raw nanopore reads. The long read-aware aligner *ngmlr* [55] was used to  
741 align the nanopore reads to the Nipponbare genome, after which the long read-aware structural  
742 variation caller *sniffles* [55] was used to call and detect inversions.

743 The number of sites aligning to the Nipponbare genome was determined using the  
744 *Mummer4* package [118]. Alignment delta files were analyzed with the *dnadiff* suite from the

745 *Mummer4* package to calculate the number of aligned sites, and the number of differences  
746 between the Nipponbare genome and the *circum*-basmati genomes.

747  
748 **Gene annotation and analysis.** Gene annotation was conducted using the *MAKER* program [52,  
749 53]. An in-depth description of running *MAKER* can be found on the website:  
750 <https://gist.github.com/darencard/bb1001ac1532dd4225b030cf0cd61ce2>. We used published  
751 *Oryza* genic sequences as evidence for the gene modeling process. We downloaded the  
752 Nipponbare cDNA sequences from RAP-DB (<https://rapdb.dna.affrc.go.jp/>) to supply as EST  
753 evidence, while the protein sequences from the 13 *Oryza* species project [37] were used as  
754 protein evidence for the *MAKER* pipeline. Repetitive regions identified from the repeat analysis  
755 were used to mask out the repeat regions for this analysis. After a first round of running *MAKER*  
756 the predicted genes were used by *SNAP* [119] and *Augustus* [120] to create a training dataset of  
757 gene models, which was then used for a second round of *MAKER* gene annotation.

758 Orthology between the genes from different rice genomes was determined with  
759 *Orthofinder* ver. 1.1.9 [59]. Ortholog statuses were visualized with the *UpSetR* package [121].

760 Gene ontology for the orthogroups that are missing specifically in the *circum*-basmati  
761 were examined by using the japonica Nipponbare gene, and conducting a gene ontology  
762 enrichment analysis on *agriGO* v2.0 [122]. Gene ontology enrichment analysis for the *circum*-  
763 basmati specific orthogroups was conducted first by predicting the function and gene ontology of  
764 each *circum*-basmati genome gene model using the eggnoG pipeline [123]. We required an  
765 ontology to have more than 10 genes as a member for further consideration, and enrichment was  
766 tested through a hypergeometric test using the *GOstat* package [124].

767

768 **Repetitive DNA annotation.** The repeat content of each genome assembly was determined  
769 using *Repeatmasker* ver. 4.0.5 (<http://www.repeatmasker.org/RMDownload.html>). We used the  
770 *Oryza*-specific repeat sequences that were identified from Choi et al. [14] (DOI:  
771 10.5061/dryad.7cr0q), who had used *Repeatmodeler* ver. 1.0.8  
772 (<http://www.repeatmasker.org/RepeatModeler.html>) to *de novo*-annotate repetitive elements  
773 across wild and domesticated *Oryza* genomes [37].

774 LTR retrotransposons were annotated using the program *LTRharvest* [125] with  
775 parameters adapted from [126]. LTR retrotransposons were classified into superfamilies [76]  
776 using the program *RepeatClassifier* from the *RepeatModeler* suite. Annotated LTR  
777 retrotransposons were further classified into specific families using the 242 consensus sequences  
778 of LTR-RTs from the RetrOryza database [83]. We used *blastn* [127] to search the RetrOryza  
779 sequences, and each of our candidate LTR retrotransposons was identified using the “80-80-80”  
780 rule [76]: two TEs belong to the same family if they were 80% identical over at least 80 bp and  
781 80% of their length.

782 Insertion times for the LTR retrotransposons were estimated using the DNA divergence  
783 between pairs of LTR sequences [75]. The L-INS-I algorithm in the alignment program *MAFFT*  
784 ver. 7.154b [128] was used to align the LTR sequences. *PAML* ver. 4.8 [129] was used to  
785 estimate the DNA divergence between the LTR sequences with the Kimura-2-parameter base  
786 substitution model [130]. DNA divergence was converted to divergence time (i.e. time since the  
787 insertion of a LTR retrotransposon) approximating a base substitution rate of  $1.3 \times 10^{-8}$  [131],  
788 which is two times higher than the synonymous site substitution rate.

789

790 **Presence/absence variation detection.** PAVs between the Nipponbare genome and the *circum-*  
791 Basmati assemblies were detected using the *Assemblytics* suites [60]. Initially, the Nipponbare  
792 genome was used as the reference to align the *circum*-basmati assemblies using the program  
793 *Minimap2*. The resulting SAM files were converted to files in delta format using the  
794 *sam2delta.py* script from the *RaGOO* suite. The delta files were then uploaded onto the online  
795 *Assemblytics* analysis pipeline (<http://assemblytics.com/>). Repetitive regions would cause  
796 multiple regions in the Nipponbare or *circum*-basmati genomes to align to one another, and in  
797 that case *Assemblytics* would call the same region as a PAV multiple times. Hence, any PAV  
798 regions that overlapped for at least 70% of their genomic coordinates were collapsed to a single  
799 region.

800 The combination of *ngmlr* and *sniffles* was also used to detect the PAVs that differed  
801 between the Nipponbare genome and the raw nanopore reads for the *circum*-basmati rices.  
802 Because *Assemblytics* only detects PAVs in the range of 50 bp to 100,000 bp, we used this  
803 window as a size limit to filter out the PAVs called by *sniffles*. Only PAVs supported by more  
804 than 5 reads by *sniffles* were analyzed.

805 *Assemblytics* and *sniffles* call the breakpoints of PAVs differently. *Assemblytics* calls a  
806 single-best breakpoint based on the genome alignment, while *sniffles* calls a breakpoint across a  
807 predicted interval. To find overlapping PAVs between *Assemblytics* and *sniffles* we added 500 bp  
808 upstream and downstream of the *Assemblytics*-predicted breakpoint positions.

809

810 **Detecting gene deletions across the *circum*-basmati population.** Genome-wide deletion  
811 frequencies of each gene were estimated using the 78-variety *circum*-basmati population  
812 genomic dataset. For each of the 78 varieties, raw sequencing reads were aligned to the *circum-*

813 basmati and Nipponbare genomes using *bwa-mem*. Genome coverage per site was calculated  
814 using *bedtools genomecov* [132]. For each variety the average read coverage was calculated for  
815 each gene, and a gene was designated as deleted if its average coverage was less than 0.05×.

816

817 **Whole-genome alignment of *Oryza* genomes assembled *de novo*.** Several genomes from  
818 published studies that were assembled *de novo* were analyzed. These include domesticated Asian  
819 rice genomes from the japonica variety group cv. Nipponbare [33]; the indica variety group cvs.  
820 93-11 [32], IR8 [37], IR64 [38], MH63 [40], R498 [41], and ZS97 [40]; the *circum*-aus variety  
821 group cvs. DJ123 [38], Kasalath [39], and N22 [37]; and the *circum*-basmati variety group cv.  
822 GP295-1 [42]. Three genomes from wild rice species were also analyzed; these were *O. barthii*  
823 [35], *O. nivara* [37], and *O. rufipogon* [37].

824 Alignment of the genomes assembled *de novo* was conducted using the approach outlined  
825 in Haudry *et al.* [133], and the alignment has been used in another rice comparative genomic  
826 study [14]. Briefly, this involved using the Nipponbare genome as the reference for aligning all  
827 other genome assemblies. Alignment between japonica and a query genome was conducted using  
828 *LASTZ* ver. 1.03.73 [134], and the alignment blocks were chained together using the UCSC Kent  
829 utilities [135]. For japonica genomic regions with multiple chains, the chain with the highest  
830 alignment score was chosen as the single-most orthologous region. This analyzes only one of the  
831 multiple regions that are potentially paralogous between the japonica and query genomes, but  
832 this was not expected to affect the downstream phylogenomic analysis of determining the origin  
833 and evolution of the *circum*-basmati rice variety group. All pairwise genome alignments between  
834 the japonica and query genomes were combined into a multi-genome alignment using *MULTIZ*  
835 [136].

836

837 **Phylogenomic analysis.** The multi-genome alignment was used to reconstruct the phylogenetic  
838 relationships between the domesticated and wild rices. Four-fold degenerate sites based on the  
839 gene model of the reference japonica genome were extracted using the *msa\_view* program from  
840 the *phast* package ver. 1.4 [137]. The four-fold degenerate sites were used by *RAxML* ver. 8.2.5  
841 [138] to build a maximum likelihood-based tree, using a general time-reversible DNA  
842 substitution model with gamma-distributed rate variation.

843 To investigate the genome-wide landscape of introgression and incomplete lineage  
844 sorting we examined the phylogenetic topologies of each gene [139]. For a three-species  
845 phylogeny using *O. barthii* as an outgroup there are three possible topologies. For each gene,  
846 topology-testing methods [140] can be used to determine which topology significantly fits the  
847 gene of interest [14]. *RAxML*-estimated site-likelihood values were calculated for each gene and  
848 the significant topology was determined using the Approximately Unbiased (AU) test [86] from  
849 the program *CONSEL* v. 0.20 [141]. Genes with AU test results with a likelihood difference of  
850 zero were omitted and the topology with an AU test support of greater than 0.95 was selected.

851

852 **Testing for evidence of admixture.** Evidence of admixture between variety groups was detected  
853 using the ABBA-BABA test D-statistics [87, 88]. In a rooted three-taxon phylogeny [i.e.  
854 “((P1,P2),P3),O” where P1, P2, and P3 are the variety groups of interest and O is outgroup *O.*  
855 *barthii*], admixture can be inferred from the combination of ancestral (“A”) and derived (“B”)   
856 allelic states of each individual. The ABBA conformation arises when variety groups P2 and P3  
857 share derived alleles, while the BABA conformation is found when P1 and P3 share derived  
858 alleles. The difference in the frequency of the ABBA and BABA conformations is measured by



859 the D-statistics, where significantly positive D-statistics indicate admixture between the P2 and  
860 P3 variety groups, and significantly negative D-statistics indicate admixture between the P1 and  
861 P3 variety groups. The genome was divided into 100,000-bp bins for jackknife resampling and  
862 calculation of the standard errors. The significance of the D-statistics was calculated using the Z-  
863 test, and D-statistics with Z-scores greater than  $|3.9|$  ( $p < 0.0001$ ) were considered significant.

864

865 **Population genomic analysis.** We downloaded FASTQ files from the 3K Rice Genome Project  
866 [7] for rice varieties that were determined to be *circum*-basmati varieties in that project. An  
867 additional 8 *circum*-basmati varieties were sequenced on the Illumina sequencing platform as  
868 part of this study. The raw reads were aligned to the scaffolded Basmati 334 genome using the  
869 program *bwa-mem*. PCR duplicates were determined computationally and removed using the  
870 program *picard* version 2.9.0 (<http://broadinstitute.github.io/picard/>). Genotype calls for each site  
871 were conducted using the *GATK HaplotypeCaller* engine using the option `-ERC GVCF`. The  
872 output files were in the genomic variant call format (gVCF), and the gVCFs from each variety  
873 were merged using the *GATK GenotypeGVCFs* engine.

874 SNP and INDEL variants from the population variant file were filtered independently  
875 using the *GATK* bestpractice hard filter pipeline [142]. SNP variants within 5 bps of an INDEL  
876 variant were filtered. *Vcftools* version 0.1.15 [143] was used to filter sites for which genotypes  
877 were not called for more than 20% of the varieties. Because domesticated rice is an inbreeding  
878 species we also implemented a heterozygosity filter by filtering out sites that had a heterozygote  
879 genotype in more than 5% of the samples using the program *vcffilterjdk.jar* from the *jvarkit* suite  
880 ([https://figshare.com/articles/JVarkit\\_java\\_based\\_utilities\\_for\\_Bioinformatics/1425030](https://figshare.com/articles/JVarkit_java_based_utilities_for_Bioinformatics/1425030)).

881 Missing genotypes were imputed and phased using *Beagle* version 4.1 [144].

882 To examine the within-*circum*-basmati variety group population structure we first  
883 randomly pruned the sites by sampling a polymorphic site every 200,000 bp using *plink* [145].  
884 *Plink* was also used to conduct a principal component analysis. Ancestry proportions of each  
885 sample were estimated using *fastSTRUCTURE* [91]. A neighbor-joining tree was built by  
886 calculating the pairwise genetic distances between samples using the Kronecker delta function-  
887 based equation [146]. From the genetic distance matrix a neighbor-joining tree was built using  
888 the program *FastME* [147].

889  
890 **Evolutionary relationships among the *circum*-basmati, *circum*-aus, and japonica**  
891 **populations.** To investigate the evolutionary origins of the *circum*-basmati population, we  
892 focused on the landrace varieties that had been sequenced with a genome-wide coverage of  
893 greater than 14×. The population data for the *circum*-aus and japonica populations were obtained  
894 from the 3K Rice Genome Project [7], from which we also analyzed only the landrace varieties  
895 that had been sequenced with a genome-wide coverage greater than 14×. For an outgroup, we  
896 obtained *O. barthii* sequencing data from previous studies [35, 148], and focused on the samples  
897 that were not likely to be feralized rices [148]. The Illumina reads were aligned to the scaffolded  
898 Basmati 334 genome and SNPs were called and filtered according to the procedure outlined in  
899 the “Population genomic analysis” section.

900 We examined the genome-wide local topological relationship using *twisst* [89]. Initially,  
901 a sliding window analysis was conducted to estimate the local phylogenetic trees in windows  
902 with a size of 100 or 500 polymorphic sites using *RAxML* with the GTRCAT substitution model.  
903 The script *raxml\_sliding\_windows.py* from the *genomics\_general* package by Simon Martin

904 ([https://github.com/simonhmartin/genomics\\_general/tree/master/phylo](https://github.com/simonhmartin/genomics_general/tree/master/phylo)) was used. The  
905 ‘complete’ option of *twisst* was used to calculate the exact weighting of each local window.  
906  
907 **□a□i demographic model.** The demography model underlying the evolution of *circum*-basmati  
908 rice was tested using the diffusion approximation method of  $\delta\delta i$  [90]. A visual representation of  
909 the 13 demographic models that were examined can be seen in Supplementary Figure S6. The  
910 population group and genotype calls used in the *twisst* analysis were also used to calculate the  
911 site allele frequencies. Polymorphic sites were polarized using the *O. barthii* reference genome.  
912 We used a previously published approach [148], which generates an *O. barthii*-ized basmati  
913 genome sequence. This was accomplished using the Basmati 334 reference genome to align the  
914 *O. barthii* genome. For every basmati genome sequence position was then changed into the  
915 aligned *O. barthii* sequence. Gaps, missing sequence, and repetitive DNA region were denoted  
916 as ‘N’.

917 We optimized the model parameter estimates using the Nelder-Mead method and  
918 randomly perturbed the parameter values for four rounds. Parameter values were perturbed for  
919 three-fold, two-fold, two-fold, and one-fold in each subsequent round, while the perturbation was  
920 conducted for 10, 20, 30, and 40 replicates in each subsequent round. In each round parameter  
921 values from the best likelihood model of the previous round were used as the starting parameter  
922 values for the next round. Parameter values from the round with the highest likelihood were  
923 chosen to parameterize each demographic model. Akaike Information Criteria (AIC) values were  
924 used to compare demography models. The demography model with the lowest AIC was chosen  
925 as the best-fitting model.

926

927 **Agronomic trait measurements.** Data on geolocation of collection as well as on seed  
928 dimensions and seed weight for each of the *circum*-basmati landrace varieties included in this  
929 study were obtained from passport data included in the online platform Genesys  
930 (<https://www.genesys-pgr.org/welcome>).

931

## 932 **DECLARATIONS**

933 **Ethics approval and consent to participate.** Not applicable.

934

935 **Consent for publication.** Not applicable.

