# Rice Pan-genome Array (RPGA): an efficient genotyping solution for pan-genome-based accelerated crop improvement in rice

Anurag Daware[1], Ankit Malik[2], Rishi Srivastava[1], Durdam Das[1], Ranjith K. Ellur[2], Ashok K. Singh[2], Akhilesh K. Tyagi[3], Swarup K. Parida[1*]

[1]National Institute of Plant Genome Research (NIPGR), Aruna Asaf Ali Marg, New Delhi 110067, India

[2]Division of Genetics, Rice Section, Indian Agricultural Research Institute (IARI), New Delhi 110012, India

[3]Department of Plant Molecular Biology, University of Delhi South Campus, New Delhi 110021, India

**\*Corresponding authors**

Prof. Akhilesh K. Tyagi

Department of Plant Molecular Biology, University of Delhi South Campus, New Delhi 110021, India

E-mail: akhilesh@genomeindia.org

Phone: 91-11-26742267; Fax: 91-11-26741759

Dr. Swarup K. Parida

National Institute of Plant Genome Research (NIPGR), Aruna Asaf Ali Marg, New Delhi - 110067, India.

E-mail: swarup@nipgr.ac.in; swarupdbt@gmail.com

Phone: 91-11-26735228; Fax: 91-11-26741658

Running Title: **Rice Pan-genome Array (RPGA)**

**Word Count:**
**Number of Figures:** 9
**Number Tables:** 00

## ABSTRACT

The advent of the pan-genome era has unraveled previously unknown genetic variation existing within diverse crop plants including rice. This untapped genetic variation is believed to account for a major portion of phenotypic variation existing in crop plants and might be responsible for missing heritability. However, the use of conventional single reference-guided genotyping often fails to capture large portion of this genetic variation leading to a reference bias. This makes it difficult to identify and utilize novel population/cultivar-specific genes for crop improvement. To overcome this challenge, we developed a rice pan-genome genotyping array (RPGA) includes 80K genome-wide SNPs which provides simple, user-friendly and cost-effective solution for rapid pan-genome-based genotyping in rice. The GWAS conducted using RPGA-SNP genotyping data of a rice diversity panel detected total of 42 loci, including previously known as well as novel genomic loci regulating grain size/weight traits in rice. Eight of the identified trait-associated loci (dispensable loci) could not be detected with conventional single reference genome-based GWAS and found to be missing from the commonly used Nipponbare reference genome. WD repeat-containing PROTEIN 12 gene underlying one of such dispensable locus on chromosome 7 (*qLWR7*) along with few other non-dispensable loci was subsequently detected using high-resolution QTL mapping confirming authenticity of RPGA-led GWAS. This demonstrates the potential of RPGA-based genotyping to overcome reference bias. Besides GWAS, the application of RPGA-based genotyping for natural allelic diversity and population structure analysis, seed purity and hybridity testing, ultra-high-density genetic map construction and chromosome level genome assembly, and marker-assisted foreground/background selection was successfully demonstrated. Based on these salient outcomes, a web application (http://www.rpgaweb.com) was also developed to provide easy to use platform for imputation of RPGA-based genotyping data using 3K Rice Reference Panel and subsequent GWAS in order to drive genetic improvement of rice.

**Key Words:** genotyping, GWAS, pan-genome, QTL mapping, SNP array

## INTRODUCTION

Conventionally, a single reference genome sequence per species is considered enough to understand the complete genetic blueprint of a species. However, availability of multiple genomes per species has revealed presence of extensive sequence diversity predominantly in the form of structural variation (SVs), within individuals of the same species. To account for intra-species variation, the concept of pan-genome was proposed (Tettelin et al., 2005). Pan-genome refers to a complete set of genes of a biological clade (for example a species), which includes both core genes (set of genes present in all individuals of a species) and dispensable genes (set of genes which are individual specific or not present in all individuals of a species). Pan-genomes for major crop species including rice, soybean, wheat, tomato, maize, and *Brassica* are now available and all these pan-genomes contains large number of previously unidentified sequence variants and dispensable genes (Hirsch et al., 2014; Montenegro et al., 2017; Sun et al., 2017; Stein et al., 2018; Zhao et al., 2018; Gao et al., 2019; Alonge et al., 2020; Walkowiak et al., 2020).

The functionally characterized dispensable genes so far have been found to perform a diverse range of crucial functions including tolerance to nutrient deficiency, submergence tolerance, organ size regulation, resistance against pathogens (Xu et al., 2006; Ashikawa et al., 2008; Fukuoka et al., 2009; Hattori et al., 2009; Gamuyao et al., 2012; Maron et al., 2013; Wang et al., 2015). Further, compared to core genes, dispensable genes evolve faster (i.e. higher mutation density and higher synonymous to non-synonymous mutation ratio) indicating the vital role of dispensable genes in providing essential diversity for a species to adapt to diverse environmental conditions. Thus, dispensable genes are believed to be a major contributor to phenotypic diversity and adaptive evolution in crop plants. This makes dispensable genes important targets for crop improvement. Some of the dispensable genes like *Sub1A, Pstol*, etc., have already been utilized in rice breeding programs and have become widely popular owing to the enormous impact of these genes on improving overall crop productivity (Xu et al., 2006; Gamuyao et al., 2012). However, the vast majority of dispensable genes remain un-annotated, making it difficult to decipher the role of these genes in regulating important traits in crop plants. The limited knowledge about the function of dispensable genes could be ascribed to sole reliance on a single reference genome for all kinds of genomic mapping studies especially genome-wide association study (GWAS) and quantitative trait loci (QTL) mapping as it leads to reference bias (Coletta et al., 2021). Leveraging pan-genome-based genotyping for genetic mapping studies can potentailly overcome the reference bias imposed by the use of a single reference genome. However, reserachers still relies on comparison of high-quality *de novo* genome assemblies for high-throughput genotyping of SVs, which makes genotyping large number of samples

93 prohebitvely expensive. Thus there is need to decouple discovery and genotyping phase to perform

94 cost-efficent pan-genome-based genotyping. The recent developemnt of pan-genome graphs provide

95 one such alternative. However, these methods still require high-depth whole genome sequencing and

96 also require extensive computational resources making them unsuitable for most plant-breeding

97 applications.

98 In the recent past, high-density SNP genotyping arrays have became widely popular for large-

99 scale genotyping in crop plants due to simple procedure, fast-turnaround time, simple data-analysis

100 and limited requirment of high-performance computing clusters (HPCCs) (Bianco et al., 2016;

101 McCouch et al., 2016; Li et al., 2019). Despite these advantages, SNP arrays developed so far in most

102 crop plants including rice, assays variants identified using a single reference genome and, therefore,

103 can not capture most of the dispensable gene variation existing in rice germplasm (Chen et al., 2014;

104 Singh et al., 2015; McCouch et al., 2016; Schmidt et al., 2017; Thomson et al., 2017; Torkamaneh et

105 al., 2019). This limits the utility of high-density SNP arrays for large-scale pan-genome-based

106 genotyping. Thus, a pan-genome-based SNP-array that can tag genetic variation (SNPs/InDels/SVs)

107 from both core as well as the dispensable genome, will be a crucial step forward toward the practical

108 utilization of pan-genomes in crop improvement.

109 Keeping the aforementioned scenario in mind, here we outline the development of "Rice Pan-

110 genome Genotyping Array" (RPGA), a novel pan-genome-based SNP genotyping array that can

111 efficiently capture haplotype variation from the entire 3K rice pan-genome representing diverse

112 population (*indica,* tropical/temperate *japonica, aus* and aromatic, etc.). We further demonstrate its

113 application for sucessfully overcoming reference bias in high-resolution GWAS and QTL mapping to

114 delineate novel genomic loci modulating traits of agronomic importance (for instance, grain

115 size/weight) in rice. Finally, a user-friendly web portal "Rice Pan-genome Genotyping Array Analysis

116 Portal" (RAP) was also developed to provide researchers with easy to use interface for performing

117 imputation and conducting pan-genome-based GWAS using genotype data generated with RPGA.

118 Thus, RPGA and RAP together provide an end-to-end genotyping solution for accelerated genomics-

119 assisted breeding and crop improvement in rice.

120 **RESULTS**

121 **Designing of Rice Pan-genome Genotyping Array (RPGA)**

122 A rice pan-genome array (80K) design is based on 3K rice pan-genome, which includes the complete

123 Nipponbare genome (IRGSP 1.0/MSU release 7: 373 Mb) and 12 pseudo-chromosomes containing

124 genomic sequence specific to different sub-groups of 4 rice subpopulations (*indica*, *japonica*, *aus* and

125 aromatic) as well as admixed accessions (Sun et al., 2017). RPGA assays total of 80504 SNPs

126    including 60026 SNPs from 12 Nipponbare chromosomes and 20478 SNPs from 12 pseudo-
127    chromosomes of 3K rice pan-genome. The 80504 RPGA SNPs display even distribution throughout
128    the rice pan-genome with ~5 SNPs in every 100 kb genomic interval. In the case of the Nipponbare
129    genome, chromosome 1 being the longest, harbors the highest number of SNPs (7348) whereas,
130    chromosome 9 being the shortest harbors the least (4704) number of SNPs (**Figure 1a, b**). About 25%
131    (22066) and 47% (37960) of the said 80504 RPGA-SNPs are from genic and intergenic regions of the
132    Nipponbare reference genome, respectively (**Figure 1c**). The remaining 28% (20478) of SNPs belong
133    to the twelve pseudo-chromosomes (dispensable/population-specific genomic sequences) from the 3K
134    rice pan-genome (**Figure 1c**).

**Large-scale validation of RPGA**

136    The large-scale validation of RPGA was performed using 188 RILs (Sonasal $\times$ Pusa Basmati 1121 $F_{12}$
137    mapping population) and a diversity panel consisting of 275 accessions representing cultivated and
138    wild Indian rice accessions. In the case of both RILs and diversity panel, very high proportion of SNPs
139    i.e. ~ 57% (46523) and ~ 60% (48133) of total 80504 RPGA-SNPsrespectively , , were classified as
140    Poly High Resolution (PHR). In addition to PHR SNPs, about 15% (12618) and ~7% of SNPs (5746)
141    were classified as Off Target Variants (OTVs) in the diversity panel and RILs, respectively, which
142    were then processed using the OTV caller provided as part of Axiom Analysis Suite to obtain
143    genotype calls (presence/absence) (**Figure 2**). Thus, total of 60751 (75.46%) and 52269 (65.68%) of
144    total 80504 SNPs from the diversity panel (275 accessions) and RILs (188 individuals) genotype
145    datasets, respectively, were selected finally for the downstream analysis (**Supplementary Table 1, 2**).
146    More than 98% concordance between RPGA- and whole genome resequencing-based SNP genotyping
147    data of four accessions.  Further, > 99% concordance was detected for technical replicates genotyped
148    with RPGA, indicates high-reproducibility of RPGA with extremely low genotyping error call rate and
149    thus, can be reliably used for large-scale genotyping in rice. The selected RPGA-SNPs found to have
150    high minor allele frequencies (MAFs > 0.15) suggesting their highly informative nature
151    (**Supplementary Figure 1a**). The design of RPGA and its genotyping potential across natural and
152    mapping population was optimized for capturing maximum level of subpopulation-specific variations
153    and thus expected to provide high-degree of polymorphism irrespective of germplasm (varietal group
154    and sub-populations) used. The SNP markers belonging to pseudo-chromosomes of different sub-
155    populations/varietal groups displayed high polymorphic potential as reflected by their MAFs
156    (**Supplementary Figure 1b**). This suggests a highly informative nature RPGA for genotyping of
157    diverse rice accessions belonging to different rice subpopulations. Further details regarding the large-
158    scale validation potential of RPGA including optimization of RPGA-SNP genotyping call success rate,