936

937 **Availability of data and materials.** Raw nanopore sequencing FAST5 files generated from this  
938 study are available at the European Nucleotide Archive under bioproject ID PRJEB28274  
939 (ERX3327648-ERX3327652) for Basmati 334 and PRJEB32431 (ERX3334790-ERX3334793)  
940 for Dom Sufid. Associated FASTQ files are available under ERX3498039-ERX3498043 for  
941 Basmati 334 and ERX3498024-ERX3498027 for Dom Sufid. Illumina sequencing generated  
942 from this study can be found under bioproject ID PRJNA422249 and PRJNA557122. A genome  
943 browser for both genome assemblies can be found at [http://purugganan-](http://purugganan-genomebrowser.bio.nyu.edu/cgi-bin/hgTracks?db=Basmati334)  
944 [genomebrowser.bio.nyu.edu/cgi-bin/hgTracks?db=Basmati334](http://purugganan-genomebrowser.bio.nyu.edu/cgi-bin/hgTracks?db=Basmati334) for Basmati 334, and  
945 <http://purugganan-genomebrowser.bio.nyu.edu/cgi-bin/hgTracks?db=DomSufid> for Dom Sufid.  
946 All data including the assembly, annotation, genome alignment, and population VCFs generated  
947 from this study can be found at <https://doi.org/10.5281/zenodo.3355330>.

948

949 **Competing interests.** XD, PR, EDH, and SJ are employees of Oxford Nanopore Technologies  
950 and are shareholders and/or share option holders.

951  
952 **Funding.** This work was supported by grants from the Gordon and Betty Moore Foundation  
953 through Grant GBMF2550.06 to S.C.G., and from the National Science Foundation Plant  
954 Genome Research Program (IOS-1546218), the Zegar Family Foundation (A16-0051) and the  
955 NYU Abu Dhabi Research Institute (G1205) to M.D.P. The funding body had no role in the  
956 design of the study and collection, analysis, and interpretation of data and in writing the  
957 manuscript.

958  
959 **Authors' contributions.** JYC, SCG, SZ, and MDP conceived the project and its components.  
960 JYC, SCG, and SZ prepared the sample material for sequencing. XD, PR, EDH, and SJ  
961 conducted the genome sequencing and assembling. JYC, ZNL, and SCG performed the data  
962 analysis. JYC and ZNL prepared the figures and tables. JYC and MDP wrote the manuscript  
963 with help from ZNL and SCG.

964  
965 **Acknowledgements.** We thank Katherine Dorph for assistance with growing and maintaining  
966 the plants, and Adrian Platts for computational support. We thank Rod Wing, David Kudrna, and  
967 Jayson Talag from Arizona Genomics Institute with the high-molecular weight DNA extraction.  
968 We thank the New York University Genomics Core Facility for sequencing support and the New  
969 York University High Performance Computing for supplying the computational resources.

970

971 **REFERENCES**

- 972 1. Gnanamanickam SS. Rice and Its Importance to Human Life. In: Biological Control of Rice  
973 Diseases. Dordrecht: Springer Netherlands; 2009. p. 1–11. doi:10.1007/978-90-481-2465-7\_1.
- 974 2. Matsuo T, Futsuhara Y, Kikuchi F, Yamaguchi H. Science of the Rice Plant. Tokyo: Food and  
975 Agriculture Policy Research Center; 1997.
- 976 3. Gross BL, Zhao Z. Archaeological and genetic insights into the origins of domesticated rice.  
977 Proc Natl Acad Sci USA. 2014;111:6190–7.
- 978 4. Nadir S, Khan S, Zhu Q, Henry D, Wei L, Lee DS, et al. An overview on reproductive  
979 isolation in *Oryza sativa* complex. AoB Plants. 2018;10:ply060.
- 980 5. Fuller DQ, Sato Y-I, Castillo C, Qin L, Weisskopf AR, Kingwell-Banham EJ, et al.  
981 Consilience of genetics and archaeobotany in the entangled history of rice. Archaeol Anthropol  
982 Sci. 2010;2:115–31.
- 983 6. Garris AJ, Tai TH, Coburn J, Kresovich S, McCouch S. Genetic structure and diversity in  
984 *Oryza sativa* L. Genetics. 2005;169:1631–8.
- 985 7. Wang W, Mauleon R, Hu Z, Chebotarov D, Tai S, Wu Z, et al. Genomic variation in 3,010  
986 diverse accessions of Asian cultivated rice. Nature. 2018;557:43–9.
- 987 8. Glaszmann JC. Isozymes and classification of Asian rice varieties. Theoret Appl Genetics.  
988 1987;74:21–30.
- 989 9. He Z, Zhai W, Wen H, Tang T, Wang Y, Lu X, et al. Two evolutionary histories in the  
990 genome of rice: the roles of domestication genes. PLoS Genet. 2011;7:e1002100.
- 991 10. Fuller DQ. Pathways to Asian Civilizations: Tracing the Origins and Spread of Rice and Rice  
992 Cultures. Rice. 2012;4:78–92.
- 993 11. Meyer RS, Purugganan MD. Evolution of crop species: genetics of domestication and  
994 diversification. Nat Rev Genet. 2013;14:840–52.
- 995 12. Huang X, Han B. Rice domestication occurred through single origin and multiple  
996 introgressions. Nat Plants. 2015;2:15207.
- 997 13. Castillo CC, Tanaka K, Sato Y-I, Ishikawa R, Bellina B, Higham C, et al. Archaeogenetic  
998 study of prehistoric rice remains from Thailand and India: evidence of early japonica in South  
999 and Southeast Asia. Archaeological and Anthropological Sciences. 2016;8:523–43.
- 1000 14. Choi JY, Platts AE, Fuller DQ, Hsing Y-I, Wing RA, Purugganan MD. The rice paradox:  
1001 Multiple origins but single domestication in Asian rice. Molecular Biology and Evolution.  
1002 2017;34:969–79.
- 1003 15. Choi JY, Purugganan MD. Multiple Origin but Single Domestication Led to *Oryza sativa*.  
1004 G3: Genes, Genomes, Genetics. 2018;8:797–803.

- 1005 16. Huang X, Kurata N, Wei X, Wang Z-X, Wang A, Zhao Q, et al. A map of rice genome  
1006 variation reveals the origin of cultivated rice. *Nature*. 2012;490:497–501.
- 1007 17. Civián P, Craig H, Cox CJ, Brown TA. Three geographically separate domestications of  
1008 Asian rice. *Nat Plants*. 2015;1:15164.
- 1009 18. Wang ZY, Zheng FQ, Shen GZ, Gao JP, Snustad DP, Li MG, et al. The amylose content in  
1010 rice endosperm is related to the post-transcriptional regulation of the waxy gene. *Plant J*.  
1011 1995;7:613–22.
- 1012 19. Sweeney MT, Thomson MJ, Pfeil BE, McCouch S. Caught red-handed: Rc encodes a basic  
1013 helix-loop-helix protein conditioning red pericarp in rice. *Plant Cell*. 2006;18:283–94.
- 1014 20. Konishi S, Izawa T, Lin SY, Ebana K, Fukuta Y, Sasaki T, et al. An SNP Caused Loss of  
1015 Seed Shattering During Rice Domestication. *Science*. 2006;312:1392–6.
- 1016 21. Kovach MJ, Calingacion MN, Fitzgerald MA, McCouch SR. The origin and evolution of  
1017 fragrance in rice (*Oryza sativa* L.). *Proceedings of the National Academy of Sciences of the*  
1018 *United States of America*. 2009;106:14444–9.
- 1019 22. Xu K, Xu X, Fukao T, Canlas P, Maghirang-Rodriguez R, Heuer S, et al. Sub1A is an  
1020 ethylene-response-factor-like gene that confers submergence tolerance to rice. *Nature*.  
1021 2006;442:705–8.
- 1022 23. Bin Rahman ANMR, Zhang J. Preferential Geographic Distribution Pattern of Abiotic Stress  
1023 Tolerant Rice. *Rice*. 2018;11:10.
- 1024 24. Singh R, Singh U, Khush G, editors. *Aromatic rices*. Oxford & IBH Publishing Co Pvt Ltd;  
1025 2000.
- 1026 25. Bradbury LMT, Gillies SA, Brushett DJ, Waters DLE, Henry RJ. Inactivation of an  
1027 aminoaldehyde dehydrogenase is responsible for fragrance in rice. *Plant Mol Biol*. 2008;68:439–  
1028 49.
- 1029 26. Chen S, Yang Y, Shi W, Ji Q, He F, Zhang Z, et al. Badh2, Encoding Betaine Aldehyde  
1030 Dehydrogenase, Inhibits the Biosynthesis of 2-Acetyl-1-Pyrroline, a Major Component in Rice  
1031 Fragrance. *Plant Cell*. 2008;20:1850–61.
- 1032 27. Zhao K, Tung C-W, Eizenga GC, Wright MH, Ali ML, Price AH, et al. Genome-wide  
1033 association mapping reveals a rich genetic architecture of complex traits in *Oryza sativa*. *Nature*  
1034 *Communications*. 2011;2:467.
- 1035 28. Heather JM, Chain B. The sequence of sequencers: The history of sequencing DNA.  
1036 *Genomics*. 2016;107:1–8.
- 1037 29. Michael TP, VanBuren R. Progress, challenges and the future of crop genomes. *Current*  
1038 *Opinion in Plant Biology*. 2015;24:71–81.

- 1039 30. Jiao W-B, Schneeberger K. The impact of third generation genomic technologies on plant  
1040 genome assembly. *Current Opinion in Plant Biology*. 2017;36:64–70.
- 1041 31. Li C, Lin F, An D, Wang W, Huang R. Genome Sequencing and Assembly by Long Reads in  
1042 Plants. *Genes*. 2017;9. doi:10.3390/genes9010006.
- 1043 32. Yu J, Hu S, Wang J, Wong GK-S, Li S, Liu B, et al. A draft sequence of the rice genome  
1044 (*Oryza sativa* L. ssp. *indica*). *Science*. 2002;296:79–92.
- 1045 33. International Rice Genome Sequencing Project. The map-based sequence of the rice genome.  
1046 *Nature*. 2005;436:793–800.
- 1047 34. Chen J, Huang Q, Gao D, Wang J, Lang Y, Liu T, et al. Whole-genome sequencing of *Oryza*  
1048 *brachyantha* reveals mechanisms underlying *Oryza* genome evolution. *Nat Commun*.  
1049 2013;4:1595.
- 1050 35. Wang M, Yu Y, Haberer G, Marri PR, Fan C, Goicoechea JL, et al. The genome sequence of  
1051 African rice (*Oryza glaberrima*) and evidence for independent domestication. *Nat Genet*.  
1052 2014;46:982–8.
- 1053 36. Zhang Y, Zhang S, Liu H, Fu B, Li L, Xie M, et al. Genome and Comparative  
1054 Transcriptomics of African Wild Rice *Oryza longistaminata* Provide Insights into Molecular  
1055 Mechanism of Rhizomatousness and Self-Incompatibility. *Molecular Plant*. 2015;8:1683–6.
- 1056 37. Stein JC, Yu Y, Copetti D, Zwickl DJ, Zhang L, Zhang C, et al. Genomes of 13 domesticated  
1057 and wild rice relatives highlight genetic conservation, turnover and innovation across the genus  
1058 *Oryza*. *Nature Genetics*. 2018;50:285.
- 1059 38. Schatz MC, Maron LG, Stein JC, Wences A, Gurtowski J, Biggers E, et al. Whole genome  
1060 de novo assemblies of three divergent strains of rice, *Oryza sativa*, document novel gene space  
1061 of *aus* and *indica*. *Genome Biology*. 2014;15:506.
- 1062 39. Sakai H, Kanamori H, Arai-Kichise Y, Shibata-Hatta M, Ebana K, Oono Y, et al.  
1063 Construction of Pseudomolecule Sequences of the *aus* Rice Cultivar Kasalath for Comparative  
1064 Genomics of Asian Cultivated Rice. *DNA Research*. 2014;21:397–405.
- 1065 40. Zhang J, Chen L-L, Xing F, Kudrna DA, Yao W, Copetti D, et al. Extensive sequence  
1066 divergence between the reference genomes of two elite *indica* rice varieties Zhenshan 97 and  
1067 Minghui 63. *Proceedings of the National Academy of Sciences of the United States of America*.  
1068 2016;113:E5163-71.
- 1069 41. Du H, Yu Y, Ma Y, Gao Q, Cao Y, Chen Z, et al. Sequencing and de novo assembly of a  
1070 near complete *indica* rice genome. *Nature Communications*. 2017;8:15324.
- 1071 42. Zhao Q, Feng Q, Lu H, Li Y, Wang A, Tian Q, et al. Pan-genome analysis highlights the  
1072 extent of genomic variation in cultivated and wild rice. *Nature Genetics*. 2018;50:278–84.



- 1073 43. Jain S, Jain RK, McCouch SR. Genetic analysis of Indian aromatic and quality rice (*Oryza*  
1074 *sativa* L.) germplasm using panels of fluorescently-labeled microsatellite markers. *Theoretical*  
1075 *and Applied Genetics*. 2004;109:965–77.
- 1076 44. Vikram P, Swamy BPM, Dixit S, Ahmed H, Cruz MTS, Singh AK, et al. Bulk segregant  
1077 analysis: “An effective approach for mapping consistent-effect drought grain yield QTLs in  
1078 rice.” *Field Crops Research*. 2012;134:185–92.
- 1079 45. McNally KL, Childs KL, Bohnert R, Davidson RM, Zhao K, Ulat VJ, et al. Genomewide  
1080 SNP variation reveals relationships among landraces and modern varieties of rice. *Proceedings*  
1081 *of the National Academy of Sciences*. 2009;106:12273–8.
- 1082 46. McNally KL, Bruskiewich R, Mackill D, Buell CR, Leach JE, Leung H. Sequencing  
1083 Multiple and Diverse Rice Varieties. *Connecting Whole-Genome Variation with Phenotypes*.  
1084 *Plant Physiology*. 2006;141:26–31.
- 1085 47. Jain M, Olsen HE, Paten B, Akeson M. The Oxford Nanopore MinION: delivery of nanopore  
1086 sequencing to the genomics community. *Genome Biology*. 2016;17:239.
- 1087 48. Chen S, Huang Z, Zeng L, Yang J, Liu Q, Zhu X. High-resolution mapping and gene  
1088 prediction of *Xanthomonas Oryzae* pv. *Oryzae* resistance gene Xa7. *Molecular Breeding*.  
1089 2008;22:433–41.
- 1090 49. Ullah I, Jamil S, Iqbal MZ, Shaheen HL, Hasni SM, Jabeen S, et al. Detection of bacterial  
1091 blight resistance genes in basmati rice landraces. *Genetics and Molecular Research*.  
1092 2012;11:1960–6.
- 1093 50. Sandhu N, Kumar A, Sandhu N, Kumar A. Bridging the Rice Yield Gaps under Drought:  
1094 QTLs, Genes, and their Use in Breeding Programs. *Agronomy*. 2017;7:27.
- 1095 51. Henry A, Gowda VRP, Torres RO, McNally KL, Serraj R. Variation in root system  
1096 architecture and drought response in rice (*Oryza sativa*): Phenotyping of the *Oryza*SNP panel in  
1097 rainfed lowland fields. *Field Crops Research*. 2011;120:205–14.
- 1098 52. Cantarel BL, Korf I, Robb SMC, Parra G, Ross E, Moore B, et al. MAKER: An easy-to-use  
1099 annotation pipeline designed for emerging model organism genomes. *Genome Res*.  
1100 2008;18:188–96.
- 1101 53. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database management  
1102 tool for second-generation genome projects. *BMC Bioinformatics*. 2011;12:491.
- 1103 54. Choi JY, Purugganan MD. Evolutionary epigenomics of retrotransposon-mediated  
1104 methylation spreading in rice. *Molecular Biology and Evolution*. 2018;35:365–82.
- 1105 55. Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, von Haeseler A, et al.  
1106 Accurate detection of complex structural variations using single-molecule sequencing. *Nature*  
1107 *Methods*. 2018;15:461–8.

- 1108 56. Alonge M, Soyk S, Ramakrishnan S, Wang X, Goodwin S, Sedlazeck FJ, et al. Fast and  
1109 accurate reference-guided scaffolding of draft genomes. *bioRxiv*. 2019;:519637.
- 1110 57. Kawahara Y, de la Bastide M, Hamilton JP, Kanamori H, McCombie WR, Ouyang S, et al.  
1111 Improvement of the *Oryza sativa* Nipponbare reference genome using next generation sequence  
1112 and optical map data. *Rice*. 2013;6:4.
- 1113 58. Sakai H, Lee SS, Tanaka T, Numa H, Kim J, Kawahara Y, et al. Rice Annotation Project  
1114 Database (RAP-DB): An Integrative and Interactive Database for Rice Genomics. *Plant and Cell  
1115 Physiology*. 2013;54:e6–e6.
- 1116 59. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome comparisons  
1117 dramatically improves orthogroup inference accuracy. *Genome Biology*. 2015;16:157.
- 1118 60. Nattestad M, Schatz MC. Assemblytics: a web analytics tool for the detection of variants  
1119 from an assembly. *Bioinformatics*. 2016;32:3021–3.
- 1120 61. Fuentes RR, Chebotarov D, Duitama J, Smith S, Hoz JFD la, Mohiyuddin M, et al. Structural  
1121 variants in 3000 rice genomes. *Genome Res*. 2019;29:870–80.
- 1122 62. Shomura A, Izawa T, Ebana K, Ebitani T, Kanegae H, Konishi S, et al. Deletion in a gene  
1123 associated with grain size increased yields during rice domestication. *Nature Genetics*.  
1124 2008;40:1023–8.
- 1125 63. Zhou Y, Zhu J, Li Z, Yi C, Liu J, Zhang H, et al. Deletion in a Quantitative Trait Gene qPE9-  
1126 1 Associated With Panicle Erectness Improves Plant Architecture During Rice Domestication.  
1127 *Genetics*. 2009;183:315–24.
- 1128 64. Lye ZN, Purugganan MD. Copy Number Variation in Domestication. *Trends in Plant  
1129 Science*. 2019;24:352–65.
- 1130 65. Hu M, Lv S, Wu W, Fu Y, Liu F, Wang B, et al. The domestication of plant architecture in  
1131 African rice. *The Plant Journal*. 2018;94:661–9.
- 1132 66. Wu Y, Zhao S, Li X, Zhang B, Jiang L, Tang Y, et al. Deletions linked to PROG1 gene  
1133 participate in plant architecture domestication in Asian and African rice. *Nature  
1134 Communications*. 2018;9:4157.
- 1135 67. Li B, Zhang Y, Li J, Yao G, Pan H, Hu G, et al. Fine Mapping of Two Additive Effect Genes  
1136 for Awn Development in Rice (*Oryza sativa* L.). *PLOS ONE*. 2016;11:e0160792.
- 1137 68. Hua L, Wang DR, Tan L, Fu Y, Liu F, Xiao L, et al. LABA1, a Domestication Gene  
1138 Associated with Long, Barbed Awns in Wild Rice. *The Plant Cell*. 2015;27:1875–88.
- 1139 69. Kumar A, Bennetzen JL. Plant Retrotransposons. *Annual Review of Genetics*. 1999;33:479–  
1140 532.

- 1141 70. Zuccolo A, Sebastian A, Talag J, Yu Y, Kim H, Collura K, et al. Transposable element  
1142 distribution, abundance and role in genome size variation in the genus *Oryza*. *BMC Evol Biol*.  
1143 2007;7:152.
- 1144 71. Lerat E. Identifying repeats and transposable elements in sequenced genomes: how to find  
1145 your way through the dense forest of programs. *Heredity*. 2010;104:520–33.
- 1146 72. Hoen DR, Hickey G, Bourque G, Casacuberta J, Cordaux R, Feschotte C, et al. A call for  
1147 benchmarking transposable element annotation methods. *Mobile DNA*. 2015;6:13.
- 1148 73. Bennetzen JL. The contributions of retroelements to plant genome organization, function and  
1149 evolution. *Trends Microbiol*. 1996;4:347–53.
- 1150 74. Voytas DF, Ausubel FM. A copia-like transposable element family in *Arabidopsis thaliana*.  
1151 *Nature*. 1988;336:242–4.
- 1152 75. SanMiguel P, Gaut BS, Tikhonov A, Nakajima Y, Bennetzen JL. The paleontology of  
1153 intergene retrotransposons of maize. *Nature Genetics*. 1998;20:43–5.
- 1154 76. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, et al. A unified  
1155 classification system for eukaryotic transposable elements. *Nature Reviews Genetics*.  
1156 2007;8:973–82.
- 1157 77. Vitte C, Panaud O, Quesneville H. LTR retrotransposons in rice (*Oryza sativa*, L.): recent  
1158 burst amplifications followed by rapid DNA loss. *BMC Genomics*. 2007;8:218.
- 1159 78. Panaud O, Vitte C, Hivert J, Muzlak S, Talag J, Brar D, et al. Characterization of  
1160 transposable elements in the genome of rice (*Oryza sativa* L.) using Representational Difference  
1161 Analysis (RDA). *Molecular genetics and genomics*: MGG. 2002;268:113–21.
- 1162 79. McCarthy EM, Liu J, Lizhi G, McDonald JF. Long terminal repeat retrotransposons of *Oryza*  
1163 *sativa*. *Genome biology*. 2002;3:RESEARCH0053.
- 1164 80. Carpentier M-C, Manfroi E, Wei F-J, Wu H-P, Lasserre E, Llauro C, et al.  
1165 Retrotranspositional landscape of Asian rice revealed by 3000 genomes. *Nat Commun*. 2019;10.  
1166 doi:10.1038/s41467-018-07974-5.
- 1167 81. Zhang Q-J, Gao L-Z. Rapid and Recent Evolution of LTR Retrotransposons Drives Rice  
1168 Genome Evolution During the Speciation of AA- Genome *Oryza* Species. *G3: Genes, Genomes,*  
1169 *Genetics*. 2017;;g3.116.037572.
- 1170 82. Kumekawa N, Ohtsubo H, Horiuchi T, Ohtsubo E. Identification and characterization of  
1171 novel retrotransposons of the gypsy type in rice. *Mol Gen Genet*. 1999;260:593–602.
- 1172 83. Chaparro C, Guyot R, Zuccolo A, Piégu B, Panaud O. RetrOryza: a database of the rice  
1173 LTR-retrotransposons. *Nucleic acids research*. 2007;35 Database issue:D66-70.

- 1174 84. Wang H, Vieira FG, Crawford JE, Chu C, Nielsen R. Asian wild rice is a hybrid swarm with  
1175 extensive gene flow and feralization from domesticated rice. *Genome research*. 2017;27:1029–  
1176 38.
- 1177 85. Li L-F, Li Y-L, Jia Y, Caicedo AL, Olsen KM. Signatures of adaptation in the weedy rice  
1178 genome. *Nature genetics*. 2017;49:811–4.
- 1179 86. Shimodaira H. An approximately unbiased test of phylogenetic tree selection. *Syst Biol*.  
1180 2002;51:492–508.
- 1181 87. Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, et al. A draft sequence of  
1182 the Neandertal genome. *Science*. 2010;328:710–22.
- 1183 88. Durand EY, Patterson N, Reich D, Slatkin M. Testing for Ancient Admixture between  
1184 Closely Related Populations. *Molecular Biology and Evolution*. 2011;28:2239–52.
- 1185 89. Martin SH, Van Belleghem SM. Exploring Evolutionary Relationships Across the Genome  
1186 Using Topology Weighting. *Genetics*. 2017;206:429–38.
- 1187 90. Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. Inferring the joint  
1188 demographic history of multiple populations from multidimensional SNP frequency data. *PLoS*  
1189 *Genet*. 2009;5:e1000695.
- 1190 91. Raj A, Stephens M, Pritchard JK. fastSTRUCTURE: variational inference of population  
1191 structure in large SNP data sets. *Genetics*. 2014;197:573–89.
- 1192 92. Michael TP, Jupe F, Bemm F, Motley ST, Sandoval JP, Lanz C, et al. High contiguity  
1193 *Arabidopsis thaliana* genome assembly with a single nanopore flow cell. *Nature*  
1194 *Communications*. 2018;9:541.
- 1195 93. Schmidt MH, Vogel A, Denton AK, Istace B, Wormit A, van de Geest H, et al. De novo  
1196 Assembly of a New *Solanum pennellii* Accession Using Nanopore Sequencing. *The Plant cell*.  
1197 2017;:tpc.00521.2017.
- 1198 94. Belser C, Istace B, Denis E, Dubarry M, Baurens F-C, Falentin C, et al. Chromosome-scale  
1199 assemblies of plant genomes using nanopore long reads and optical maps. *Nature Plants*.  
1200 2018;4:879.
- 1201 95. Howe K, Wood JM. Using optical mapping data for the improvement of vertebrate genome  
1202 assemblies. *GigaScience*. 2015;4:10.
- 1203 96. Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of next-generation  
1204 sequencing technologies. *Nature Reviews Genetics*. 2016;17:333–51.
- 1205 97. Udall JA, Dawe RK. Is It Ordered Correctly? Validating Genome Assemblies by Optical  
1206 Mapping. *The Plant cell*. 2018;30:7–14.

- 1207 98. Sedlazeck FJ, Lee H, Darby CA, Schatz MC. Piercing the dark matter: bioinformatics of  
1208 long-range sequencing and mapping. *Nature Reviews Genetics*. 2018;19:329–46.
- 1209 99. Wendel JF, Jackson SA, Meyers BC, Wing RA. Evolution of plant genome architecture.  
1210 *Genome Biology*. 2016;17:37.
- 1211 100. vonHoldt BM, Takuno S, Gaut BS. Recent Retrotransposon Insertions Are Methylated and  
1212 Phylogenetically Clustered in Japonica Rice (*Oryza sativa* spp. japonica). *Molecular Biology and*  
1213 *Evolution*. 2012;29:3193–203.
- 1214 101. Baucom RS, Estill JC, Leebens-Mack J, Bennetzen JL. Natural selection on gene function  
1215 drives the evolution of LTR retrotransposon families in the rice genome. *Genome Research*.  
1216 2009;19:243–54.
- 1217 102. Naito K, Zhang F, Tsukiyama T, Saito H, Hancock CN, Richardson AO, et al. Unexpected  
1218 consequences of a sudden and massive transposon amplification on rice gene expression. *Nature*.  
1219 2009;461:1130–4.
- 1220 103. Civián P, Ali S, Batista-Navarro R, Drosou K, Ihejieta C, Chakraborty D, et al. Origin of the  
1221 Aromatic Group of Cultivated Rice (*Oryza sativa* L.) Traced to the Indian Subcontinent. *Genome*  
1222 *Biol Evol*. 2019;11:832–43.
- 1223 104. Molina J, Sikora M, Garud N, Flowers JM, Rubinstein S, Reynolds A, et al. Molecular  
1224 evidence for a single evolutionary origin of domesticated rice. *Proc Natl Acad Sci USA*.  
1225 2011;108:8351–6.
- 1226 105. Fuller DQ. Finding Plant Domestication in the Indian Subcontinent. *Current Anthropology*.  
1227 2011;52:S347–62.
- 1228 106. Filiault DL, Ballerini ES, Mandáková T, Aköz G, Derieg NJ, Schmutz J, et al. The  
1229 *Aquilegia* genome provides insight into adaptive radiation and reveals an extraordinarily  
1230 polymorphic chromosome with a unique history. *eLife*. 2018;7:e36426.
- 1231 107. Liakat Ali M, McClung AM, Jia MH, Kimball JA, McCouch, Georgia CE. A Rice Diversity  
1232 Panel Evaluated for Genetic and Agro-Morphological Diversity between Subpopulations and its  
1233 Geographic Distribution. *Crop Science*. 2011;51:2021–35.
- 1234 108. Travis AJ, Norton GJ, Datta S, Sarma R, Dasgupta T, Savio FL, et al. Assessing the genetic  
1235 diversity of rice originating from Bangladesh, Assam and West Bengal. *Rice (N Y)*. 2015;8:35.
- 1236 109. Zhang H-B, Zhao X, Ding X, Paterson AH, Wing RA. Preparation of megabase-size DNA  
1237 from plant nuclei. *The Plant Journal*. 1995;7:175–84.
- 1238 110. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and  
1239 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome*  
1240 *research*. 2017;27:722–36.

- 1241 111. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using  
1242 repeat graphs. *Nature Biotechnology*. 2019;37:540.
- 1243 112. Vaser R, Sović I, Nagarajan N, Šikić M. Fast and accurate de novo genome assembly from  
1244 long uncorrected reads. *Genome research*. 2017;27:737–46.
- 1245 113. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
1246 arXiv. 2013;:1303.3997v2.
- 1247 114. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An  
1248 Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly  
1249 Improvement. *PLoS ONE*. 2014;9:e112963.
- 1250 115. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*.  
1251 2018;34:3094–100.
- 1252 116. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:  
1253 assessing genome assembly and annotation completeness with single-copy orthologs.  
1254 *Bioinformatics*. 2015;31:3210–2.
- 1255 117. Cabanettes F, Klopp C. D-GENIES: dot plot large genomes in an interactive, efficient and  
1256 simple way. *PeerJ*. 2018;6. doi:10.7717/peerj.4958.
- 1257 118. Marçais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A. MUMmer4: A  
1258 fast and versatile genome alignment system. *PLOS Computational Biology*. 2018;14:e1005944.
- 1259 119. Korf I. Gene finding in novel genomes. *BMC Bioinformatics*. 2004;5:59.
- 1260 120. Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically mapped  
1261 cDNA alignments to improve de novo gene finding. *Bioinformatics*. 2008;24:637–44.
- 1262 121. Conway JR, Lex A, Gehlenborg N, Hancock J. UpSetR: an R package for the visualization  
1263 of intersecting sets and their properties. *Bioinformatics*. 2017;33:2938–40.
- 1264 122. Tian T, Liu Y, Yan H, You Q, Yi X, Du Z, et al. agriGO v2.0: a GO analysis toolkit for the  
1265 agricultural community, 2017 update. *Nucleic Acids Res*. 2017;45 Web Server issue:W122–9.
- 1266 123. Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von Mering C, et al. Fast  
1267 Genome-Wide Functional Annotation through Orthology Assignment by eggNOG-Mapper. *Mol*  
1268 *Biol Evol*. 2017;34:2115–22.
- 1269 124. Falcon S, Gentleman R. Using GOSTats to test gene lists for GO term association.  
1270 *Bioinformatics*. 2007;23:257–8.
- 1271 125. Ellinghaus D, Kurtz S, Willhoeft U. LTRharvest, an efficient and flexible software for de  
1272 novo detection of LTR retrotransposons. *BMC bioinformatics*. 2008;9:18.