159 and reproducibility/informativeness of valid RPGA-SNP genotyping data across accessions/RILs are
160 provided as **Supplementary Results.**

161 **Genetic mapping of sub-population specific (dispensable) sequences from 3K rice pan-genome to**
162 **twelve rice chromosomes**

163 To develop a framework map for 3K rice pan-genome, an ultra-high-density genetic linkage map using
164 the aforesaid RPGA-based SNP genotyping data of 188 RILs (Sonasal $\times$ PB 1121 F$_{12}$) was generated.
165 The RIL genotype contained a total of 46523 PHR SNPs evenly distributed throughout the 12
166 Nipponbare chromosomes as well as in 12 subpopulation-specific pseudo-chromosomes from 3K rice
167 pan-genome. These SNPs provide uniform genomic coverage of the entire rice pan-genome with ~7
168 SNPs/100 kb. Further, filtering for missingness, segregation distortion, the error-rate and duplicate
169 marker was performed to obtain 13792 unique SNPs across 12 Nipponbare chromosomes and 12 sub-
170 population specific pseudo-chromosomes. Initially, only 9100 SNPs belonging to 12 Nipponbare
171 chromosomes were utilized for genetic linkage map construction. This yielded a genetic linkage map
172 spanning a total genetic distance of 1416.4 cM with an average inter-maker distance of 0.1 cM
173 (**Supplementary Table 3**). Chromosomes 1 (LG1) and 7 (LG7) was found to be the longest and
174 shortest spanning 179.4 and 82.6 cM, respectively. .

175 Further, the genetic map was re-constructed using 13793 SNPs, which includes 9100 SNPs
176 from 12 Nipponbare chromosomes as well as 4693 SNPs from 12 pseudo-chromosomes. The genetic
177 linkage map spanned a total genetic distance of 2312.2 cM with an average inter-maker distance of 0.2
178 cM reflecting the ultra-high-density nature of the constructed genetic map (**Figure 3; Table 4**).
179 Chromosomes 1 (LG1) and 7 (LG7) was found to be the longest and shortest spanning 311.5 and
180 137.6 cM, respectively (**Supplementary Table 4**). This ultra-high-density genetic map generated
181 using markers from the entire 3K rice pan-genome enabled the identification of genetic positions of
182 4693 SNPs belonging to subpopulation/pseudo-chromosomes relative to the SNPs from the
183 Nipponbare reference genome. These 4693 SNPs tagged > 93% contigs from 12 sub-population
184 specific pseudo-chromosomes, harboring thousands of important dispensable genes. Therefore,
185 RPGA-based ultra-high-density genetic linkage map provides a good estimate of genetic locations of
186 these contigs from 12 sub-population specific pseudo-chromosomes.

187 Further, approximate physical locations of 4693 SNPs from 12 sub-population specific pseudo-
188 chromosomes relative to Nipponbare reference genome were determined using an integrated approach.
189 This involves BLAST search against six genome assemblies, namely Nipponbare (*japonica*), Nagina
190 22 (*aus*), IR 64 (*indica*), Sonasal, Basmati 334 and Dom Sufid (aromatic), pair-wise LD estimates and
191 genetic positions of SNPs derived from aforementioned high-density genetic linkage map. The

6

192    information on the physical positions of contigs and dispensable genes within these contigs is vital for

193    a diverse range of genomic applications. This is especially important in both GWAS and QTL

194    mapping studies, to identify dispensable novel genes (unable to detect from Nipponbare reference

195    genome) underlying QTLs governing important agronomic traits.

**Validation of RPGA-based ultra-high-density genetic map using *de novo* genome assembly of**

197    **"Sonasal"**

198    To assess accuracy of RPGA-based Sonsal × PB 1121 ultra-high-density genetic map, the *de novo*

199    whole genome assembly of one of the parental accession i.e., "Sonasal" was generated using Oxford

200    Nanopore long-read sequencing. The Sonasal genome assembly (304 contigs) spanning 368.2 Mb with

201    41020 annotated genes. The Sonasal genome assembly revealed high contiguity with N50 of 4.02 Mb

202    and 96.5% recovery of Benchmarking Universal Single-Copy Orthologs (BUSCO). These BUSCO

203    statistics are comparable to previously assembled rice genomes including *japonica* cultivar

204    Nipponbare (98.4%), *indica* cultivar R498 (98.0%), Basmati cultivar Basmati 334 (97%) and *Sadri*

205    cultivar Dom Sufid (97%) (Choi et al., 2020). About 93% (~ 343 Mb) of assembled 368.2 Mb contig

206    sequences of 12 chromosomes of Sonasal genome was efficiently anchored onto 12 linkage groups of

207    an ultra-high-density genetic linkage map. Further, whole-genome sequence alignment revealed high

208    degree of microsynteny between Sonasal and diverse rice genomes including Nipponbare. This

209    ascertains the high-accuracy and usefulness of RPGA-based ultra-high density genetic linkage map for

210    anchoring *de novo* assembled contigs (**Figure 4**). The comparison of high-quality *de novo* Sonasal

211    contigs with RPGA-based ultra-high-density genetic linkage map revealed high concordance in marker

212    order with > 99.8% markers agreeing with each other. This confirmed the high-quality of RPGA-based

213    ultra-high density genetic linkage map for their immense use in developing the high-quality

214    chromosome-level genome assemblies as well as identification/mapping of sub-population specific

215    dispensable novel sequences (genes) from 3K rice pan-genome on chromosomes of rice.

**Assessment of natural allelic diversity, evolutionary pattern and population structure in a**

217    **diversity panel**

218    The principal component analysis (PCA) was conducted using the RPGA-SNP genotyping data of 271

219    rice accessions alone as well as along with 3045 accessions from 3K rice genome ( Wang *et al.*,

220    2018b). The first two principal components (PCs) explained > 90% of genetic variation. The first PC

221    separated *indica* distinctly from *japonica* and aromatic/traditional Basmati accessions, whereas the

222    second PC separated *aus* and *indica* accessions (**Figure 5a, b**). The majority of these rice accessions

223    belonged to *indica*, *aus* and aromatic subpopulations, whereas only few are represented from the

224    *japonica* subpopulation. Similar to previous reports, *indica* accessions were found to be genetically

7

225    closer to *aus*. Interestingly, both Indian aromatic*/*traditional Basmati and evolved Basmati accessions

226    clustered distinctly from aromatic/Basmati accessions present in the 3K rice reference panel. The

227    Indian aromatic*/*traditional Basmati accessions were found genetically closer with *japonica* and *aus*

228    accessions.    Whereas    the    evolved    Basmati    were    found    to    cluster    between    *indica*    and

229    aromatic*/*traditional Basmati accessions, following their evolution from cross-hybridization between

230    traditional Basmati varieties with superior grain qualities and *indica* varieties with dwarf height, early

231    flowering and high yield (Singh et al., 2018).

232          Subsequently, the RPGA-SNP genotyping data generated for 271 rice accessions (60751

233    SNPs) was further analyzed using fastSTRUCTURE to unravel the fine population structure among

234    rice accessions under study. The results revealed the existence of five distinct population groups in

235    rice accessions. These population groups corresponded to two *indica* sub-groups (*INDI* and *INDII*),

236    and one population group each of *aus*, aromatic/traditional Basmati (ARO*/*TR-BAS) and evolved

237    Basmati (EV-BAS), respectively. Based on ancestry cutoff of 65%, accessions were classified into one

238    of the five distinct population groups (*INDI*, *INDII*, *AUS*, ARO*/*TR-BAS and EV-BAS) and admixture

239    classes, admixed *Indica* (*INDI - INDII*), admixed *Indica-aus*, and admixed *Indica*-aromatic (**Figure**

240    **6a**). The two *indica* subpopulation groups corresponded to *Xian*/*Indica-2* (*XI*-2) and *XI-3* from South

241    Asia and Southeast Asia, respectively, which are previously reported along with two other *indica*

242    subpopulation groups (*XI-1A* from East Asia and*XI-1B* constitutes modern varieties of diverse origins)

243    (Wang et al., 2018b). Further, the *indica* rice germplasm accessions especially from India belong to

244    either of the two aforesaid *indica* subpopulations or evolved as a result of cross-hybridization within

245    these subpopulations or with any of the remaining subpopulations. This suggests the indigenous nature

246    of the most widely cultivated Indian *indica* rice accessions. Apart from this, traditional Basmati

247    accessions and aromatic landraces from north-eastern India were found to cluster together as a group

248    close to traditional Basmati distinct from modern evolved Basmati/aromatic accessions. Contrary to

249    this, evolved Basmati/aromatic accessions were closely related to the *IND1* subpopulation, confirming

250    their origin from cross-hybridization between *IND1* and traditional Basmati accessions. The neighbor-

251    joining (NJ) phylogenetic tree generated also confirmed the existence of five distinct population

252    groups apart from three admixed populations (**Figure 6b**).

**The utility of RPGA for conducting pan-genome-based GWAS**

254    To establish the utility of RPGA for conducting pan-genome-based GWAS, the analysis was

255    performed using the RPGA-SNP genotyping data of 203 germplasm accessions (selected out of the

256    aforesaid 275 rice accessions), phenotyped for four different grain size/weight traits, i.e. grain length,

257    grain width, length-to-width ratio and thousand-grain weight. These 203 accessions displayed a wide

phenotypic variation for all four grain size/weight traits measured. The grain length of accessions ranged from 4.4 to 8.3 mm with a mean ± S.D of 5.9 ± 0.8 mm, whereas, grain width varied from 1.3 to 2.9 mm with a mean ± S.D of 1.9 ± 0.3 mm. The thousand-grain weight of accessions ranged from 11.6 to 41.3 g with a mean ± S.D of 23.1 ± 4.5 g (**Supplementary Figure 2, 3**). Pan-genome-based GWAS was performed using the RPGA genotyping data of 63002 SNPs that includes 48256 SNPs from the Nipponbare reference genome and 14746 SNPs from the subpopulation-specific pseudo-chromosomes, across 206 rice accessions. The pan-genome-based GWAS detected a number of significant associations for all four studied grain size/weight traits. These associations belong to Nipponbare chromosomes as well as sub-population-specific pseudo-chromosomes, which together constitute the 3K rice pan-genome.

The pan-genome-based GWAS detected many previously known gene loci regulating grain size and weight of rice. These include major grain size gene *GRAIN SIZE 3* (*GS3*) which is known to regulate grain size i.e., grain length, grain width and length-to-width ratio (Fan et al., 2009; Takano-Kai et al., 2009; Lu et al., 2013).  Apart from *GS3, POSITIVE REGULATOR OF GRAIN LENGTH 1 (PGL1)*, a known positive regulator of grain length in rice is also found to be associated with grain length (Heang and Sassa, 2012a,b). Similarly, previously known major grain size/weight genes, *GW5* and *OsPPKL1* were also found strongly associated with both grain width as well as thousand-grain weight (Weng et al., 2008; Gao et al. 2019). These results not only reaffirmed an important role of previously known grain size/weight genes such as *GS3*, *GW5* and *OsPPKL1* genes in regulating grain size and grain weight in rice accessions but also validated the authenticity of RPGA-based marker-trait association study (**Figure** 7; **Supplementary Figure 4, 5, 6, 7; Supplementary Table 5, 6, 7, 8, 9**).

In addition to the aforementioned known genes, many novel genomic loci including 10 loci for grain length, 5 loci for grain width, 3 loci for length-to-width ratio and 9 loci for thousand-grain weight, were detected from the Nipponbare reference genome (**Figure 7; Supplementary Figure 4, 5, 6, 7; Supplementary Table 5, 6, 7, 8, 9**). These include transcription regulators like *OsGAMYB1* (*LOC_Os01g59660*) and *OsRUB1* (*LOC_Os10g11260*) genes involved in ubiquitin proteasome pathway as well as the serine/threonine-protein kinase gene (*LOC_Os11g40970*) modulating grain size/weight in rice (**Details in Supplementary Results**).

**RPGA-based GWAS successfully uncovers novel sub-population specific (dispensable) genes regulating grain size/weight in rice**

Apart from loci detected from the 12 Nipponbare chromosomes, multiple SNPs/contigs of pseudo-chromosomes were found to be significantly associated with one or more of the 4 target grain size/weight traits that are repeatedly detected with different statistical models of GWAS. These

291     include 2 loci each for grain length, grain width and grain length-to-width ratio, and 3 loci for

292     thousand-grain weight in rice. Since these SNPs/contigs are known to be partially or completely

293     missing from the Nipponabre reference genome, the physical positions of these contigs were first

294     determined using an integrated genomic approach which leverages information from RPGA-based

295     genetic linkage map, alignment of sub-population specific contigs to multiple reference genome

296     assemblies and pair-wise LD across rice pan-genome (Detailed strategy as per Materials and methods).

297     The candidate genes underlying these loci associated with grain size/weight were subsequently

298     identified. Here three of these vital trait associations for rice grain size/weight are discussed in further

299     details.

### *GW5* associated with rice thousand grain weight

301     The prominent locus detected from the pseudo-chromosome IG5 (MSP-unaln_IG5~14354843) was

302     found to be strongly associated with thousand-grain weight. Further, physical location of this locus

303     was identified using the aforesaid integrated approach. Based on this analysis, the physical position of

304     target locus coincides with the 1212 bp InDel (SV) located upstream of *GW5* (calmodulin binding

305     protein-encoding gene) that has been identified previously as a causal variation responsible for

306     regulating grain size/weight (Liu et al., 2017). This further proves the efficacy of the RPGA-based

307     GWAS approach to efficiently identify SVs associated with important agronomic traits in rice.

### *WDR12* associated with rice grain length and grain length-to-width ratio

309     Apart from this, a SNP from *japonica* group 10 (JG10) pseudo-chromosome (MSP-

310     unaln_JG10~6950498) was repeatedly detected for grain length as well as grain length-to-width ratio

311     (**Figure 7; Supplementary Figure 4, 5, 6, 7; Supplementary Table 5, 6, 7, 8, 9**). The location of

312     SNP unaln_JG10~6950498 was detected on chromosome 7 with high confidence. Further, to identify

313     underlying candidate genes for this QTL, 5 kb flanking sequences from the location of SNPs were

314     retrieved from genome assemblies of seven different rice accessions belonging to different rice

315     subpopulations. This includes *japonica* (Nipponbare), *indica* (IR 64 and R 498), *aus* (Nagina 22) and

316     aromatic (Sonasal, Basmati 334 and Dom Sufid) subpopulations. The extracted sequences were then

317     subjected to multiple sequence alignment along with the full contig sequences harboring the lead SNPs

318     to determine any variation present within this region, between the aforementioned seven rice

319     accessions. The SNP-unaln_JG10~6950498 spanned the upstream/promoter regions of two

320     consecutive genes encoding WD repeat-containing PROTEIN 12 (*LOC_Os07g40930*) and X8

321     domain-containing PLASMODESMATA CALLOSE-BINDING PROTEIN (*LOC_Os07g40940*). A

322     closer look at multiple sequence alignment revealed the presence of multiple InDels within 1.5 kb

323     upstream region of *LOC_Os07g40930*. These InDels are likely to impact the promoter of

10

324 *LOC_Os07g40930* possibly modulating its expression. Further, the InterPro scan detected multiple

325 WD40 repeats confirming *LOC_Os07g40930* as a member of the WD-repeat WDR12/Ytm1 family.

326 Interestingly, PESCADILLO (PES), a conserved nucleolar protein involved in ribosome biogenesis, is

327 known to perform its function by interacting with BLOCK OF PROLIFERATION 1 (BOP1) and WD

328 REPEAT DOMAIN 12 (*WDR12*) to form the PeBoW (PES-BOP1-WDR12) complex. The reduction

329 in PeBoW proteins was found to suppress cell proliferation and cell expansion in *Arabidopsis*

330 underling its vital role in the regulation of cell cycle for growth and development in plants (Cho et al.,

331 2013; Ahn et al., 2016). Apart from this, members of the WD40 repeat superfamily are known F-box

332 proteins, a structural component of SKP1/Cullin/F-box (SCF) E3 ubiquitin ligase complex. These

333 include *Arabidopsis STERILE APETALA* (*SAP*) gene and its ortholog from cucumber *LITTLELEAF*

334 *(LL)*. Both of these genes encode WD40 repeat domain-containing F-box proteins and regulate organ

335 size further emphasizing the vital role of WD40 repeat domain-containing proteins in plant growth and

336 development (Wang et al., 2016c; Yang et al., 2018). This evidence indicates the probable role of

337 *LOC_Os07g40930* in cell number/cell size regulation in rice and, therefore, *LOC_Os07g40930* seems

338 to be a probable candidate gene regulating grain length in rice.

**GRAS family transcription factor associated with grain length**

340 For grain length, another important locus was detected on the pseudo-chromosome JG10 (MSP-

341 unaln_JG10~4245875) (**Figure 7; Supplementary Figure 4, 5, 6, 7; Supplementary Table 5, 6, 7, 8,**

342 **9**). However, the contig harboring SNP, unaln_JG10~4245875 (2.19 kb) was found to be completely

343 missing from the Nipponbare (*japonica*) reference genome but present on the chromosome 6 in

344 genomes of multiple other rice accessions including *indica* (IR 64 and R 498), *aus* (Nagina 22) and

345 aromatic (Sonasal, Basmati 334, and Dom Sufid). The physical position of this SNP locus (QTLs

346 thereof) was further determined with reference to the Nipponbare genome using an integrated genomic

347 approach as described earlier. Similar to SNP-unaln_JG10~6950498, 5 kb flanking sequences from the

348 location of SNP unaln_JG10~4245875 in all the aforesaid genomes except Nipponbare (as the

349 sequence is missing from the Nipponbare genome) were retrieved and compared using multiple

350 sequence alignment. The contig harboring SNP-unaln_JG10~4245875 (2.19 kb) was found to overlap

351 with the GRAS family transcription factor gene (*OsR498G0612839600.01*) in R 498 (*indica*) reference

352 genome. This gene predominantly expresses in leaf, panicle and endosperm of R 498 <u>*indica*</u> rice

353 accession and encodes for 498 amino acid protein. Subsequent, comparison of an

354 *OsR498G0612839600.01* genomic sequence among five rice accessions, revealed 1669 bp insertion in

355 all three aromatic accessions (Sonasal, Basmati 334 and Dom Sufid) compared to both *indica* (IR 64

356 and R 498) and an *aus* (Nagina 22), leading to truncated protein with 108 amino acids. Interestingly,

11

357    the GRAS family proteins are previously reported to play a vital role in plant development, including

358    male gametogenesis predominantly by regulating cell division, proliferation and maintenance. One of

359    such GRAS family rice gene, *GRAIN SIZE 6* (*GS6*) */OsGRAS-32/ DLT* is reported to regulate grain

360    width and thousand-grain weight by promoting cell expansion in spikelet hull in rice through

361    modulating both brassinosteroid and gibberellin signaling (Sun et al., 2013). However, a homology

362    study revealed that protein encoded by *OsR498G0612839600.01* is homolog to scarecrow-like protein

363    3 (*SCL3*) which belongs scarecrow subgroup of the GRAS family. In *Arabidopsis*, SCL3 is known to

364    function downstream of DELA transcription factor (another subgroup of GRAS family) and are known

365    to negatively regulate of gibberellin signaling. *SCL3* plays an important a role in *Arabidopsis* root

366    elongation by regulating cell division in roots (Lee et al., 2016). However, the effect of *SCL3* on seed

367    development is not yet studied in plants. Thus, it would be interesting to study the probable role of

368    *SCL3* in regulating grain size/weight of rice.

369    Besides, diverse sub-population specific (dispensable) novel candidate genes underlying the

370    grain size/weight loci mapped on pseudo-chromosomes were identified. These include genes encoding

371    multi-domain protein (*LOC_Os10g31770*) for grain width and grain length-to-width ratio as well as

372    dirigent family (*LOC_Os12g12600*) and unknown expressed (*LOC_Os12g12610*) genes for grain

373    weight in rice (**Details in Supplementary Results**).

**Imputation of RPGA-SNP genotyping data using 3K Rice Reference Panel (RICE-RP)**

375    Imputation of high-quality reference SNP genotyping dataset dramatically increases marker density

376    and has been shown to positively impact GWAS results. As the RPGA is designed to efficiently tag

377    almost all haplotype variation existing in 3K rice panel, the imputation of RPGA-based genotyping

378    data of diversity panel was performed using 3K RICE-RP (Wang et al., 2018a). This increased 48257

379    to 5231751 SNPs from 12 Nipponbare chromosomes. These SNPs tags > 3500 cloned and functionally

380    characterized genes in rice. Further, 2760746 of these SNPs displayed minor allele frequency (MAF)

381    >0.05, which translates average marker density of ~ 7.4 SNPs/kb compared to pre-imputation average

382    marker density of ~ 0.12 SNPs/kb.This result suggests that the SNP genotyping data generated using

383    RPGA can be efficiently imputed to obtain a manifold increase in marker density which can further

384    result in higher resolution and enhanced power of QTL detection in GWAS.

385    Further, GWAS analysis was repeated for grain length and grain width traits using imputed

386    RPGA genotyping data with a total of ~2.76 M SNPs to assess the effect of increased marker density

387    on power and resolution of trait association in rice. The post-imputation GWAS for grain length trait

388    detected *Grain Size 3* (*GS3*) locus; however, with a much smaller p-value. Another major locus

389    detected on chromosome 12 also crossed a stringent significance threshold which previously failed to

12

390  do so (**Supplementary Figure 8**). This indicates the increased power of QTL detection in GWAS

391  conducted using the imputed RPGA-based SNP genotyping data. Interestingly, post-imputation

392  GWAS for grain width trait detected a novel locus on chromosome 8 which remained undetected

393  before imputation (**Supplementary Figure 8**). This highlights the potential of imputed RPGA-based

394  SNP genotyping data to detect novel associations which otherwise would be missed due to low

395  marker-density.