- 1273 126. Copetti D, Zhang J, El Baidouri M, Gao D, Wang J, Barghini E, et al. RiTE database: a  
1274 resource database for genus-wide rice genomics and evolutionary biology. *BMC Genomics*.  
1275 2015;16:538.
- 1276 127. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:  
1277 architecture and applications. *BMC Bioinformatics*. 2009;10:421.
- 1278 128. Katoh K, Standley DM. MAFFT Multiple Sequence Alignment Software Version 7:  
1279 Improvements in Performance and Usability. *Molecular Biology and Evolution*. 2013;30:772–  
1280 80.
- 1281 129. Yang Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Molecular Biology and*  
1282 *Evolution*. 2007;24:1586–91.
- 1283 130. Kimura M. A simple method for estimating evolutionary rates of base substitutions through  
1284 comparative studies of nucleotide sequences. *Journal of molecular evolution*. 1980;16:111–20.
- 1285 131. Ma J, Bennetzen JL. Rapid recent growth and divergence of rice nuclear genomes. *Proc*  
1286 *Natl Acad Sci USA*. 2004;101:12404–10.
- 1287 132. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic  
1288 features. *Bioinformatics*. 2010;26:841–2.
- 1289 133. Haudry A, Platts AE, Vello E, Hoen DR, Leclercq M, Williamson RJ, et al. An atlas of over  
1290 90,000 conserved noncoding sequences provides insight into crucifer regulatory regions. *Nat*  
1291 *Genet*. 2013;45:891–8.
- 1292 134. Harris RS. Improved pairwise alignment of genomic dna. PhD Thesis, The Pennsylvania  
1293 State University. 2007.
- 1294 135. Kent WJ, Baertsch R, Hinrichs A, Miller W, Haussler D. Evolution's cauldron: duplication,  
1295 deletion, and rearrangement in the mouse and human genomes. *Proc Natl Acad Sci USA*.  
1296 2003;100:11484–9.
- 1297 136. Blanchette M, Kent WJ, Riemer C, Elnitski L, Smit AFA, Roskin KM, et al. Aligning  
1298 multiple genomic sequences with the threaded blockset aligner. *Genome Res*. 2004;14:708–15.
- 1299 137. Hubisz MJ, Pollard KS, Siepel A. PHAST and RPHAST: phylogenetic analysis with  
1300 space/time models. *Brief Bioinformatics*. 2011;12:41–51.
- 1301 138. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
1302 phylogenies. *Bioinformatics*. 2014;30:1312–3.
- 1303 139. Martin SH, Jiggins CD. Interpreting the genomic landscape of introgression. *Current*  
1304 *Opinion in Genetics & Development*. 2017;47:69–74.
- 1305 140. Goldman N, Anderson JP, Rodrigo AG, Olmstead R. Likelihood-Based Tests of Topologies  
1306 in Phylogenetics. *Systematic Biology*. 2000;49:652–70.

- 1307 141. Shimodaira H, Hasegawa M. CONSEL: for assessing the confidence of phylogenetic tree  
1308 selection. *Bioinformatics*. 2001;17:1246–7.
- 1309 142. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, del Angel G, Levy-Moonshine A, et  
1310 al. From FastQ Data to High-Confidence Variant Calls: The Genome Analysis Toolkit Best  
1311 Practices Pipeline. In: *Current Protocols in Bioinformatics*. Hoboken, NJ, USA: John Wiley &  
1312 Sons, Inc.; 2013. p. 11.10.1-11.10.33. doi:10.1002/0471250953.bi1110s43.
- 1313 143. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call  
1314 format and VCFtools. *Bioinformatics*. 2011;27:2156–8.
- 1315 144. Browning BL, Browning SR. Genotype Imputation with Millions of Reference Samples.  
1316 *The American Journal of Human Genetics*. 2016;98:116–26.
- 1317 145. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: A  
1318 Tool Set for Whole-Genome Association and Population-Based Linkage Analyses. *The*  
1319 *American Journal of Human Genetics*. 2007;81:559–75.
- 1320 146. Freedman AH, Gronau I, Schweizer RM, Ortega-Del Vecchyo D, Han E, Silva PM, et al.  
1321 Genome sequencing highlights the dynamic early history of dogs. *PLoS Genet*.  
1322 2014;10:e1004016.
- 1323 147. Lefort V, Desper R, Gascuel O. FastME 2.0: A Comprehensive, Accurate, and Fast  
1324 Distance-Based Phylogeny Inference Program. *Molecular Biology and Evolution*.  
1325 2015;32:2798–800.
- 1326 148. Choi JY, Zaidem M, Gutaker R, Dorph K, Singh RK, Purugganan MD. The complex  
1327 geography of domestication of the African rice *Oryza glaberrima*. *PLOS Genetics*.  
1328 2019;15:e1007414.
- 1329



## Figure Legend

**Figure 1. Dot plot comparing the assembly contigs of Basmati 334 and Dom Sufid** to (A) all chromosomes of the Nipponbare genome assembly and (B) only chromosome 6 of Nipponbare. Only alignment blocks with greater than 80% overlap in sequence identity are shown.

**Figure 2. *Circum-basmati* gene sequence evolution.** (A) The deletion frequency of genes annotated from the Basmati 334 and Dom Sufid genomes. Frequency was estimated from sequencing data on a population of 78 *circum-basmati* varieties. (B) Groups of orthologous and paralogous genes (*i.e.*, orthogroups) identified in the reference genomes of N22, Nipponbare (NPB), and R498, as well as the *circum-basmati* genome assemblies Basmati 334 (B334) and Dom Sufid (DS) of this study. (C) Visualization of the genomic region orthologous to the Nipponbare gene Os03g0418600 (*Awn3-1*) in the N22, Basmati 334, and Dom Sufid genomes. Regions orthologous to *Awn3-1* are indicated with a dotted box.

**Figure 3. Presence/absence variation across the *circum-basmati* rice genome assemblies.** (A) Distribution of presence/absence variant sizes compared to the japonica Nipponbare reference genome. (B) Number of presence/absence variants that are shared between or unique for the *circum-basmati* genomes. (C) Chromosome-wide distribution of presence/absence variation for each *circum-basmati* rice genome, relative to the Nipponbare genome coordinates.

**Figure 4. Repetitive DNA landscape of the Basmati 334 and Dom Sufid genomes.** (A) Proportion of repetitive DNA content in the *circum-basmati* genomes represented by each repeat family. (B) Distribution of insert times for the *gypsy* and *copia* LTR retrotransposons. (C)

Phylogeny of *gypsy* and *copia* LTR retrotransposons based on the *rve* gene. LTR retrotransposons were annotated from the reference genomes of domesticated and wild rices.

**Figure 5. Comparative genomic analysis of *circum-basmati* rice evolution.** (A) Maximum-likelihood tree based on four-fold degenerate sites. All nodes had over 95% bootstrap support. (B) Percentage of genes supporting the topology involving japonica (J; Nipponbare, NPB), *circum-basmati* (cB, *circum-basmati*; Basmati 334, B334; Dom Sufid, DS), and *O. rufipogon* (R) after an Approximately Unbiased (AU) test. (C) Results of ABBA-BABA tests. Shown are median Patterson's D-statistics with 95% confidence intervals determined from a bootstrapping procedure. For each tested topology the outgroup was always *O. barthii*. (D) Percentage of genes supporting the topology involving *circum-aus* (cA; N22), *circum-basmati*, and indica (I; R498) after an Approximately Unbiased (AU) test. (E) Per-chromosome distribution of D-statistics for the trio involving R498, N22, and each *circum-basmati* genome. Genome-wide D-statistics with 95% bootstrap confidence intervals are indicated by the dark and dotted lines. (F) Model of admixture events that occurred within domesticated Asian rice. The direction of admixture has been left ambiguous as the ABBA-BABA test cannot detect the direction of gene flow. The *Oryza sativa* variety groups are labeled as *circum-aus* (cA), *circum-basmati* (cB), indica (I), and japonica (J), and the wild relative is *O. rufipogon* (R).

**Figure 6. Population relationships among the *circum-aus* (cA), *circum-basmati* (cB), and japonica rices (J).** (A) Sum of genome-wide topology weights for a three-population topology involving trios of the *circum-aus*, *circum-basmati*, and japonica rices. Topology weights were estimated across windows with 100 SNPs. (B) Chromosomal distributions of topology weights

involving trios of the *circum*-aus, *circum*-basmati, and japonica rices (left), and the sum of the topology weights (right). (C) Best-fitting  $\delta a\delta i$  model for the *circum*-aus, *circum*-basmati, and japonica rices. See Supplemental Table 17 for parameter estimates.

**Figure 7. Population structure within the *circum*-basmati rices.** (A) PCA plot for the 78-variety *circum*-basmati rice population genomic dataset. The three genetic groups designated by this study can be seen in the color-coded circles with dashed lines. (B) *ADMIXTURE* plot of  $K=2, 3, 4,$  and  $5$  for the 78 landraces. The color-coding from (A) is indicated above each sample's ancestry proportion. (C) Geographic distribution of the 78 *circum*-basmati rice varieties with their grouping status color-coded according to (A). (D) Agronomic measurements for the 78 *circum*-basmati rice varieties sorted into the three groups designated by this study. \*\* indicate p-value  $< 0.01$  and \*\*\* indicate p-value  $< 0.001$ .

### Supplemental Figures

**Supplemental Figure 1. Dot plot comparing chromosome 6 of japonica variety Nipponbare to *circum*-aus variety N22 and indica variety R498.**

**Supplemental Figure 2. Distribution of the proportion of missing nucleotides for japonica variety Nipponbare gene models across the orthologous non-japonica genomic regions.**

**Supplemental Figure 3. Effect of coverage threshold to call a deletion and the total number of deletion calls for samples with various genome coverage.**

**Supplemental Figure 4. Insertion time of LTR retrotransposon in various *Oryza* variety group genomes.** Number of annotated LTR retrotransposons is shown above boxplot. The variety group genomes that do not have a significantly different insertion time after a Tukey's range test are indicated with the same letter.

**Supplemental Figure 5. Density of presence-absence variation (PAV) per 500,000 bp window for each chromosome.**

**Supplemental Figure 6. Genome-wide topology weight from 500 SNP size window.** Chromosomal distribution of topology weights involving trios of the *circum-aus*, *circum-basmati*, and *japonica* rices (left), and the sum of the topology weights (right).

**Supplemental Figure 7. 13 demographic models tested by  $\square a \square i$ .**

**Supplemental Figure 8.  $\square a \square i$  model fit for the best-fitting demographic model.** Above row shows the observed and model fit folded site frequency spectrum. Below shows the map and histogram of the residuals.

**Supplemental Figure 9. Neighbor-joining phylogenetic tree of the 78 *circum-basmati* population sample.**

**Supplemental Table 1. Inversion detect by *sniffles* in the Nipponbare reference genome.**

**Supplemental Table 2. The 78 *circum*-basmati samples with Illumina sequencing result used in this study.**

**Supplemental Table 3. Names of the Basmati 334 and Dom Sufid genome gene models that had a deletion frequency of zero across the population.**

**Supplemental Table 4. Names of the Basmati 334 and Dom Sufid genome gene models that had a deletion frequency of above 0.3 and omitted from down stream analysis.**

**Supplemental Table 5. Orthogroup status for the Basmati 334, Dom Sufid, R498, Nipponbare, and N22 genome gene models.**

**Supplemental Table 6. Count and repeat types of the presence-absence variation (PAV) in the Basmati 334 or Dom Sufid genome in comparison to the Nipponbare genome.**

**Supplemental Table 7. Gene ontology results for orthogroups where gene members from the *circum*-basmati are missing.**

**Supplemental Table 8. Gene ontology results for orthogroups where gene members from *circum*-aus, indica, and japonica are missing.**

**Supplemental Table 9. Population frequency across the 78 *circum*-basmati samples for orthogroups that were specifically missing a gene in the Basmati 334 and Dom Sufid**

**genome gene models.**

**Supplemental Table 10. Genome coordinates of the LTR retrotransposons of the Basmati 334 genomes.**

**Supplemental Table 11. Genome coordinates of the LTR retrotransposons of the Dom Sufid genomes.**

**Supplemental Table 12. Genome coordinates of the Gypsy elements indicated with a single star in Figure 3.**

**Supplemental Table 13. Genome coordinates of the Copia elements indicated with a single star in Figure 3.**

**Supplemental Table 14. Genome coordinates of the Gypsy elements indicated with a double star in Figure 3.**

**Supplemental Table 15. Genome coordinates of the Copia elements indicated with a triple star in Figure 3.**

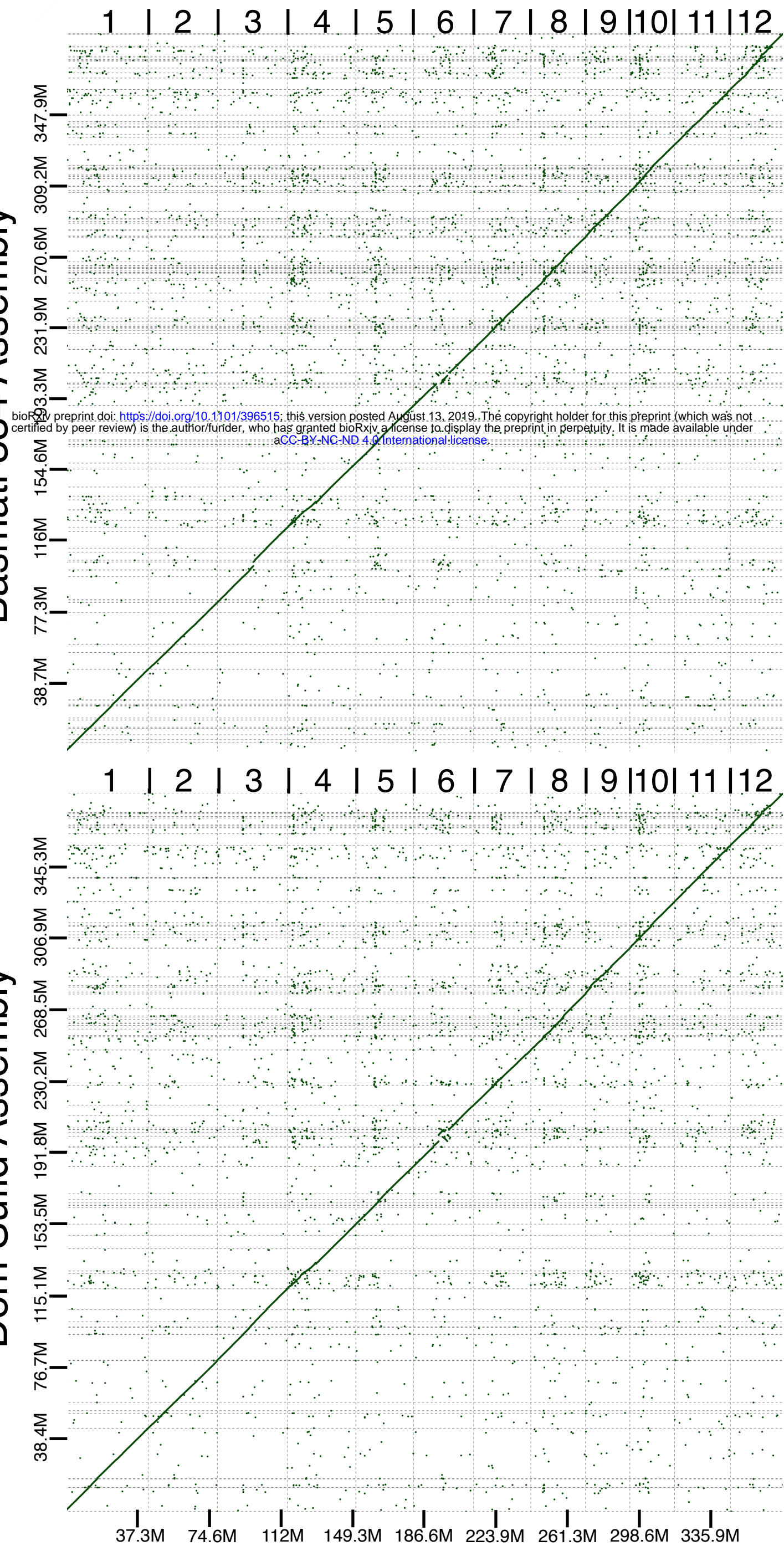
**Supplemental Table 16. The 82 *Oryza* population samples with Illumina sequencing result used in this study.**

**Supplemental Table 17.  $\alpha_i$  parameter estimates for the 13 different demographic models.** See supplemental figure 7 for visualization of the estimating parameters.

# A Nipponbare Chromosome

Basmati 334 Assembly

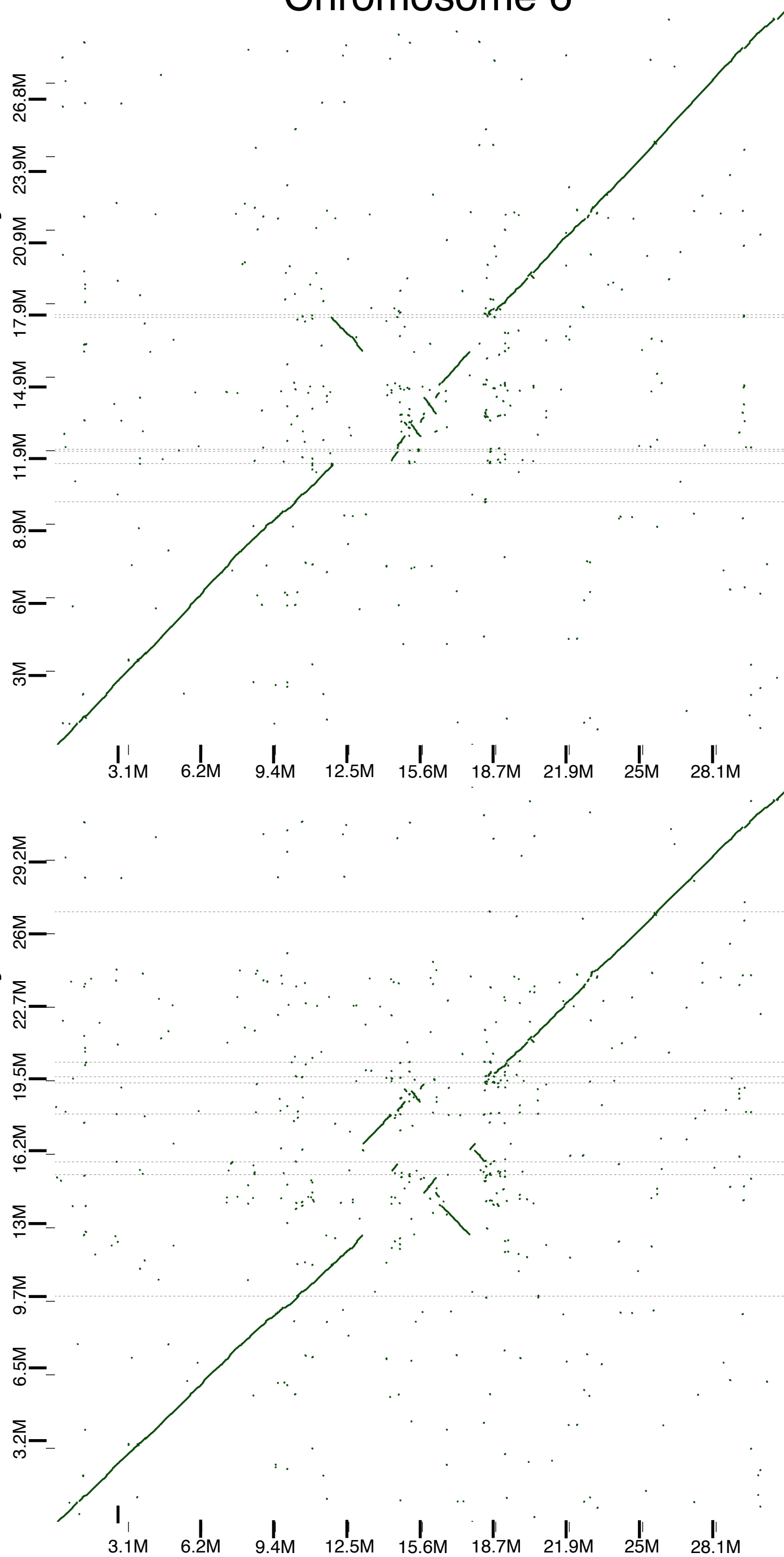
Dom Sufid Assembly



# B Nipponbare Chromosome 6

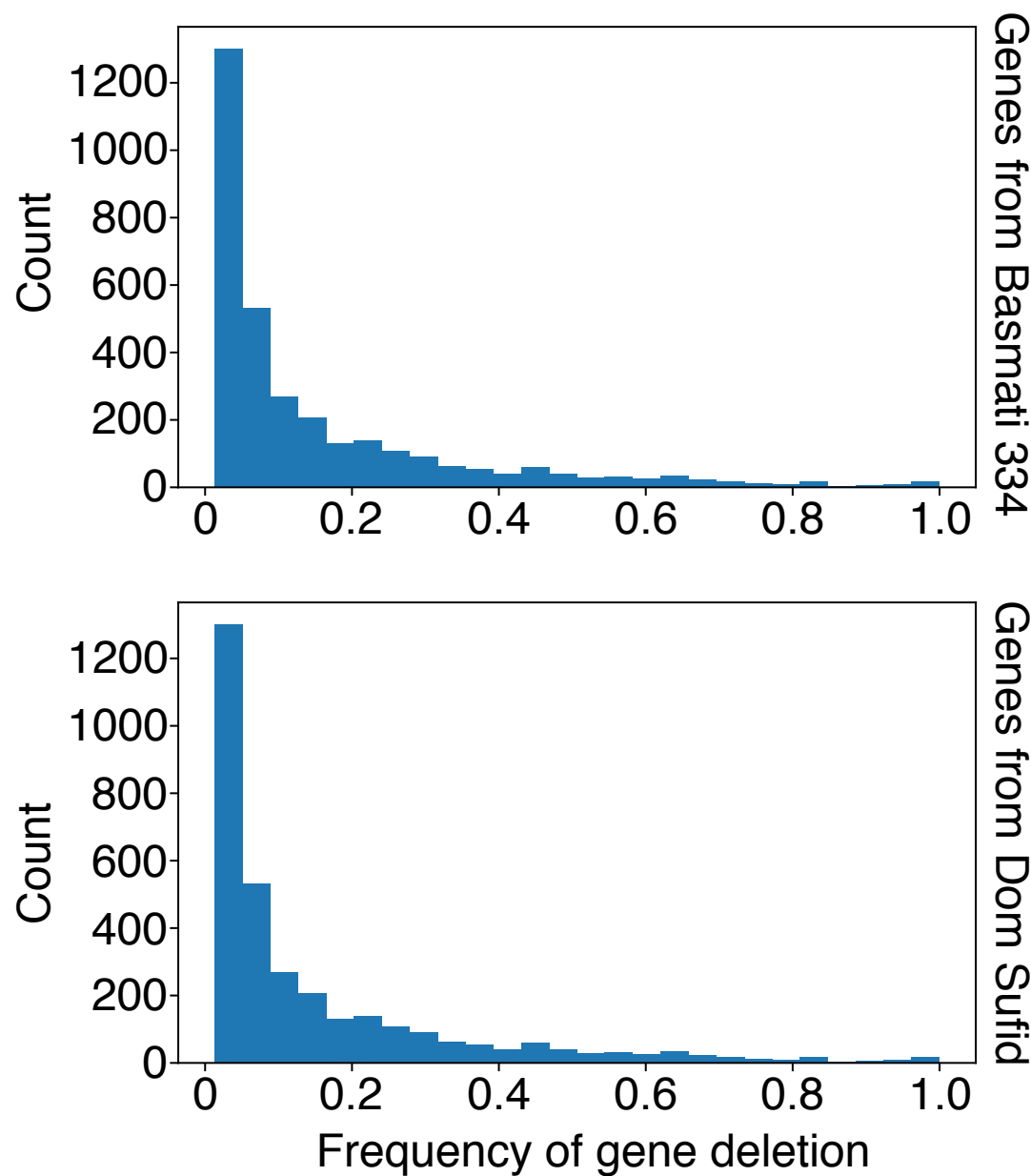
Basmati 334 Assembly

Dom Sufid Assembly

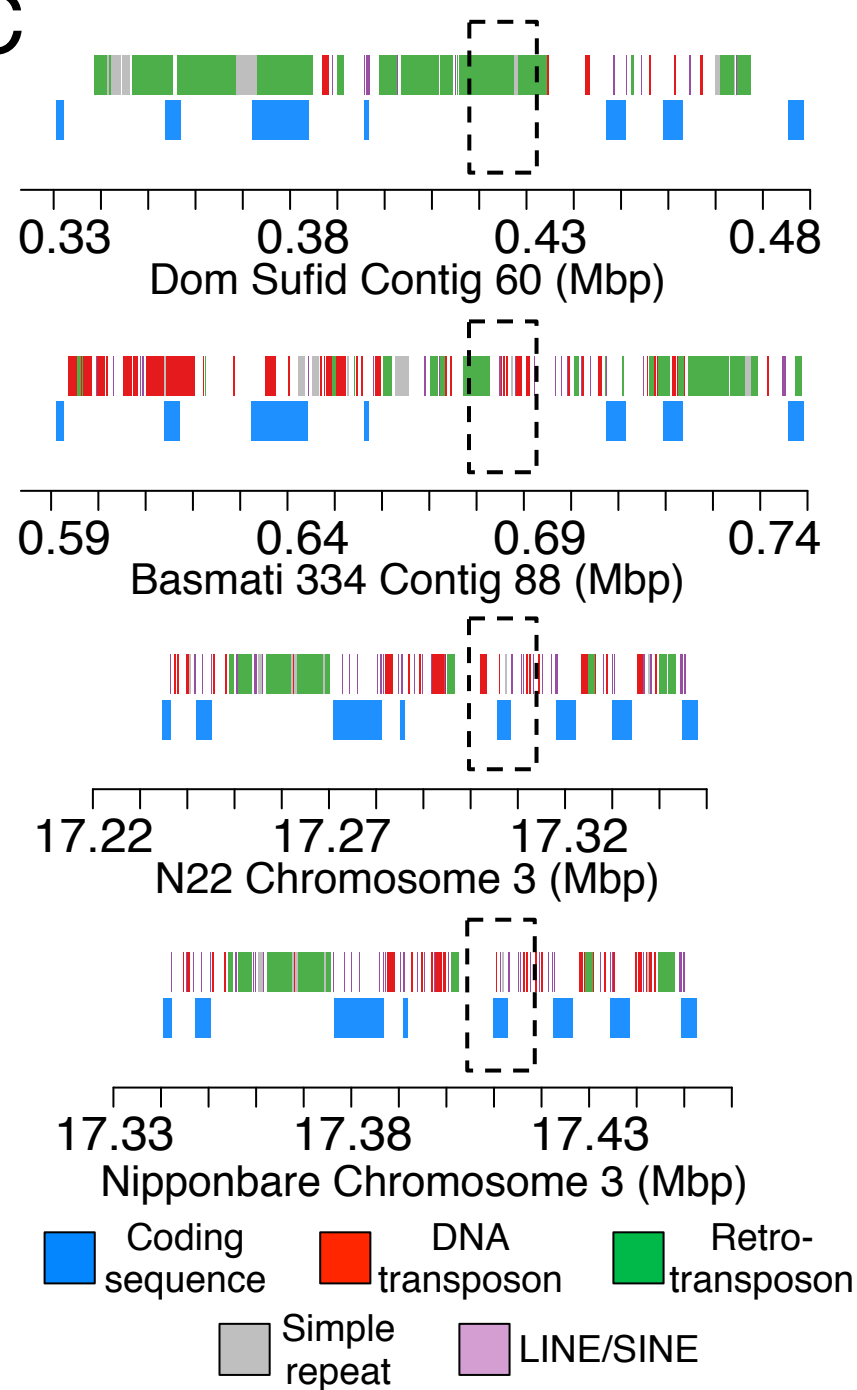




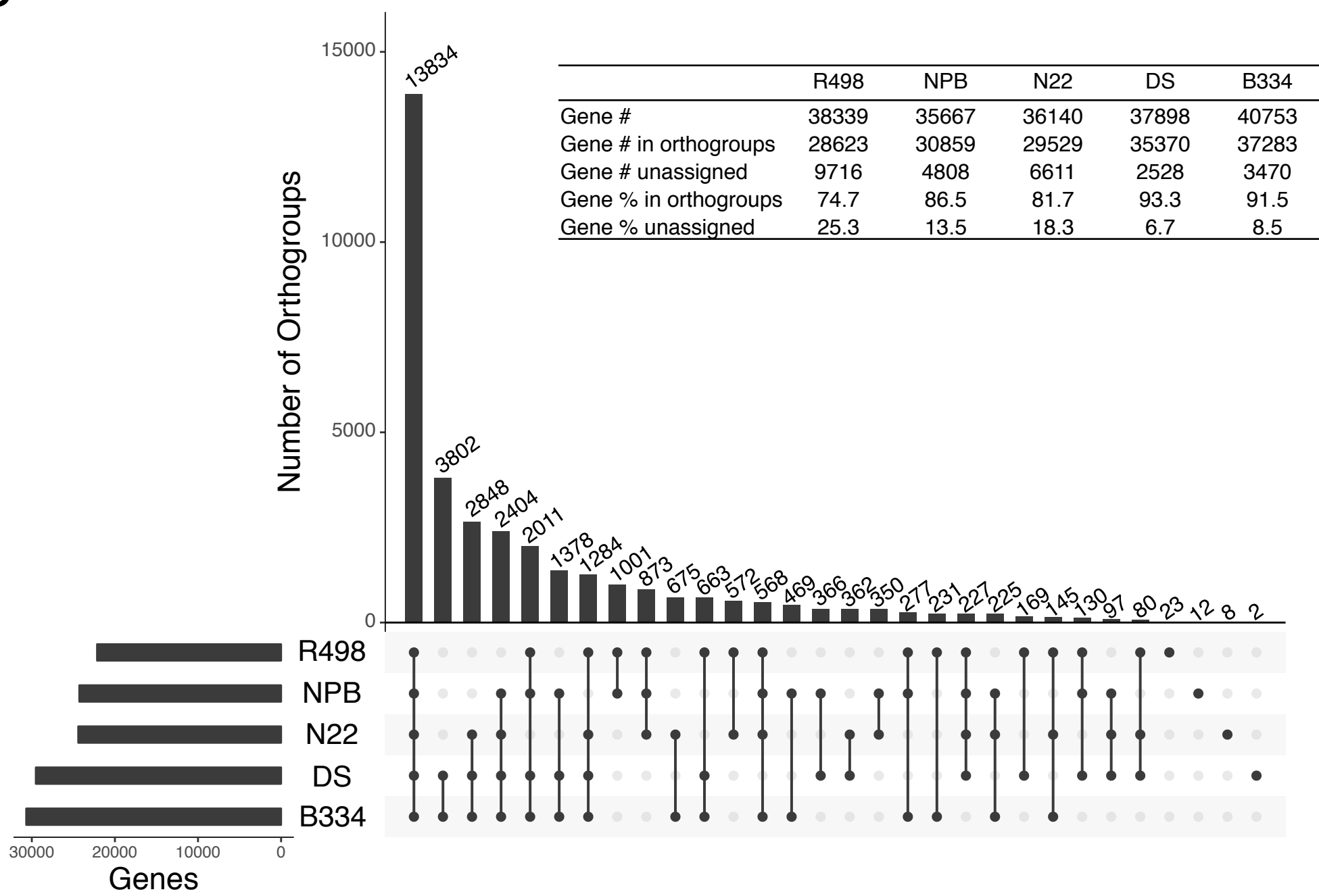
A

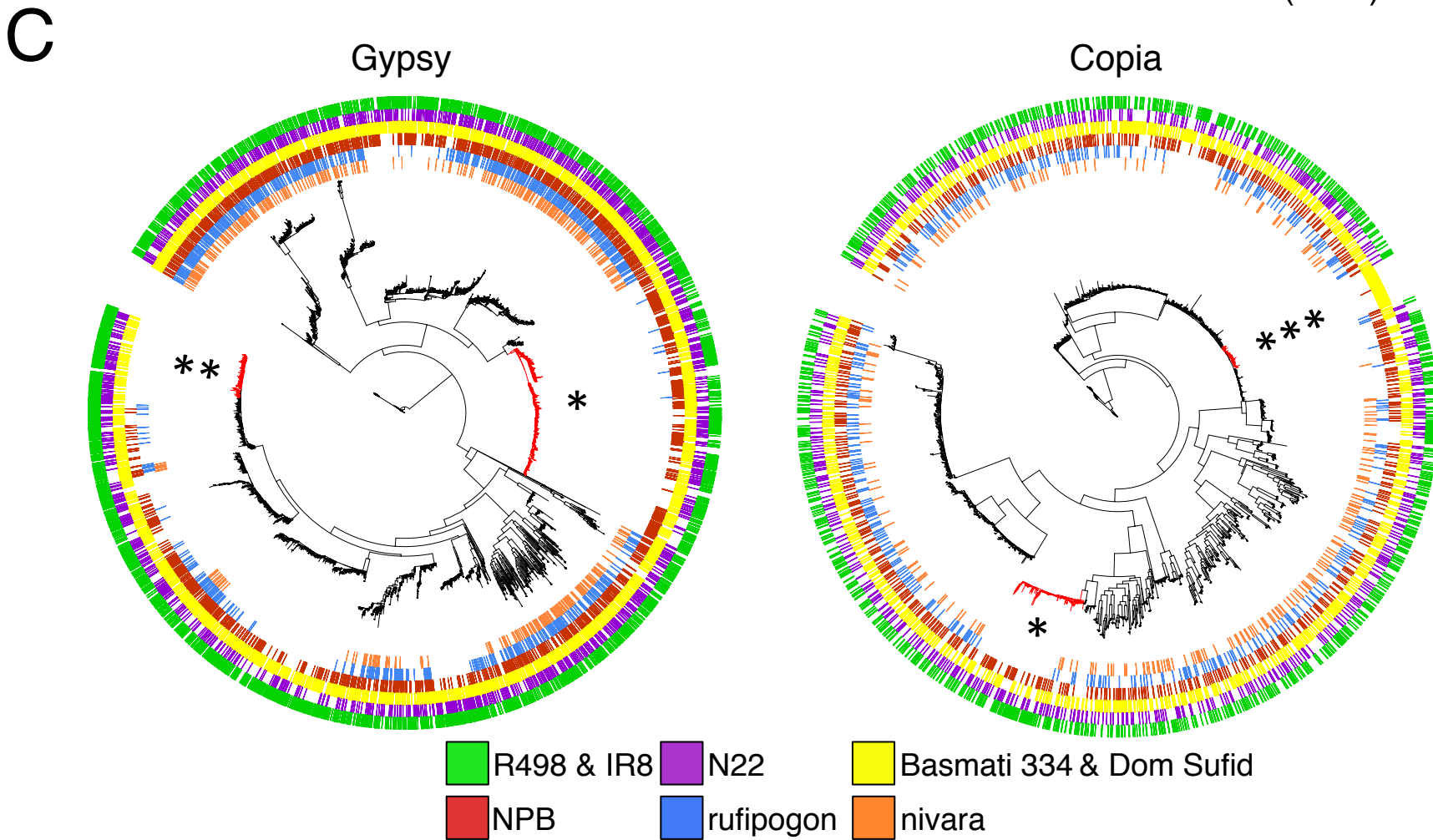