**RPGA-based high-resolution QTL mapping to dissect the genetic basis of grain size/weight**

**variation in aromatic rice**

398  To conduct QTL mapping, 190 mapping individuals and parental accessions selected from a

399  developed $F_{12}$ RIL population (Sonasal × PB 1121) RIL () were phenotyped for grain length, grain

400  width, length-to-width ratio and thousand grain weight. The grain length of RILs ranged from 3.5 to

401  8.1 mm with a mean ± standard deviation (SD) of 5.4 ± 1.4 mm. The grain width of RILs varied from

402  1.5 to 2.5 mm with a mean ± SD of 1.9 ± 0.17 mm. Finally, the thousand-grain weight for RILs ranged

403  from 8.8 to 31.3 g with a mean ± SD of 17.2 ± 4.4 g. Further, grain length positively correlated with

404  thousand-grain weight (Pearson's correlation coefficient, r = 0.74) whereas grain width was found to

405  be negatively correlated with both grain length (r = -0.29) and thousand-grain weight (r = -0.032)

406  (**Supplementary Figure 9, 10**). This suggests thousand-grain weight is predominantly influenced by

407  grain length, whereas, grain width has either little or negligible influence on the thousand-grain weight

408  in the said mapping population.

409  The previously generated ultra-high-density genetic linkage map (Sonasal × PB 1121 RILs)

410  and high-quality SNP genotyping data (13793 SNPs) were integrated with grain size/weight (grain

411  length, grain width, length-to-width ratio and thousand-grain weight) phenotype data of 190 RILs and

412  parents , to identify potential QTLs and their underlying candidate genes governing said target traits in

413  rice. A total of six QTLs which are genetically mapped on either chromosome 3 or 7 governing four

414  grain size/weight traits were identified **(Figure 8; Supplementary Table 10)**. The phenotypic

415  variation explained (PVE) by these grain size/weight QTLs varied from 13.3 to 42 % **(Supplementary**

416  **Table 10**). Five of these QTLs showed correspondence with previously known grain size/weight

417  genes. Among these, major grain length QTL, *qGL3* and thousand-grain weight QTL, *qTGW3*,

418  coincides with *GS3* locus (Fan et al., 2006). The grain length QTL, *qGL7* found to harbor previously

419  known grain length gene *GRAIN LENGTH 7* (*GL7*) (Wang et al., 2015; Zhou et al., 2015). Similarly,

420  grain width QTL, *qGW7* contains Os*BZR1* which is a known regulator of grain size in rice (Tong et

421  al., 2009). The QTL detected for length-to-width ratio *qLWR7* harbors *GL3.2,* a previously known

422  grain length regulating gene encoding cytochrome P450  **(Figure 8; Supplementary Table 10)**.

423  Interestingly, the only novel QTL identified on chromosome 7 (*qLWR7*) was found to coincide
424  with dispensable locus associated with grain length as well as grain length-to-width ratio as identified
425  in RPGA-based GWAS (**Figure 7; Figure 8c**). Thus revalidating the authenticity of previously
426  detected dispensable loci using RPGA-based GWAS. As previously discussed, the locus contains
427  WDR12 protein encoding gene (*LOC_Os07g40930*) can be a candidate gene modulating grain
428  size/weight in rice, as WDR12 protein is a vital component of the PeBoW (PES-BOP1-
429  WDR12) complex, which is known to regulate cell proliferation in *Arabidopsis* (Cho et al., 2013; Ahn
430  et al., 2016).

431  **The utility of RPGA in hybridity testing and genetic background recovery analysis**

432  To demonstrate the utility of RPGA for hybridity testing, multiple $F_1$ hybrids generated by cross-
433  hybridization of PB 1121 with different *indica* rice accessions (Vandana, Phule Radha, Pusa 1927,
434  C101A51, and Kalibagh) were genotyped using RPGA. The genome-wide SNP data consisting of
435  42565 SNP markers were then converted into the Graphical Genotype to visually detect the hybrid
436  lines. Comparison of rice accessions using the parental polymorphic SNP markers across 12 rice
437  chromosomes showed heterozygous nature of analyzed $F_1$ lines revealing the usefulness of RPGA in
438  hybridity testing of rice.

439  Similarly, to demonstrate utility of RPGA for background selection commonly required for
440  marker-assisted breeding and crop improvement of rice, two NILs with Sonasal (recurrent parent) and
441  LGR (donor parent) were genotyped using RPGA-SNP. The RPGA SNP genotyping data of two NILs
442  were then compared with recurrent and donor parental accessions. The graphical genotype for
443  genotyping data of 42565 SNPs in the NILs was made to assess the genomic contribution from both
444  parental accessions. This graphical representation ascertained the advantages of RPGA in higher
445  genetic background recovery analysis, further highlighting the importance of RPGA for marker-
446  assisted crop improvement in rice (**Supplementary Figure 11, 12**).

447  **Rice Genome Genotyping Array Analysis Portal (RAP) provides a user-friendly interface for**
448  **analyzing RPGA SNP genotyping data**

449  A web-based application RAP was developed to provide the user with the ability to conduct pan-
450  genome-based GWAS using RPGA-based SNP genotyped data without the need for programming
451  skills. Using RAP, pan-genome-based GWAS can be performed either with existing RPGA
452  genotyping data of diverse rice accessions or user can upload their own preferred RPGA SNP
453  genotyping data. The RAP further enables the user to identify the physical locations of pseudo-
454  chromosome-specific SNPs concerning different available rice reference genomes, which is crucial for
455  the identification of dispensable genes regulating the trait of interest. Further, RAP also enables users

14

456 to impute RPGA SNP genotyping data using a 3K rice reference panel. The GWAS performed using
457 the imputed SNP genotyping data is especially useful for the identification of novel causal mutations
458 underlying the trait-associated loci in rice. In addition to GWAS, the RAP also enables the user to
459 retrieve SNP genotyping information based on locus IDs of the rice genes (Rice Genome Annotation
460 Project, http://rice.uga.edu/) as well as using the genomic coordinates of 12 chromosomes. Thus, RAP
461 complements the utility of RPGA, and these two altogether provide end-to-end solutions from
462 genotyping to downstream data analysis for their eventual deployment in genomics-assisted crop
463 improvement in rice (**Figure 9**).

464 **DISCUSSION**

465 The recent emergence of pan-genome studies in rice has made evident the presence of extensive
466 dispensable genome diversity in rice (Sun et al., 2017; Wang et al., 2018b; Zhao et al., 2018).
467 Dispensable genes are now known to be the major contributors to adaptive evolution, domestication
468 and overall diversity of agronomically important traits in all major crop plants including rice (Golicz et
469 al., 2016; Tao et al., 2019; Tranchant-Dubreuil et al., 2019). Despite the tremendous potential,
470 leveraging the dispensable gene variation for crop improvement remains a challenging task, especially
471 due to constraints involved in the efficient genotyping of dispensable gene variants in a large number
472 of accessions. Keeping this in mind, we designed "Rice Pan-genome Genotyping Array" (RPGA) to
473 perform rapid, user-friendly and cost-efficient pan-genome-based genotyping in rice. The RPGA is
474 based on a 3K rice pan-genome and, therefore, can efficiently capture variation from both core and
475 dispensable genes, across diverse rice sub-populations (Sun et al., 2017). We demonstrated the utility
476 of RPGA for highly accurate large-scale genotyping of experimental mapping population and diverse
477 natural rice accessions of a diversity panel. The majority of RPGA markers displayed high MAFs
478 suggesting a highly informative nature of RPGA compared to other existing SNP genotyping arrays
479 (Singh et al., 2015; McCouch et al., 2016; Thomson et al., 2017). This suggests suitability of RPGA
480 for diverse high-throughput genotyping applications in rice. Further, a framework map for 3K rice
481 pan-genome generated by integrating RPGA-based ultra-high-density genetic linkage map, a multi-
482 genome alignment of the un-anchored pan-genome contigs and pairwise-LD information between
483 RPGA SNP markers. The importance of such framework map for leveraging crop pan-genomes for
484 multiple genomic-assisted breeding applications have been highlighted previously in maize (Lu et al.,
485 2015). Therefore, the framework-map is expected to be an important step forward toward leveraging
486 3K rice pan-genome for genomics assisted breeding applications and crop improvement in rice.

487 Indian rice genepool possesses a large number of highly diverse rice accessions representing
488 different ecogeographical regions across the globe with wide phenotypic variations for multiple

economically important agronomic traits. However, comprehensive effort to decipher the natural allelic diversity and population genetic structure among rice accessions was still not completely understood (Kumbhar et al., 2015; Roy et al., 2015; Singh et al., 2016). In this context, RPGA-based SNP genotyping was used to assess the natural allelic diversity and population genetic structure among 271 diverse rice accessions. PCA grouped these accessions into five distinct clusters representing *indica, japonica, aus* and aromatic/Basmati subpopulations which is in accordance with previous studies conducted with global rice accessions (McCouch et al., 2016; Wang et al., 2018b). Further, as expected most of the rice accessions represented from India belonged to either *indica* or *aus* subpopulations, as expected from their morphological characteristics and geographical locations. Interestingly, Indian traditional Basmati accessions were found to cluster distinctly from aromatic rice accessions belonging to both north-eastern India and other parts of the world. Thus, it will be interesting to further identify the precise genes/alleles differentiating Indian Basmati accessions from other aromatic rice accessions. Further, evolved Basmati accessions were found to cluster closer to the *indica* accessions in contrast to traditional Basmati which displayed a closer genetic relationship with *japonica* and *aus* accessions. This can be explained by the breeding history of evolved Basmati varieties, which were developed by cross-hybridization between traditional Basmati varieties with superior grain quality and *indica* varieties possessing favorable traits like dwarf height, early flowering and high yield (Singh et al., 2018). Further, assessment of fine population genetic structure existing within rice accessions revealed the existence of two further sub-groups within *indica* subpopulation i.e. *INDI* and *INDII* corresponding to *Xian*/*Indica-2* (*XI-2*) and *XI-3* from South Asia and Southeast Asia, respectively which are previously reported along with two other *indica* subpopulation groups (*XI-1A* from East Asia, *XI-1B* of modern varieties of diverse origins) (Wang et al., 2018b). This suggests the indigenous nature of the most widely cultivated *indica* rice accessions from India. The evolved Basmati accessions were found to be closely related to the *IND1* subpopulation, confirming their origin from cross-hybridization between *IND1* and traditional Basmati accessions. This suggests the utility of RPGA-based SNP genotyping for efficiently decoding the natural allelic diversity and population genetic structure in order to understand the domestication pattern in rice genepool.

The sole reliance on reference genomes creates a reference bias, which negatively impacts many genomic applications including GWAS (Paten et al., 2017; Coletta et al., 2021). The RPGA-based GWAS conducted in this study detected associations (for diverse grain size/weight traits) from both Nipponbare reference genome and sub-population specific pseudo-chromosomes present in 3K rice pan-genome. The later represents loci partially or completely absent from the Nipponbare reference genome. The previously generated pan-genome framework map allowed us to locate the physical locations of these pseudo-chromosome loci relative to the Nipponbare reference genome and

16

523 subsequent identification of underlying candidate genes. This demonstrates the potential of RPGA-
524 based GWAS to overcome the reference bias imposed by traditional genotyping based on single
525 reference genome. The RPGA-based GWAS detected many previously known major grain size/weight
526 genes like *GS3* and *PGL1* (grain length and length-to-width ratio) and *GW5* (grain width, length-to-
527 width ratio, and thousand-grain weight) etc., (Fan et al., 2006; Weng et al., 2008; Heang & Sassa,
528 2012a, 2012b). This not only highlights the major role of these previously known genes in regulating
529 grain size/weight in rice accessions but also validates the ability of pan-genome-based GWAS to
530 detect true associations. The GWAS identified many novel candidate genes including those that are
531 partially or completely absent form Nipponbare genome. These candidate genes similar to previous
532 studies found to be predominantly involved in few major pathways such as ubiquitin-proteasome
533 pathway, brassinosteroid signaling as well as in transcriptional regulation (Li et al. 2018).
534 Subsequently, high-resolution QTL mapping conducted using the RPGA-based ultra-high-density
535 genetic linkage map ( Sonasal × PB 1121 RILs) once again identified *GS3* as major locus regulating
536 grain length and thousand grain weight in aromatic rice. Interestingly, a novel chromosome 7 locus
537 (*qLWR7*) regulating length-to-width ratio detected in QTL mapping was also found to overlap with
538 pseudo-chromosome specific QTL detected in RPGA-based GWAS for grain width. This concordance
539 between GWAS and QTL mapping outcomes validates the authenticity of QTLs (from the Nipponbare
540 genome as well as sub-population specific pseudo-chromosomes) identified using RPGA-based
541 GWAS. Thus, RPGA-based GWAS can be efficiently identified using trait-associated genomic loci
542 which otherwise would have been missed using conventional GWAS relying on a single reference
543 genome. The pan genome-based GWAS and QTL mapping overall delineated two known rice grain
544 size/weight genes, namely *GS3* regulating grain length and thousand grain weight and *GW5* governing
545 thousand-grain weight in rice. Besides, multiple promising sub-population specific (dispensable) novel
546 candidate genes modulating grain size/weight for the loci mapped on psudochromosomes were
547 identified. These include *WDR12* (*LOC_Os07g40930*) for grain length and grain length-to-width ratio,
548 multi-domain protein (*LOC_Os10g31770*) for grain width, GRAS family transcription factor
549 (*OsR498G0612839600*.01) for grain length and dirigent family gene (*LOC_Os12g12600*) for grain
550 weight in rice. Among these, the *WDR12* gene underlying the *qLWR7* QTL, validated by both RPGA-
551 based GWAS and QTL mapping, is known to modulate cell number/cell size during growth and
552 development of crop plants and thus appears to be a promising candidate gene regulating grain length
553 and grain length-to-width ratio in rice.

554 These aforementioned results confirm the utility of RPGA for efficiently leveraging 3K rice
555 pan-genome for diverse genomics-assisted breeding applications including molecular diversity
556 analysis, ultra-high-density genetic linkage map constitution, high-resolution QTL mapping, and

17

557   GWAS in rice. Further, novel grain size genes/QTLs identified with RPGA-based QTL mapping and

558   GWAS will not only expedite the genetic improvement of Indian rice varieties but also unravel many

559   previously missing links in the genetic and molecular regulation of grain size/weight in rice. Lastly,

560   the "Rice Pan-genome Genotyping Array Analysis Portal" (RAP) developed in this study provides the

561   user with easy to use web interface for handling and analyzing RPGA genotype data. This includes

562   imputation using the 3K Rice Reference Panel and performing pan-genome-based GWAS analysis.

563   This makes the RPGA as an ideal SNP genotyping solution for pan-genome based genetic studies

564   including complex trait dissection and diverse applications in Indian trade and commerce including

565   DNA fingerprinting, genetic purity and hybridity testing. Thus, RPGA and RAP together provide an

566   end-to-end solution for leveraging 3K rice pan-genome for accelerated molecular breeding

567   applications including marker-assisted selection and genomic selection thus enabling genetic

568   enhancement of rice (**Figure 10**).

569   **MATERIAL AND METHODS**

570   **Development of a bi-parental mapping population**

571   Pusa Basmati 1121 (PB 1121) and Sonasal, two accessions with contrasting grain size/weight, were

572   utilized for the generation of a bi-parental $F_{12}$ RIL mapping population (190 individuals) varying for

573   grain size and weight traits. PB 1121 is a long-grain Basmati variety with grain length: 9.0 mm, grain

574   width: 2.78 mm and thousand-grain weight: 19.2 g. Sonasal is a short-grain aromatic landrace with

575   grain length: 5.2 mm and grain width: 2.7 mm and thousand-grain weight: 11 g.

576   **Constitution of the diversity panel**

577   A diversity panel comprising of 271 diverse rice accessions belonging to different rice population

578   (*indica,* tropical/temperate *japonica, aus*, long/short-grain aromatic and wild etc.) was constituted.

579   This panel includes 241 landraces and high-yielding mega-varieties  representing different agro-

580   climatic regions of India. In addition, the diversity panel included 24 traditional/ evolved Basmati rice

581   varieties notified by the Government of India and accessions and six accessions of wild rice (three

582   accessions of *Oryza rufipogon* and *Oryza nivara* each). The diversity panel harbors a high level of

583   diversity for a diverse range of traits of agronomic importance including yield (grain size/weight) and

584   a/biotic stress tolerance, etc. Therefore, this diversity panel is ideal for GWAS and genomic selection,

585   especially for grain size/weight trait in rice.

586   **Phenotypic evaluation**

587   The RIL mapping population (190 individuals) along with parental accessions (PB 1121 and Sonasal)

588   and diversity panel (271 accessions) were grown at the field during the *Kharif*  for three consecutive

589 years (2016, 2017 and 2018) in a complete randomized block design with three replications. For the

590 measurement of grain size/dimension traits (grain length and grain width), 50 representative grains

591 were scanned using the Epson Expression 1000XL seed scanner. The mean grain length and mean

592 grain width in millimeters (mm) of each mapping individual/accession were then considered for

593 further analysis. The grain length to grain width ratio of each mapping individual/accession was

594 calculated by dividing mean grain length with mean grain width.

595 The grain weight (g) was then estimated by measuring the mean weight of 1000 mature dried grains

596 (at 10% moisture content) selected from each accession/individual with three replications. The various

597 statistical parameters, including frequency distribution, coefficient of variation (CV) and broad-sense

598 heritability ($H^2$) of grain size/weight traits among accessions/individuals were estimated using

599 SPSSv17.0.

**Designing of a RPGA for SNP genotyping**

601 A rice pan-genome array (80K) was designed to tag haplotype variation from the entire 3K rice pan-

602 genome which includes the complete Nipponbare genome (IRGSP 1.0/MSU release 7: 373 Mbps) and

603 genomic sequence specific to rice pan-genomes representing different sub-groups of 5 rice sub-

604 populations (*indica,* tropical/temperate *japonica, aus* and aromatic, etc.) (Sun et al., 2017)

605 **(Supplementary Figure 13)**. To select SNPs representing the Nipponbare genome, SNPs were

606 retrieved from all the major publicly available rice sequence variation datasets which include: I) 3,000

607 Rice Genomes Project: ~ 40 M SNPs from 3243 resequenced global rice accessions (Wang et al.,

608 2018b), II) RiceVarMap: ~ 6.5 million SNPs from 1479 resequenced Chinese rice accessions (Zhao et

609 al., 2015), and III) High-Density Rice Array (HDRA): ~ 0.7 M SNPs (700K SNPs) from 1953

610 genotyped rice accessions (McCouch et al., 2016). Further, to identify SNPs from remaining rice pan-

611 genome sequences (different rice sub-group specific genomic contigs), the publicly available  short-

612 read sequence data for 400 Indian rice accessions (NCBI-SRA) including PB 1121 and Sonasal

613 (Kumar et al. 2021) were aligned against the rice 3K pan-genome sequence (Sun et al., 2017) using

614 BWA program (Li and Durbin, 2009). The high-quality SNPs (~ 0.5 M) were then identified utilizing

615 standard GATK best practices workflow for variant calling (McKenna et al., 2010; Poplin et al.,

616 2017). From this constituted SNP resource, 80504 SNPs distributed uniformly (50 kb genomic

617 intervals) across 3K rice pan-genomes were finally screened following diverse selection criteria and

618 quality filter parameters for tiling on the RPGA **(Supplementary Method)**.  The RPGA includes 2000

619 SNPs from 164 functionally characterized known genes associated with grain size/weight, grain

620 quality and grain aroma traits in rice.

**RPGA genotyping, genotype data analysis and imputation**

622 The genotyping of both RIL population (190 individuals and two parental accessions) and the diversity
623 panel (271 accessions) was performed using the 80K RPGA on GeneTitan® Multi Channel (MC)
624 instrument (Affymetrix, USA) following manufacturer's instructions (**Supplementary Method**).
625 Subsequent quality control and genotype calling was carried using Axiom Analysis Suite 2.0
626 (Affymetrix, USA) following Axiom best practices workflow (**Supplementary Method**). The
627 genotype data of the diversity panel was imputed based on rice reference panel (RICE-RP) generated
628 by combining 3K Rice Genome dataset (18M SNPs from 3024 accessions) and HDRA dataset (700K
629 SNPs from 1568 accessions) using IMPUTE2 (Howie et al., 2009) as per Wang et al. (2018a).

630 **High-density genetic linkage map construction**

631 To construct a high-density genetic linkage map, 80K RPGA-derived genotype data of SNPs showing
632 polymorphism between two parental accessions as well as 190 mapping individuals of a RIL
633 population (Sonasal × PB 1121) were analyzed in the R/qtl program (Broman et al., 2003). The SNP
634 genotyping data of RILs were used to construct a preliminary genetic map with the ordered groupings
635 of marker-pairs in the individual linkage groups at a minimum logarithm of odds (LOD) of 8 and
636 maximum pair-wise recombination fraction (RF) of 0.40 (**Supplementary Method**). After removing
637 the erroneous SNP genotyping data of RILs and subsequent optimal reordering of markers in the
638 linkage groups, a RPGA-based high-density genetic map was constructed finally.

639 *De novo* **genome assembly and annotation of a draft genome of Sonasal**

640 For long-read Nanopore and short-read Illumina sequencing, high-quality genomic DNA isolated from
641 the young leaf samples of a short-grain aromatic rice variety Sonasal was used for constructing
642 genomic DNA libraries using the SQK-LSK109 ligation kit (company name, country) and Illumina
643 TruSeq DNA sample prep kit (Illumina, USA), respectively following the manufacturer's protocol.
644 The libraries prepared were sequenced using the FLO-MIN106 (R9.4) flowcells on PromethION 24
645 (P24) Nanopore platform and following the paired-end sequencing workflow of Illumina HiSeq2000
646 platform (Illumina, USA) (**Supplementary Method**).

647 The high-quality Nanopore long-reads were corrected using Canu (Koren et al., 2014) and
648 further assembled using Flye genome assembler (Kolmogorov et al., 2019), and finally polished with
649 Recon v1.4.11 (https://github.com/isovic/racon) and Medaka
650 (https://github.com/nanoporetech/medaka) for different rounds. BUSCO v4.0.5 (Simão et al., 2015)
651 was used to assess the completeness of draft genome assembly of Sonasal (**Supplementary Method**).
652 The synteny between the assembled Sonasal draft genome and the Nipponbare reference genome
653 (IRGSP 1.0) was analyzed using D-GENIES program (Cabanettes and Klopp, 2018). Subsequently,
654 comprehensive structural annotation of Sonasal draft genomic sequences was performed using

20

655  MAKER annotation pipeline (Holt and Yandell, 2011), RepeatMasker v4.0.7 (Ref), SNAP (Korf,
656  2004) and Augustus (Stanke et al., 2008). The high-density genetic linkage maps developed from
657  Sonasal × PB 1121 RIL mapping population were finally used for the scaffolding of annotated contigs
658  with Chromonomer v1.10 (Catchen et al., 2020) (**Supplementary Method**).