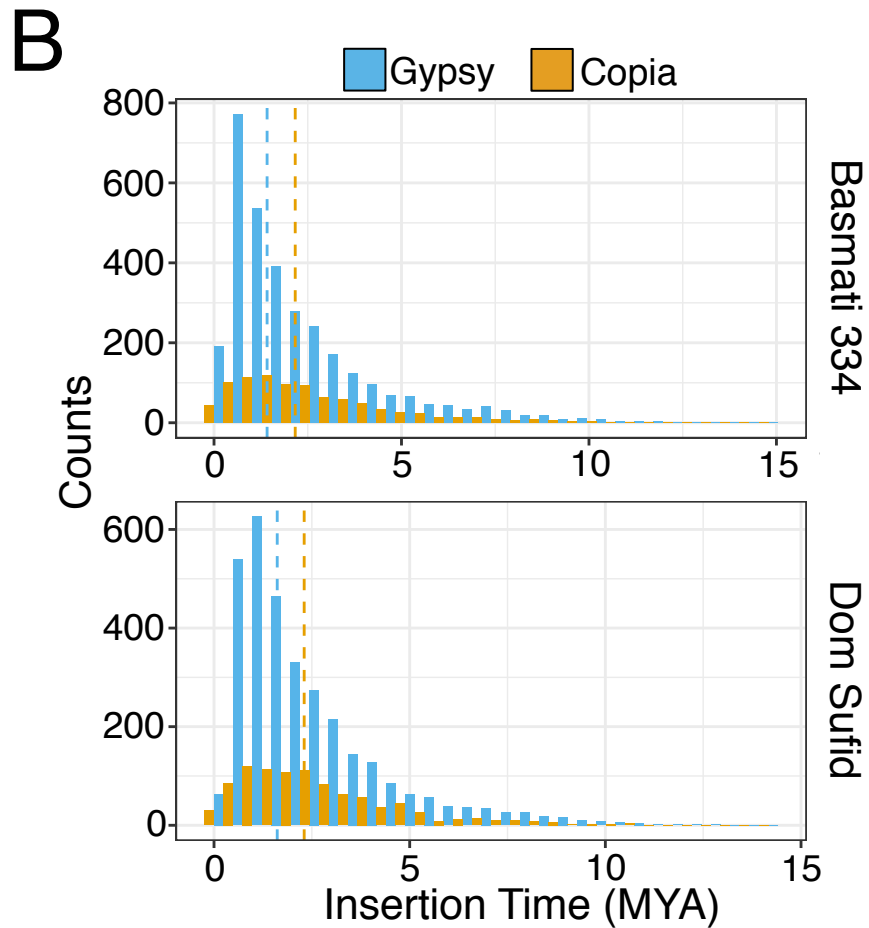
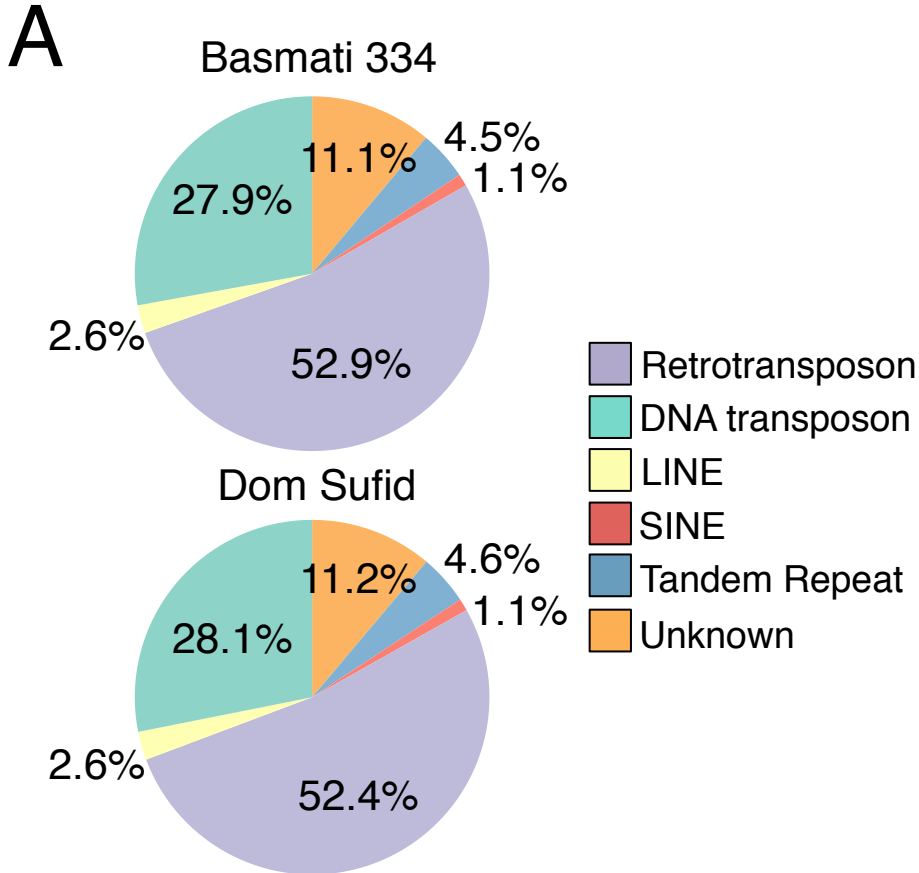


C



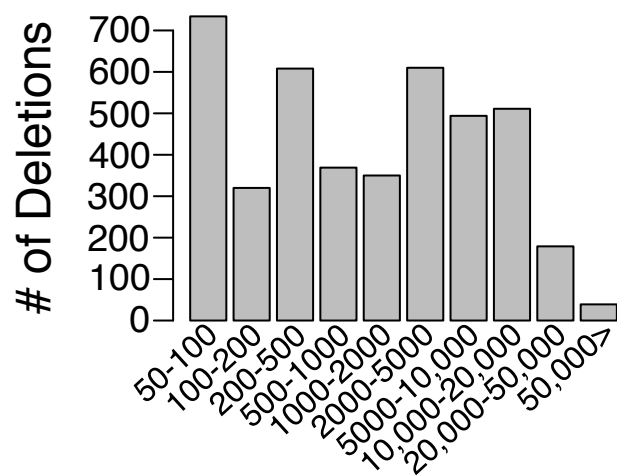
B



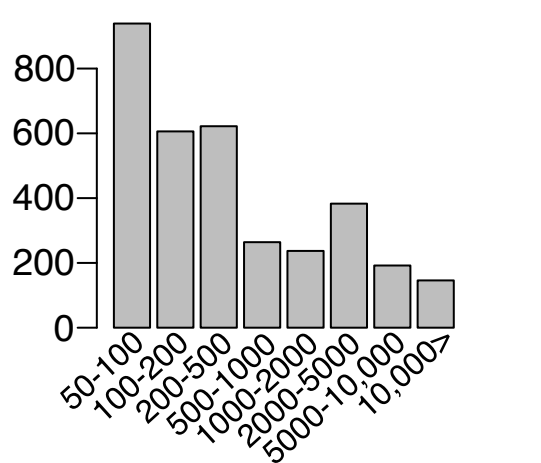
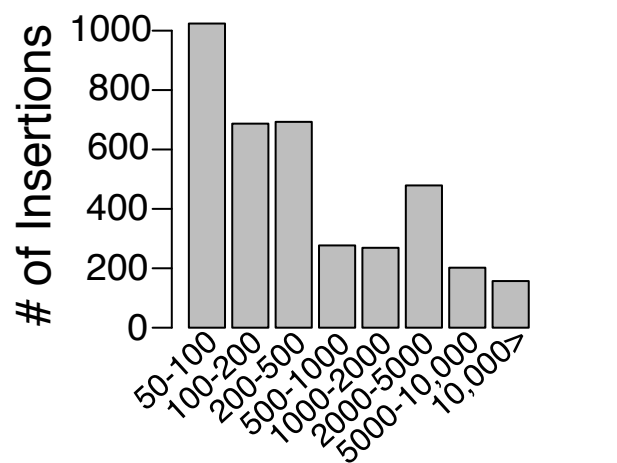
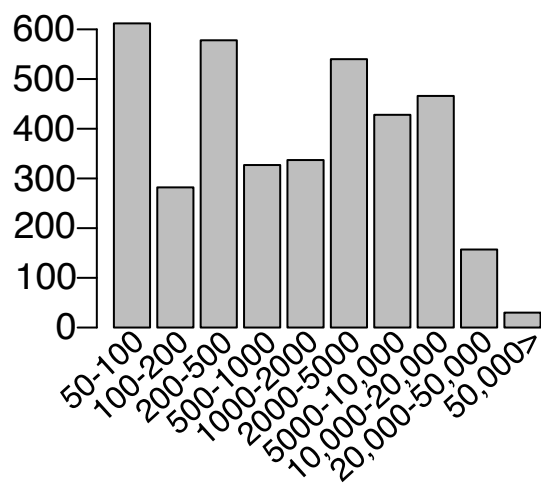


**A**

Basmati 334



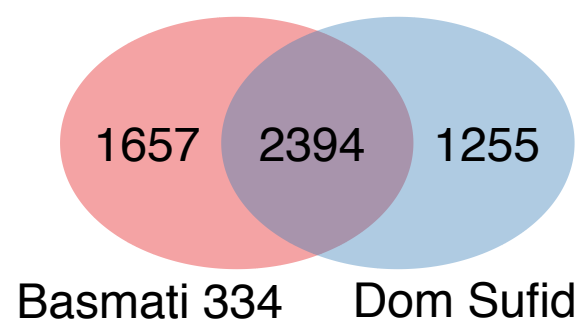
Dom Sufid



Size of INDEL (bp)