659  **PCA and LD analysis**

660  For PCA, the genotype data of SNPs common to both 271 Indian rice accessions (RPGA dataset) and
661  rice reference panel with 3032 global rice accessions (3000 Rice Genomes Project, 2014) were
662  retrieved and combined. The PCA was performed on combined as well as aforesaid individual SNP
663  genotype dataset with TASSEL v5.0 (Bradbury et al., 2007). The PCA plot was then generated using
664  the ggfortify R package. LD statistics for the RPGA-derived SNP genotype dataset were calculated
665  using PLINK v2.0 (Chang et al., 2015) with a window size of 50 kb. The LD-decay plot was then
666  generated with R version 3.4.1.

667  **Population structure and molecular diversity analysis**

668  To decipher the population structure in 271 diverse Indian accessions, a core SNP dataset  was
669  generated with the LD pruning procedure implemented in PLINK v2.0 (Chang et al., 2015)  following
670  Supplemental Methods. The core SNP dataset (5812 SNPs) was analyzed using the variational
671  Bayesian model of fastSTRUCTURE v1.0 with default parameters and number of clusters (K) ==1 to
672  15. Finally, the replicates of chosen clusters (most-likely K value) once finalized by best chose
673  function of fastSTRUCTURE were summarized with the CLUMPK (http://clumpak.tau.ac.il/).

674  **GWAS of grain size/weight**

675  GWAS for grain size/weight traits was performed using five different models (CMLM, MLMM,
676  SUPER, FarmCPU and BLINK) implemented in GAPIT v.3 (Wang and Zang 2019). Before
677  performing GWAS, SNP genotype data were filtered to ensure less than 10% missing genotyping data
678  both for accessions and markers. Besides, markers with less than 5% minor allele frequency (MAF)
679  were also eliminated. For all the five methods, the kinship matrix internally calculated in GAPIT was
680  utilized to account for relatedness between accessions. Further, the top two principal components were
681  used as covariates in statistical models to account for the population structure existing within the
682  diversity panel. Additionally, to balance the trade-off between stringency and false negatives two
683  different significance thresholds [a stringent Benjamini-Hochberg threshold ($P \leq 1.5E^{-7}$) and a less
684  stringent suggestive threshold ($P \leq 5E^{-5}$)] were utilized. A grain size/weight QTL detected in more
685  than one method with GWAS was considered true positive.

686  **Identification of candidate genes underlying trait associated genomic loci**

687 For identification of candidate genes underlying grain size/weight associated loci from the Nipponbare
688 reference genome, haplotype blocks associated with the most significant SNPs were identified with
689 PLINK v2.0 (Chang et al., 2015) and all genes within the haplotype blocks considered as potential
690 candidate genes. Further, candidate genes were prioritized based on the presence of major effect
691 mutations within genes as well as the putative function of these genes.

692 However, candidate genes underlying associated loci from subpopulation-specific pseudo-
693 chromosomes (MSPs either absent or have unknown physical positions in Nipponbare reference
694 genome) were determined using an integrated approach based on multiple rice reference genomes as
695 per the Supplemental Methods and considered further as novel MSPs associated with grain size/weight
696 in rice.

**QTL mapping**

698 QTL mapping was performed by integrating the high-quality SNP genotyping data across 144 RIL
699 mapping individuals (Sonasal × PB 1121) and high-density genetic linkage map with grain size/weight
700 trait phenotype data of respective RILs utilizing r/qtl package (Broman et al., 2003). Genotype
701 probabilities for these markers were then calculated with calc.genoprob() function utilizing the.
702 Further, interval mapping with Kosambi mapping function was conducted utilizing the scanone()
703 function with Haley-Knott Regression, assuming normal phenotypic distribution. The LOD
704 significance threshold was further determined with 1000 permutations for each trait independently.
705 Further, the QTL model refinement was performed for the detection of additional linked QTLs. The
706 QTL model was then finalized by stepwise forward selection and backward elimination to identify
707 best fit the QTL model for given data. The phenotypic variation explained (PVE) and estimated effect
708 sizes of each QTL was determined by performing *ANOVA* on the final selected model. The 95%
709 confidence interval for each QTL peak was then determined as 1.5 LOD units flanking the peak
710 marker (Mangin et al., 1994; Dupuis and Siegmund, 1999).

**Development of RPGA Analysis Portal (RAP)**

712 RAP (http://www.rpgaweb.com) was developed based on Apache HTTP server (version 2.4.6)
713 integrated with PHP (version 7.3.3) and MySQL (version 8.0.15), on a server machine with Centos 8
714 operating system. HTML and CSS were used to create a responsive front-end interface whereas PHP
715 JavaScript and R were used to develop the backend of RAP. The imputation pipeline previously
716 described by Wang et al. (2018a) was further integrated to conduct efficient imputation of RPGA
717 genotype data based on rice reference panel (RICE-RP) generated by combining 3K Rice Genome
718 dataset (18M SNPs from 3024 accessions) and High-Density Rice Array dataset (700K SNPs from

22

719    1568 accessions). In addition, the GAPIT pipeline (Lipka et al., 2012) and JBrowse (Buels et al.,

720    2016) was integrated to conduct GWAS analysis and visualize variants across the rice genome.

## ACKNOWLEDGMENTS

## CONFLICT OF INTERESTS

726    The authors declare that they have no competing interests.

## REFERENCES

729    3,000 Rice Genomes Project, 2014. The 3,000 rice genomes project. GigaScience, 3(1), pp.2047-
730      217X.

732    Abe, Y., Mieda, K., Ando, T., Kono, I., Yano, M., Kitano, H. and Iwasaki, Y., 2010. The SMALL
733      AND ROUND SEED1 (SRS1/DEP2) gene is involved in the regulation of seed size in rice.
734      Genes & genetic systems, 85(5), pp.327-339.

736    Ahn, C.S., Cho, H.K., Lee, D.H., Sim, H.J., Kim, S.G. and Pai, H.S., 2016. Functional characterization
737      of the ribosome biogenesis factors PES, BOP1, and WDR12 (PeBoW), and mechanisms of
738      defective cell growth and proliferation caused by PeBoW deficiency in Arabidopsis. Journal of
739      experimental botany, 67(17), pp.5217-5232.

741    Alonge, M., Wang, X., Benoit, M., Soyk, S., Pereira, L., Zhang, L., Suresh, H., Ramakrishnan, S.,
742      Maumus, F., Ciren, D. and Levy, Y., 2020. Major impacts of widespread structural variation on
743      gene expression and crop improvement in tomato. Cell, 182(1), pp.145-161.

745    Ashikawa, I., Hayashi, N., Yamane, H., Kanamori, H., Wu, J., Matsumoto, T., Ono, K. and Yano, M.,
746      2008. Two adjacent nucleotide-binding site–leucine-rich repeat class genes are required to confer
747      Pikm-specific rice blast resistance. Genetics, 180(4), pp.2267-2276.

749    Bianco, L., Cestaro, A., Linsmith, G., Muranty, H., Denancé, C., Théron, A., Poncet, C., Micheletti,
750      D., Kerschbamer, E., Di Pierro, E.A. and Larger, S., 2016. Development and validation of the
751      Axiom® Apple480K SNP genotyping array. The Plant Journal, 86(1), pp.62-74.

753    Bradbury, P.J., Zhang, Z., Kroon, D.E., Casstevens, T.M., Ramdoss, Y. and Buckler, E.S., 2007.
754      TASSEL: software for association mapping of complex traits in diverse samples. Bioinformatics,
755      23(19), pp.2633-2635.

757    Bradbury, P.J., Zhang, Z., Kroon, D.E., Casstevens, T.M., Ramdoss, Y. and Buckler, E.S., 2007.

758      TASSEL: software for association mapping of complex traits in diverse samples. Bioinformatics,
759      23(19), pp.2633-2635.

761 Broman, K.W., Wu, H., Sen, Ś. and Churchill, G.A., 2003. R/qtl: QTL mapping in experimental
762      crosses. bioinformatics, 19(7), pp.889-890.

764 Buels, R., Yao, E., Diesh, C.M., Hayes, R.D., Munoz-Torres, M., Helt, G., Goodstein, D.M., Elsik,
765      C.G., Lewis, S.E., Stein, L. and Holmes, I.H., 2016. JBrowse: a dynamic web platform for
766      genome visualization and analysis. Genome biology, 17(1), pp.1-12.

768 Cabanettes, F. and Klopp, C., 2018. D-GENIES: dot plot large genomes in an interactive, efficient and
769      simple way. PeerJ, 6, p.e4958.

771 Catchen, J., Amores, A. and Bassham, S., 2020. Chromonomer: a tool set for repairing and enhancing
772      assembled genomes through integration of genetic maps and conserved synteny. G3: Genes,
773      Genomes, Genetics, 10(11), pp.4115-4128.

775 Catchen, J., Amores, A. and Bassham, S., 2020. Chromonomer: a tool set for repairing and enhancing
776      assembled genomes through integration of genetic maps and conserved synteny. G3: Genes,
777      Genomes, Genetics, 10(11), pp.4115-4128.

779 Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M. and Lee, J.J., 2015. Second-
780      generation PLINK: rising to the challenge of larger and richer datasets. Gigascience, 4(1),
781      pp.s13742-015.

783 Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M. and Lee, J.J., 2015. Second-
784      generation PLINK: rising to the challenge of larger and richer datasets. Gigascience, 4(1),
785      pp.s13742-015.

787 Chen, H., Xie, W., He, H., Yu, H., Chen, W., Li, J., Yu, R., Yao, Y., Zhang, W., He, Y. and Tang, X.,
788      2014. A high-density SNP genotyping array for rice biology and molecular breeding. Molecular
789      plant, 7(3), pp.541-553.

791 Cho, H.K., Ahn, C.S., Lee, H.S., Kim, J.K. and Pai, H.S., 2013. Pescadillo plays an essential role in
792      plant cell growth and survival by modulating ribosome biogenesis. The Plant Journal, 76(3),
793      pp.393-405.

795 Choi, J.Y., Lye, Z.N., Groen, S.C., Dai, X., Rughani, P., Zaaijer, S., Harrington, E.D., Juul, S. and
796      Purugganan, M.D., 2020. Nanopore sequencing-based genome assembly and evolutionary
797      genomics of circum-basmati rice. Genome biology, 21(1), pp.1-27.

799 Della Coletta, R., Qiu, Y., Ou, S., Hufford, M.B. and Hirsch, C.N., 2021. How the pan-genome is
800      changing crop genomics and improvement. Genome biology, 22(1), pp.1-19.

802  Didion, J.P., Yang, H., Sheppard, K., Fu, C.P., McMillan, L., de Villena, F.P.M. and Churchill, G.A.,
803      2012. Discovery of novel variants in genotyping arrays improves genotype retention and reduces
804      ascertainment bias. BMC genomics, 13(1), pp.1-18.
805

806  Dupuis, J. and Siegmund, D., 1999. Statistical methods for mapping quantitative trait loci from a dense
807      set of markers. Genetics, 151(1), pp.373-386.
808

809  Fan, C., Xing, Y., Mao, H., Lu, T., Han, B., Xu, C., Li, X. and Zhang, Q., 2006. GS3, a major QTL for
810      grain length and weight and minor QTL for grain width and thickness in rice, encodes a putative
811      transmembrane protein. Theoretical and applied genetics, 112(6), pp.1164-1171.
812

813  Fan, C., Yu, S., Wang, C. and Xing, Y., 2009. A causal C–A mutation in the second exon of GS3
814      highly associated with rice grain length and validated as a functional marker. Theoretical and
815      Applied Genetics, 118(3), pp.465-472.
816  Fukuoka, S., Saka, N., Koga, H., Ono, K., Shimizu, T., Ebana, K., Hayashi, N., Takahashi, A.,
817      Hirochika, H., Okuno, K. and Yano, M., 2009. Loss of function of a proline-containing protein
818      confers durable disease resistance in rice. Science, 325(5943), pp.998-1001.
819

820  Gamuyao, R., Chin, J.H., Pariasca-Tanaka, J., Pesaresi, P., Catausan, S., Dalid, C., Slamet-Loedin, I.,
821      Tecson-Mendoza, E.M., Wissuwa, M. and Heuer, S., 2012. The protein kinase Pstol1 from
822      traditional rice confers tolerance of phosphorus deficiency. Nature, 488(7412), pp.535-539.
823

824  Gao, L., Gonda, I., Sun, H., Ma, Q., Bao, K., Tieman, D.M., Burzynski-Chang, E.A., Fish, T.L.,
825      Stromberg, K.A., Sacks, G.L. and Thannhauser, T.W., 2019. The tomato pan-genome uncovers
826      new genes and a rare allele regulating fruit flavor. Nature genetics, 51(6), pp.1044-1051.
827

828  Gao, X., Zhang, J.Q., Zhang, X., Zhou, J., Jiang, Z., Huang, P., Tang, Z., Bao, Y., Cheng, J., Tang, H.
829      and Zhang, W., 2019. Rice qGL3/OsPPKL1 functions with the GSK3/SHAGGY-like kinase
830      OsGSK3 to modulate brassinosteroid signaling. The Plant Cell, 31(5), pp.1077-1093.
831

832  Golicz, A.A., Bayer, P.E., Barker, G.C., Edger, P.P., Kim, H., Martinez, P.A., Chan, C.K.K., Severn-
833      Ellis, A., McCombie, W.R., Parkin, I.A. and Paterson, A.H., 2016. The pangenome of an
834      agronomically important crop plant Brassica oleracea. Nature communications, 7(1), pp.1-8.
835

836  Hattori, Y., Nagai, K., Furukawa, S., Song, X.J., Kawano, R., Sakakibara, H., Wu, J., Matsumoto, T.,
837      Yoshimura, A., Kitano, H. and Matsuoka, M., 2009. The ethylene response factors SNORKEL1
838      and SNORKEL2 allow rice to adapt to deep water. Nature, 460(7258), pp.1026-1030.
839

840  Heang, D. and Sassa, H., 2012a. An atypical bHLH protein encoded by POSITIVE REGULATOR OF
841      GRAIN LENGTH 2 is involved in controlling grain length and weight of rice through interaction
842      with a typical bHLH protein APG. Breeding science, 62(2), pp.133-141.
843

844  Heang, D. and Sassa, H., 2012b. Antagonistic actions of HLH/bHLH proteins are involved in grain
845      length and weight in rice. PloS one, 7(2), p.e31325.

846

847  Hirsch, C.N., Foerster, J.M., Johnson, J.M., Sekhon, R.S., Muttoni, G., Vaillancourt, B., Peñagaricano,
848      F., Lindquist, E., Pedraza, M.A., Barry, K. and de Leon, N., 2014. Insights into the maize pan-
849      genome and pan-transcriptome. The Plant Cell, 26(1), pp.121-135.

850

851  Holt, C. and Yandell, M., 2011. MAKER2: an annotation pipeline and genome-database management
852      tool for second-generation genome projects. BMC bioinformatics, 12(1), pp.1-14.

853

854  Howie, B.N., Donnelly, P. and Marchini, J., 2009. A flexible and accurate genotype imputation
855      method for the next generation of genome-wide association studies. PLoS genetics, 5(6),
856      p.e1000529.

857

858  Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H. and Phillippy, A.M., 2017. Canu:
859      scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation.
860      Genome research, 27(5), pp.722-736.

861  Korf, I., 2004. Gene finding in novel genomes. BMC bioinformatics, 5(1), pp.1-9.

862

863  Kumar, A., Daware, A., Kumar, A., Kumar, V., Gopala Krishnan, S., Mondal, S., Patra, B.C., Singh,
864      A.K., Tyagi, A.K., Parida, S.K. and Thakur, J.K., 2020. Genome‐wide analysis of polymorphisms
865      identified domestication‐associated long low‐diversity region carrying important rice grain
866      size/weight quantitative trait loci. The Plant Journal, 103(4), pp.1525-1547.

867

868  Kumbhar, S.D., Kulwal, P.L., Patil, J.V., Sarawate, C.D., Gaikwad, A.P. and Jadhav, A.S., 2015.
869      Genetic diversity and population structure in landraces and improved rice varieties from India.
870      Rice science, 22(3), pp.99-107.

871

872  Lee, S.A., Jang, S., Yoon, E.K., Heo, J.O., Chang, K.S., Choi, J.W., Dhar, S., Kim, G., Choe, J.E.,
873      Heo, J.B. and Kwon, C., 2016. Interplay between ABA and GA modulates the timing of
874      asymmetric cell divisions in the Arabidopsis root ground tissue. Molecular plant, 9(6), pp.870-
875      884.

876

877  Li, N., Xu, R., Duan, P. and Li, Y., 2018. Control of grain size in rice. Plant Reproduction, 31(3),
878      pp.237-251.

879

880  Li, R., Li, Y., Zheng, H., Luo, R., Zhu, H., Li, Q., Qian, W., Ren, Y., Tian, G., Li, J. and Zhou, G.,
881      2010. Building the sequence map of the human pan-genome. Nature biotechnology, 28(1), pp.57-
882      63.

883

884  Li, X., Singh, J., Qin, M., Li, S., Zhang, X., Zhang, M., Khan, A., Zhang, S. and Wu, J., 2019.
885      Development of an integrated 200K SNP genotyping array and application for genetic mapping,
886      genome assembly improvement and genome wide association studies in pear (Pyrus). Plant
887      biotechnology journal, 17(8), pp.1582-1594.

888

889  Li, Y.H., Zhou, G., Ma, J., Jiang, W., Jin, L.G., Zhang, Z., Guo, Y., Zhang, J., Sui, Y., Zheng, L. and

Zhang, S.S., 2014. De novo assembly of soybean wild relatives for pan-genome analysis of diversity and agronomic traits. Nature biotechnology, 32(10), pp.1045-1052.

Lipka, A.E., Tian, F., Wang, Q., Peiffer, J., Li, M., Bradbury, P.J., Gore, M.A., Buckler, E.S. and Zhang, Z., 2012. GAPIT: genome association and prediction integrated tool. Bioinformatics, 28(18), pp.2397-2399.

Liu, J., Chen, J., Zheng, X., Wu, F., Lin, Q., Heng, Y., Tian, P., Cheng, Z., Yu, X., Zhou, K. and Zhang, X., 2017. GW5 acts in the brassinosteroid signalling pathway to regulate grain width and weight in rice. Nature Plants, 3(5), pp.1-7.

Lu, F., Romay, M.C., Glaubitz, J.C., Bradbury, P.J., Elshire, R.J., Wang, T., Li, Y., Li, Y., Semagn, K., Zhang, X. and Hernandez, A.G., 2015. High-resolution genetic mapping of maize pan-genome sequence anchors. Nature Communications, 6(1), pp.1-8.

Lu, L., Shao, D., Qiu, X., Sun, L., Yan, W., Zhou, X., Yang, L., He, Y., Yu, S. and Xing, Y., 2013. Natural variation and artificial selection in four genes determine grain shape in rice. New Phytologist, 200(4), pp.1269-1280.

Mabire, C., Duarte, J., Darracq, A., Pirani, A., Rimbert, H., Madur, D., Combes, V., Vitte, C., Praud, S., Rivière, N. and Joets, J., 2019. High throughput genotyping of structural variations in a complex plant genome using an original Affymetrix® axiom® array. BMC genomics, 20(1), pp.1-25.

Mangin, B., Goffinet, B. and Rebai, A., 1994. Constructing confidence intervals for QTL location. Genetics, 138(4), pp.1301-1308.

Maron, L.G., Guimarães, C.T., Kirst, M., Albert, P.S., Birchler, J.A., Bradbury, P.J., Buckler, E.S., Coluccio, A.E., Danilova, T.V., Kudrna, D. and Magalhaes, J.V., 2013. Aluminum tolerance in maize is associated with higher MATE1 gene copy number. Proceedings of the National Academy of Sciences, 110(13), pp.5241-5246.

McCouch, S.R., Wright, M.H., Tung, C.W., Maron, L.G., McNally, K.L., Fitzgerald, M., Singh, N., DeClerck, G., Agosto-Perez, F., Korniliev, P. and Greenberg, A.J., 2016. Open access resources for genome-wide association mapping in rice. Nature communications, 7(1), pp.1-14.

McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M. and DePristo, M.A., 2010. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. Genome research, 20(9), pp.1297-1303.

Montenegro, J.D., Golicz, A.A., Bayer, P.E., Hurgobin, B., Lee, H., Chan, C.K.K., Visendi, P., Lai, K., Doležel, J., Batley, J. and Edwards, D., 2017. The pangenome of hexaploid bread wheat. The Plant Journal, 90(5), pp.1007-1013.

934

935 Paten, B., Novak, A.M., Eizenga, J.M. and Garrison, E., 2017. Genome graphs and the evolution of
936      genome inference. Genome research, 27(5), pp.665-676.

937

938 Poplin, R., Ruano-Rubio, V., DePristo, M.A., Fennell, T.J., Carneiro, M.O., Van der Auwera, G.A.,
939      Kling, D.E., Gauthier, L.D., Levy-Moonshine, A., Roazen, D. and Shakir, K., 2017.
940 Scaling accurate genetic variant discovery to tens of thousands of samples. BioRxiv, p.201178.

941

942 Read, B.A., Kegel, J., Klute, M.J., Kuo, A., Lefebvre, S.C., Maumus, F., Mayer, C., Miller, J., Monier,
943      A., Salamov, A. and Young, J., 2013. Pan genome of the phytoplankton Emiliania underpins its
944      global distribution. Nature, 499(7457), pp.209-213.

945

946 Schmidt, M., Van Bel, M., Woloszynska, M., Slabbinck, B., Martens, C., De Block, M., Coppens, F.
947      and Van Lijsebettens, M., 2017. Plant-RRBS, a bisulfite and next-generation sequencing-based
948      methylome profiling method enriching for coverage of cytosine positions. BMC plant biology,
949      17(1), pp.1-14.

950 Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. and Zdobnov, E.M., 2015. BUSCO:
951      assessing genome assembly and annotation completeness with single-copy orthologs.
952      Bioinformatics, 31(19), pp.3210-3212.