**B**

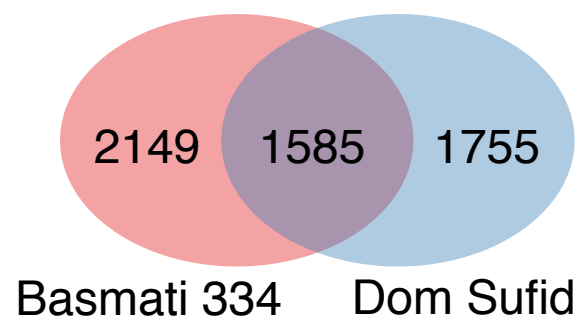
Deletions



Basmati 334

Dom Sufid

Insertions

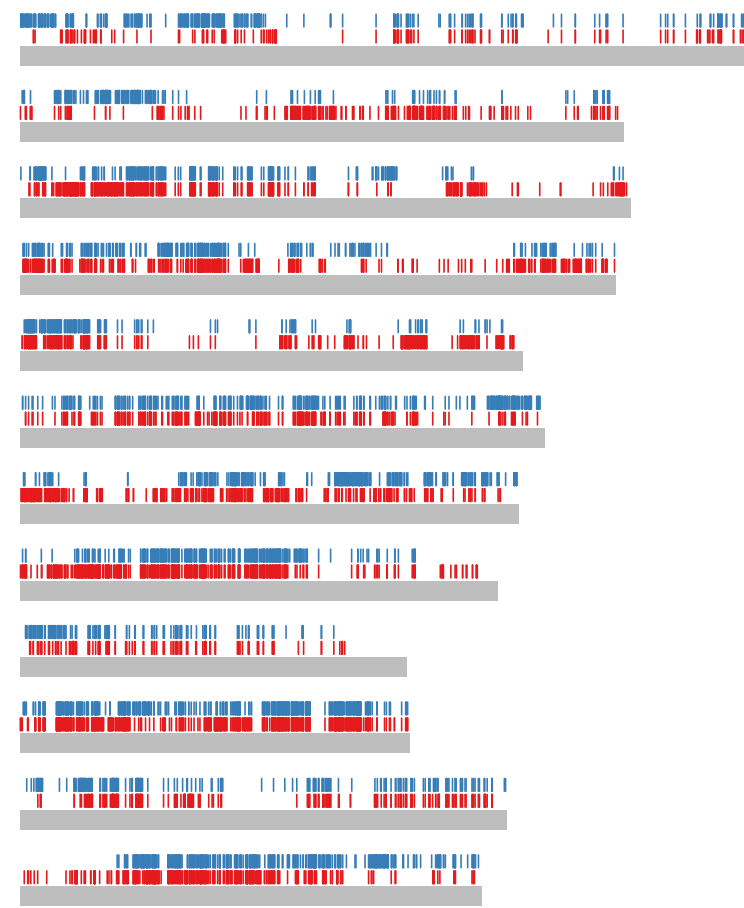


Basmati 334

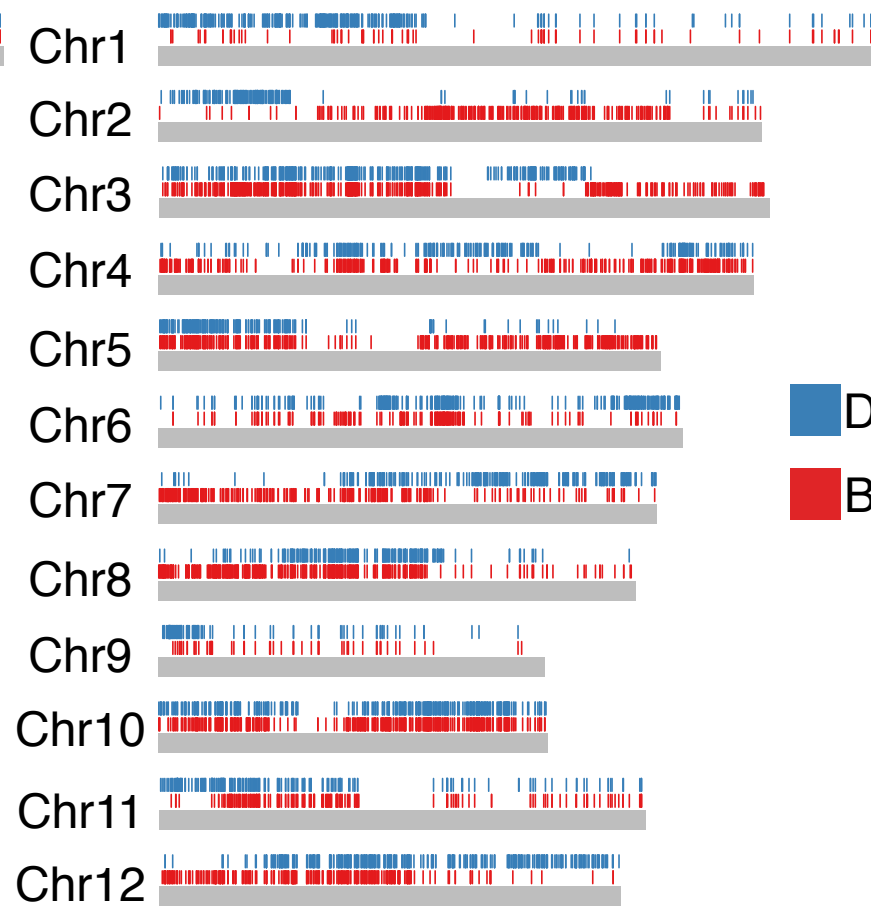
Dom Sufid

**C**

Deletions



Insertions



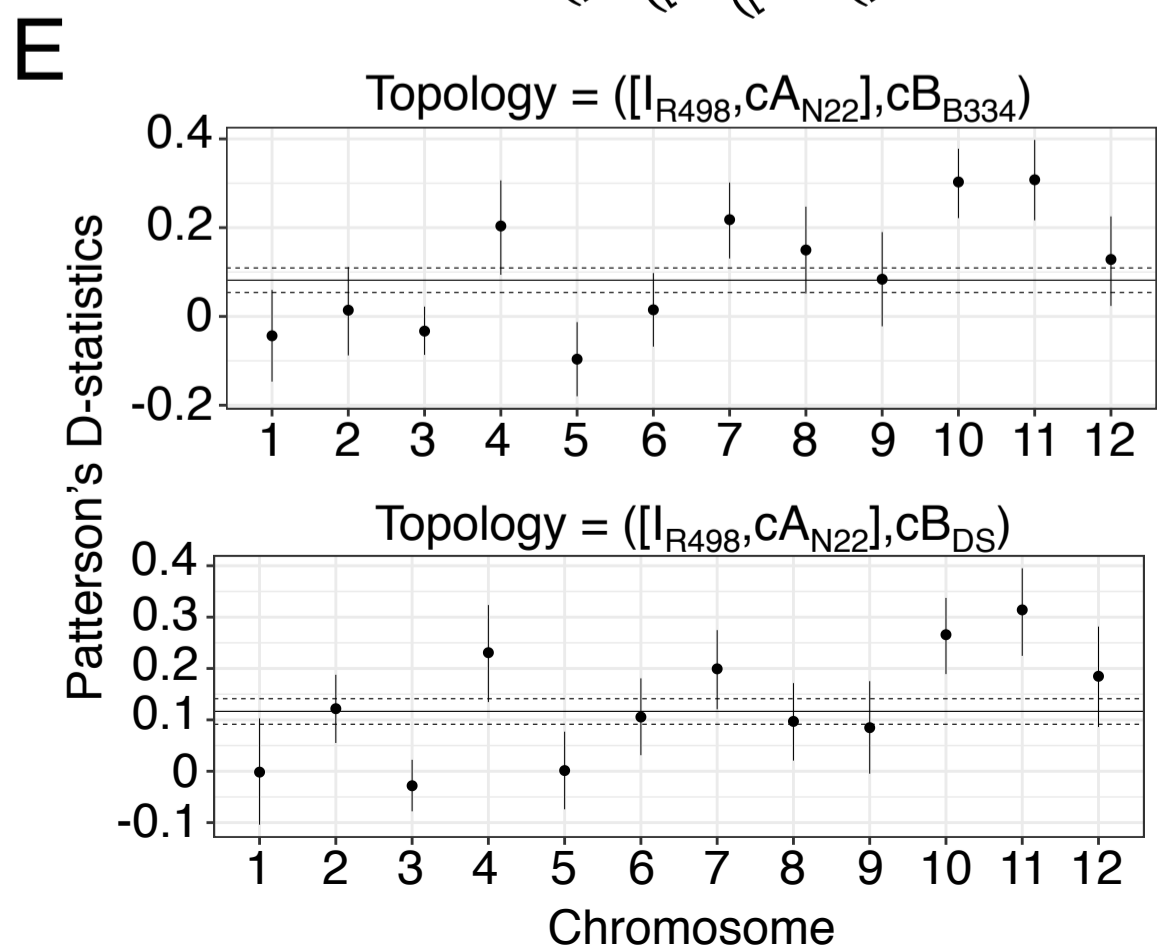
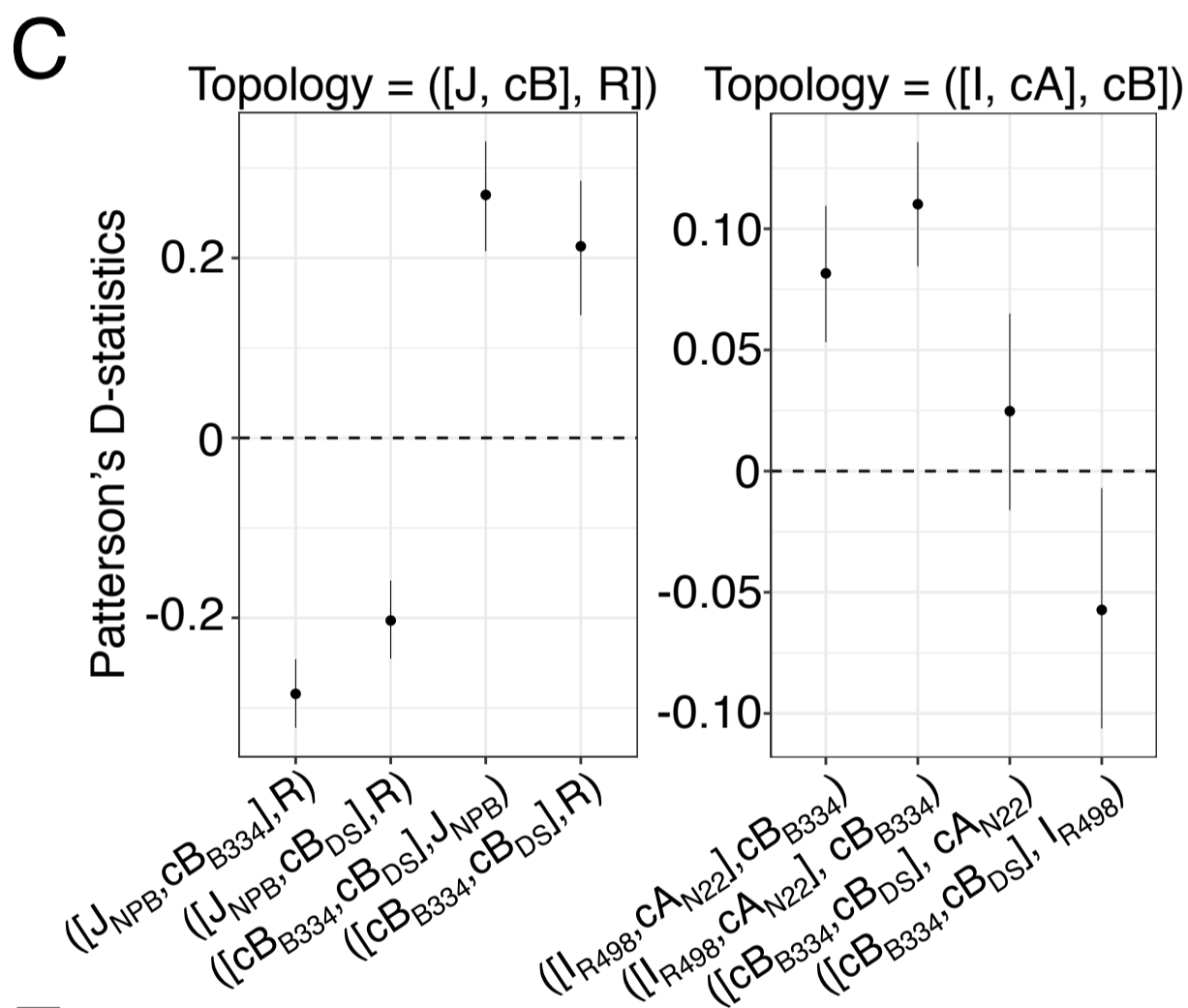
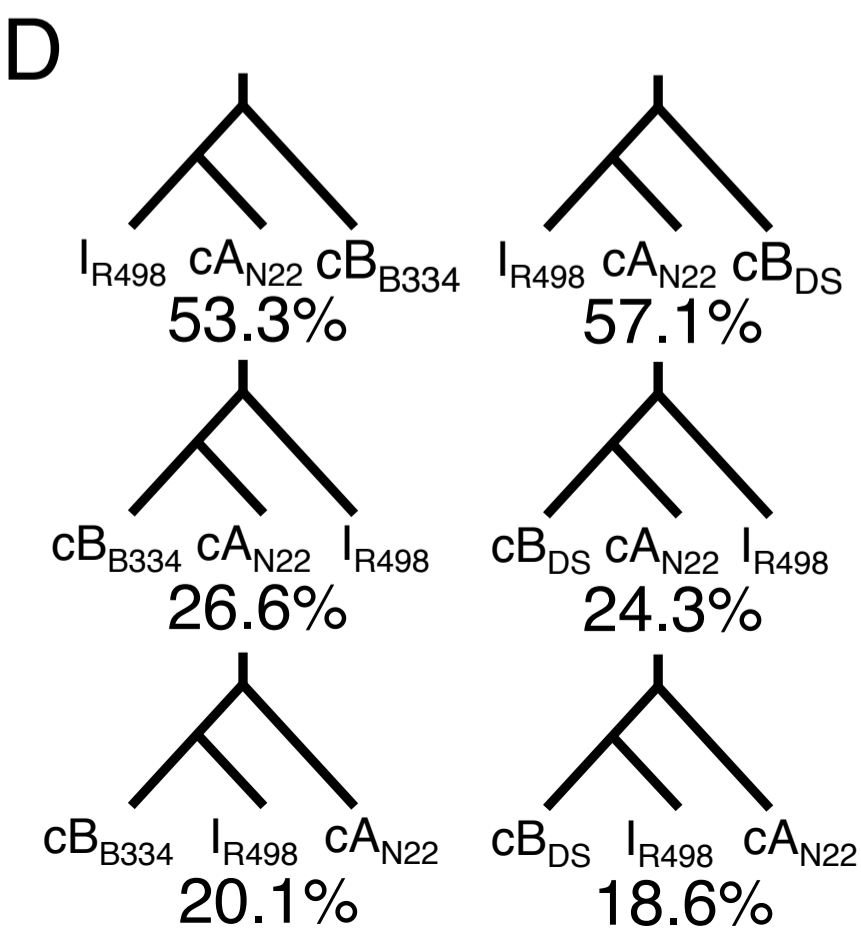
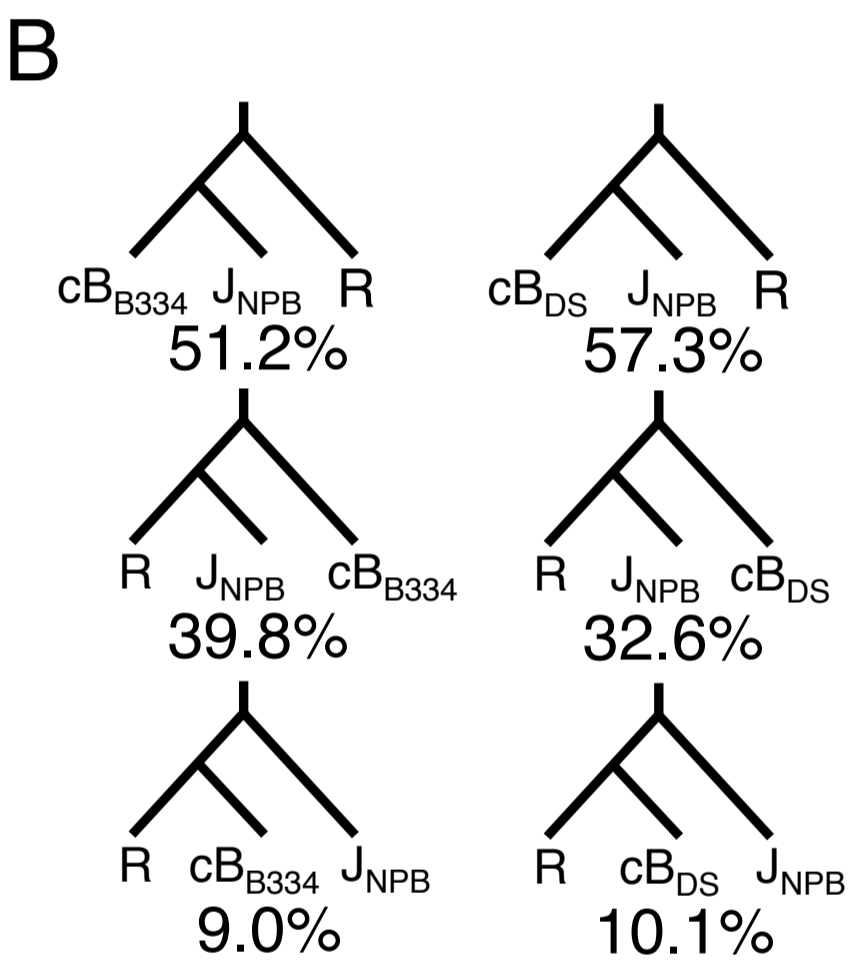
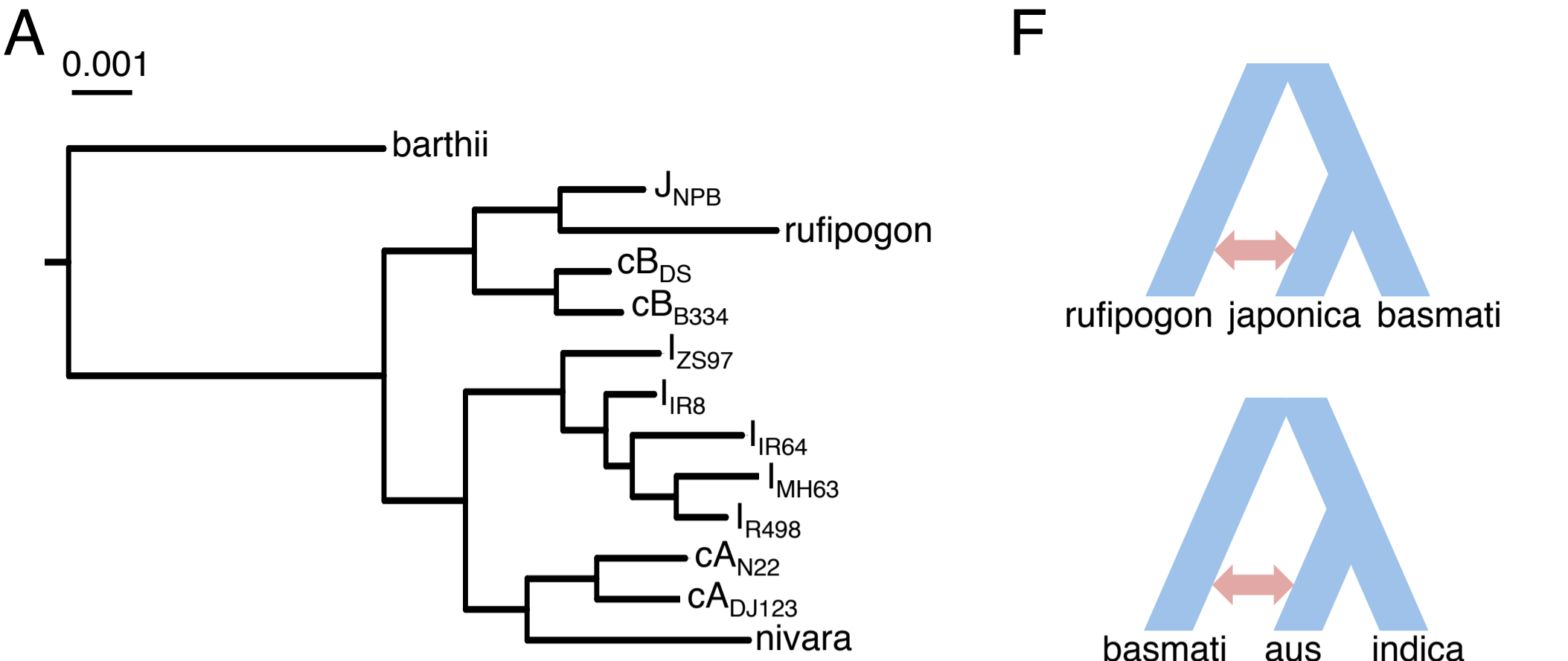
Dom Sufid

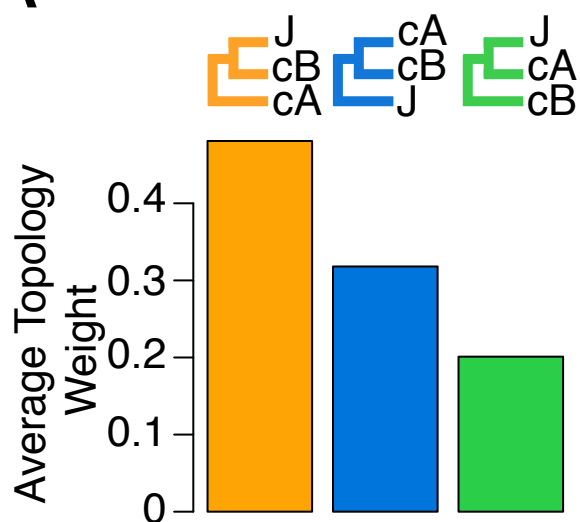
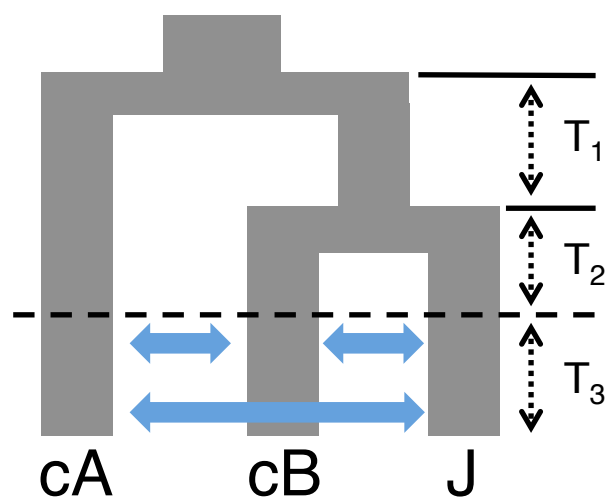
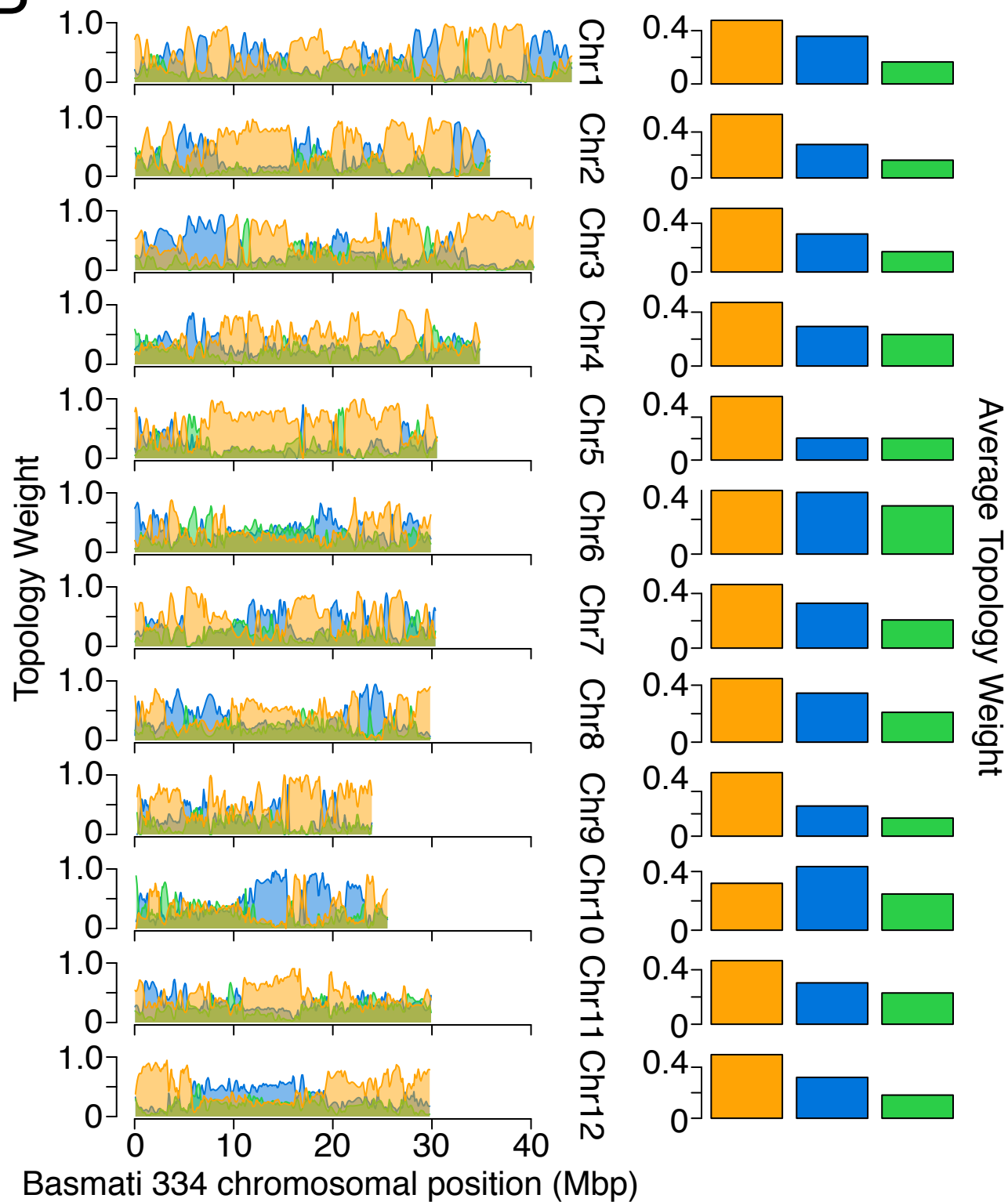
Basmati 334

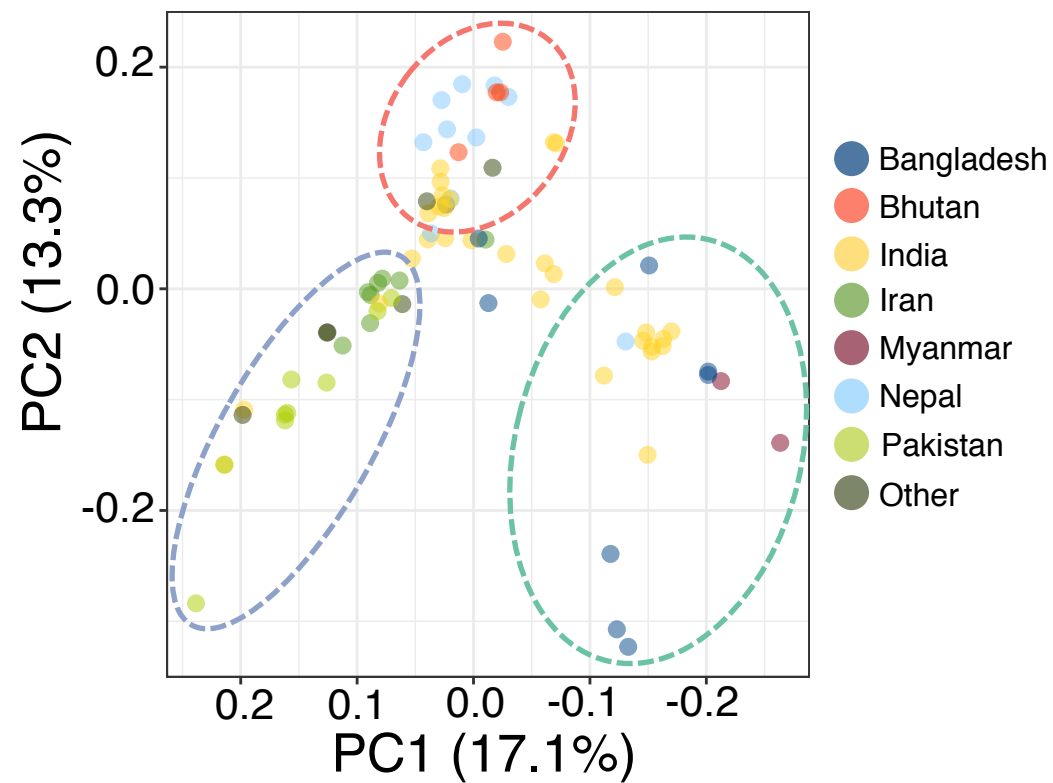
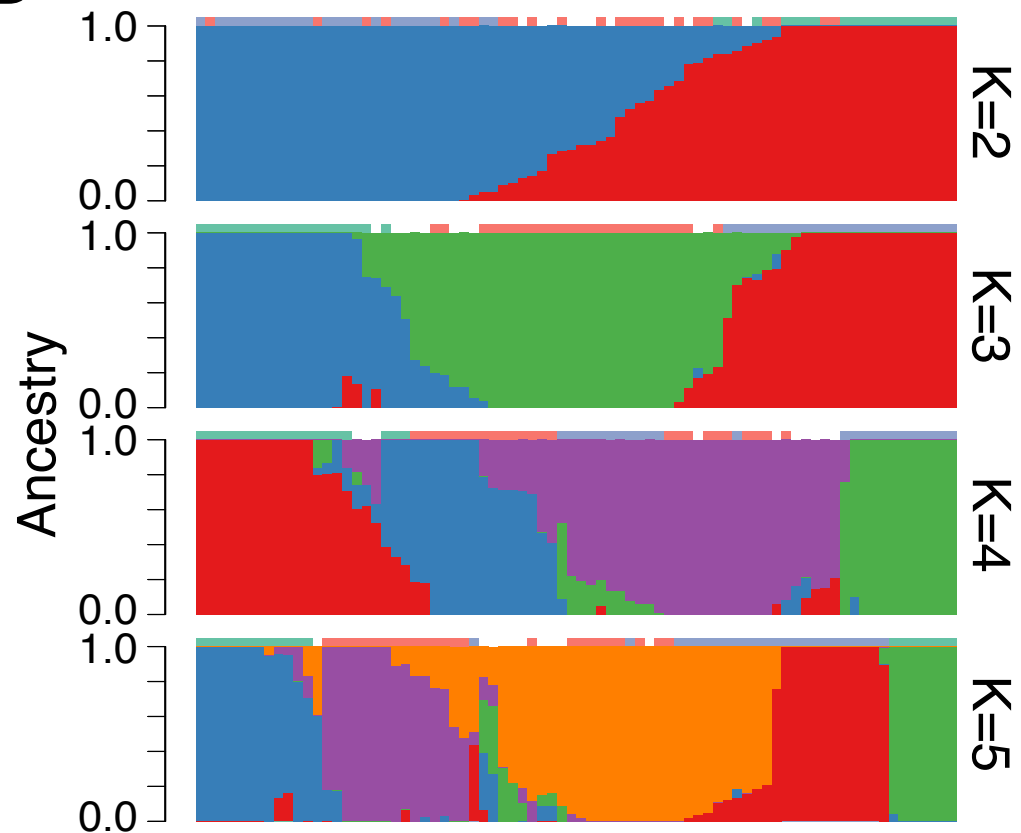
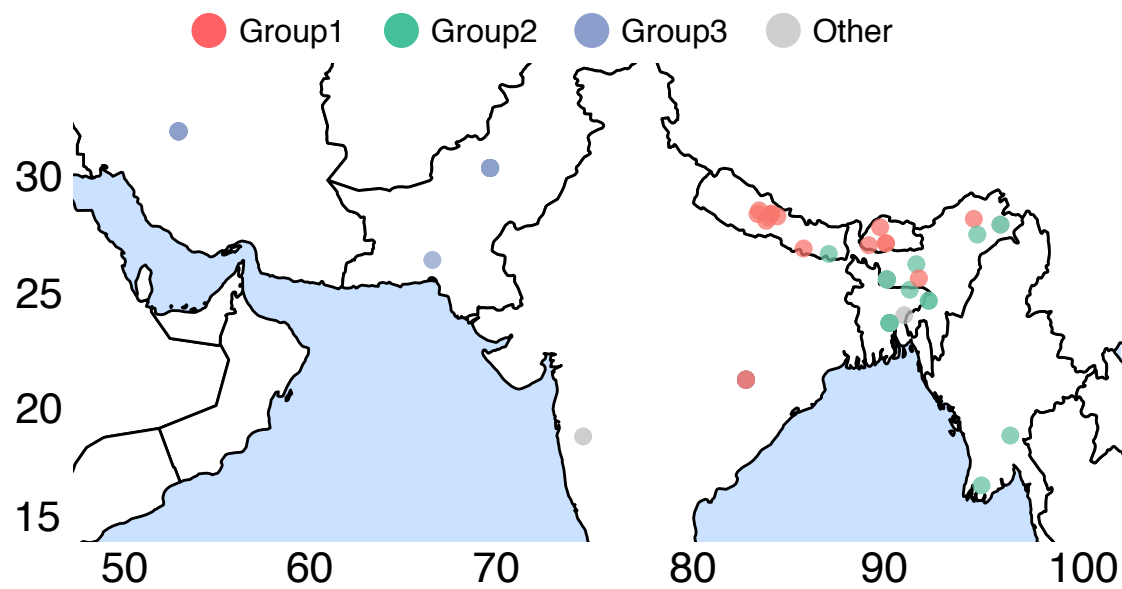
0 8 16 24 32 40

0 8 16 24 32 40

Chromosomal position of Nipponbare reference genome (Mbp)



**A****C****B**

**A****B****C****D**