953

954 Singh, N., Choudhury, D.R., Tiwari, G., Singh, A.K., Kumar, S., Srinivasan, K., Tyagi, R.K., Sharma,
955      A.D., Singh, N.K. and Singh, R., 2016. Genetic diversity trend in Indian rice varieties: an analysis
956      using SSR markers. BMC genetics, 17(1), pp.1-13.

957

958 Singh, N., Jayaswal, P.K., Panda, K., Mandal, P., Kumar, V., Singh, B., Mishra, S., Singh, Y., Singh,
959      R., Rai, V. and Gupta, A., 2015. Single-copy gene based 50 K SNP chip for genetic studies and
960      molecular breeding in rice. Scientific reports, 5(1), pp.1-9.

961

962 Singh, V., Singh, A.K., Mohapatra, T. and Ellur, R.K., 2018. Pusa Basmati 1121–a rice variety with
963      exceptional kernel elongation and volume expansion after cooking. Rice, 11(1), pp.1-10.

964

965 Stanke, M., Diekhans, M., Baertsch, R. and Haussler, D., 2008. Using native and syntenically mapped
966      cDNA alignments to improve de novo gene finding. Bioinformatics, 24(5), pp.637-644.

967 Stein, J.C., Yu, Y., Copetti, D., Zwickl, D.J., Zhang, L., Zhang, C., Chougule, K., Gao, D.,

968

969 Iwata, A., Goicoechea, J.L. and Wei, S., 2018. Genomes of 13 domesticated and wild rice relatives
970      highlight genetic conservation, turnover and innovation across the genus Oryza. Nature genetics,
971      50(2), pp.285-296.

972

973 Sun, C., Hu, Z., Zheng, T., Lu, K., Zhao, Y., Wang, W., Shi, J., Wang, C., Lu, J., Zhang, D. and Li, Z.,
974      2017. RPAN: rice pan-genome browser for∼ 3000 rice genomes. Nucleic acids research, 45(2),
975      pp.597-605.

976

977 Sun, L., Li, X., Fu, Y., Zhu, Z., Tan, L., Liu, F., Sun, X., Sun, X. and Sun, C., 2013. GS 6, a member

978    of the GRAS gene family, negatively regulates grain size in rice. Journal of integrative plant
979    biology, 55(10), pp.938-949.
980
981    Takano-Kai, N., Jiang, H., Kubo, T., Sweeney, M., Matsumoto, T., Kanamori, H., Padhukasahasram,
982        B., Bustamante, C., Yoshimura, A., Doi, K. and McCouch, S., 2009. Evolutionary history of GS3,
983        a gene conferring grain length in rice. Genetics, 182(4), pp.1323-1334.
984
985    Tao, Y., Zhao, X., Mace, E., Henry, R. and Jordan, D., 2019. Exploring and exploiting pan-genomics
986        for crop improvement. Molecular Plant, 12(2), pp.156-169.
987    Tettelin, H., Masignani, V., Cieslewicz, M.J., Donati, C., Medini, D., Ward, N.L., Angiuoli,
988
989    S.V., Crabtree, J., Jones, A.L., Durkin, A.S. and DeBoy, R.T., 2005. Genome analysis of multiple
990        pathogenic isolates of Streptococcus agalactiae: implications for the microbial "pan-genome".
991        Proceedings of the National Academy of Sciences, 102(39), pp.13950-13955.
992
993    Thomson, M.J., Singh, N., Dwiyanti, M.S., Wang, D.R., Wright, M.H., Perez, F.A., DeClerck, G.,
994        Chin, J.H., Malitic-Layaoen, G.A., Juanillas, V.M. and Dilla-Ermita, C.J., 2017. Large-scale
995        deployment of a rice 6 K SNP array for genetics and breeding applications. Rice, 10(1), pp.1-13.
996
997    Tong, H., Jin, Y., Liu, W., Li, F., Fang, J., Yin, Y., Qian, Q., Zhu, L. and Chu, C., 2009. DWARF
998        AND LOW-TILLERING, a new member of the GRAS family, plays positive roles in
999        brassinosteroid signaling in rice. The Plant Journal, 58(5), pp.803-816.
1000
1001   Torkamaneh, D., Laroche, J., Valliyodan, B., O'Donoughue, L., Cober, E., Rajcan, I., Abdelnoor,
1002       R.V., Sreedasyam, A., Schmutz, J., Nguyen, H.T. and Belzile, F., 2019. Soybean haplotype map
1003       (GmHapMap): a universal resource for soybean translational and functional genomics. BioRxiv,
1004       p.534578.
1005
1006   Tranchant-Dubreuil, C., Rouard, M. and Sabot, F., 2019. Plant pangenome: impacts on phenotypes
1007       and evolution. Annual Plant Reviews.
1008
1009   Walkowiak, S., Gao, L., Monat, C., Haberer, G., Kassa, M.T., Brinton, J., Ramirez-Gonzalez, R.H.,
1010       Kolodziej, M.C., Delorean, E., Thambugala, D. and Klymiuk, V., 2020. Multiple wheat genomes
1011       reveal global variation in modern breeding. Nature, 588(7837), pp.277-283.
1012
1013   Wang, D.R., Agosto-Pérez, F.J., Chebotarov, D., Shi, Y., Marchini, J., Fitzgerald, M., McNally, K.L.,
1014       Alexandrov, N. and McCouch, S.R., 2018. An imputation platform to enhance integration of rice
1015       genetic resources. Nature communications, 9(1), pp.1-10.
1016
1017   Wang, J. and Zhang, Z., 2021. GAPIT Version 3: boosting power and accuracy for genomic
1018       association and prediction. Genomics, Proteomics & Bioinformatics.
1019
1020   Wang, W., Mauleon, R., Hu, Z., Chebotarov, D., Tai, S., Wu, Z., Li, M., Zheng, T., Fuentes, R.R.,
1021       Zhang, F. and Mansueto, L., 2018. Genomic variation in 3,010 diverse accessions of Asian

1022    cultivated rice. Nature, 557(7703), pp.43-49.

1024    Wang, Y., Xiong, G., Hu, J., Jiang, L., Yu, H., Xu, J., Fang, Y., Zeng, L., Xu, E., Xu, J. and Ye, W.,
1025        2015. Copy number variation at the GL7 locus contributes to grain size diversity in rice. Nature
1026        genetics, 47(8), pp.944-948.

1028    Weng, J., Gu, S., Wan, X., Gao, H., Guo, T., Su, N., Lei, C., Zhang, X., Cheng, Z., Guo, X. and Wang,
1029        J., 2008. Isolation and initial characterization of GW5, a major QTL associated with rice grain
1030        width and weight. Cell research, 18(12), pp.1199-1209.

1032    Xu, F., Fang, J., Ou, S., Gao, S., Zhang, F., Du, L., Xiao, Y., Wang, H., Sun, X., Chu, J. and Wang,
1033        G., 2015. Variations in CYP 78 A 13 coding region influence grain size and yield in rice. Plant,
1034        cell & environment, 38(4), pp.800-811.

1036    Xu, K., Xu, X., Fukao, T., Canlas, P., Maghirang-Rodriguez, R., Heuer, S., Ismail, A.M., Bailey-
1037        Serres, J., Ronald, P.C. and Mackill, D.J., 2006. Sub1A is an ethylene-response-factor-like gene
1038        that confers submergence tolerance to rice. Nature, 442(7103), pp.705-708.
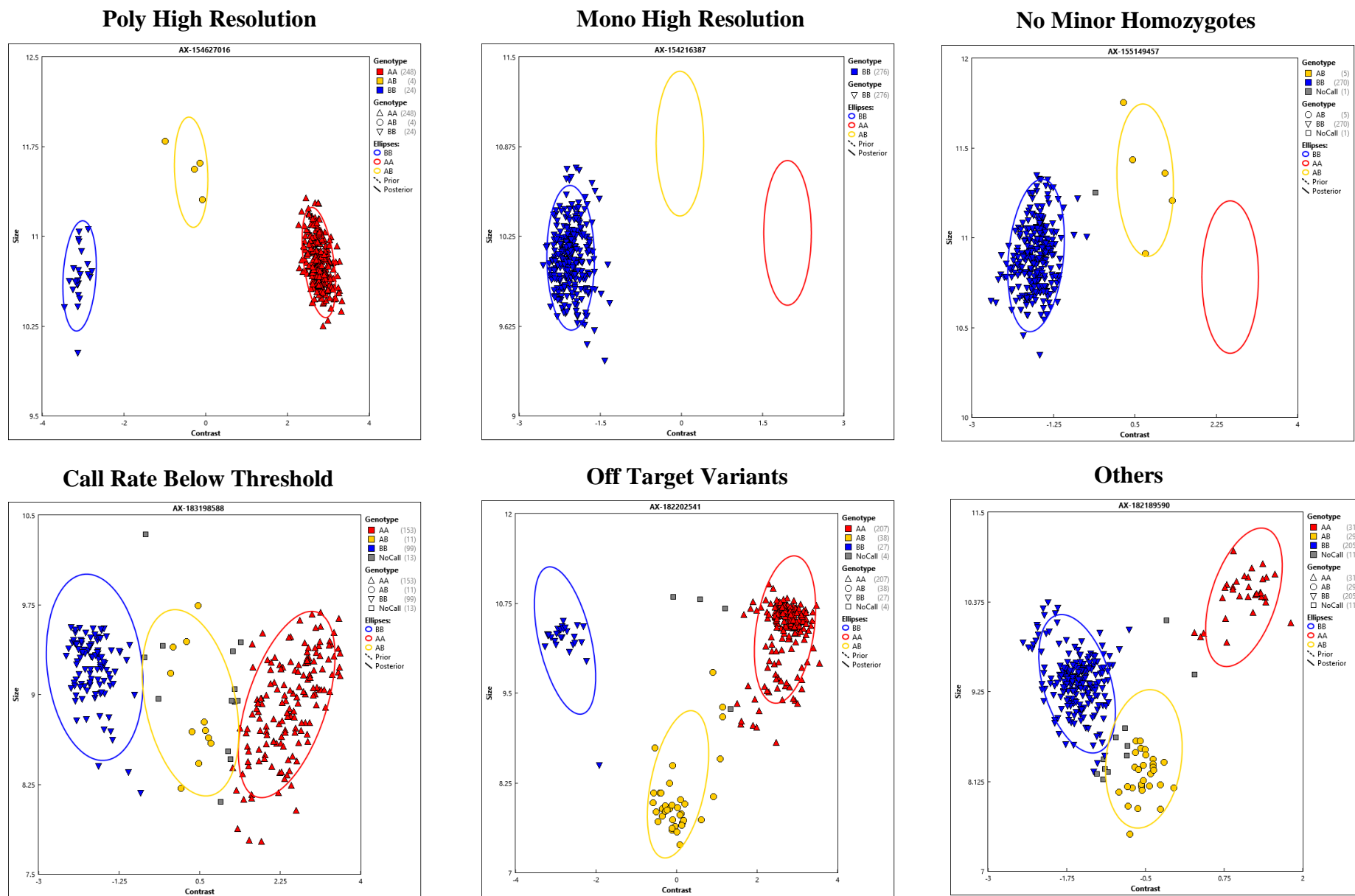1039    Zhao, Q., Feng, Q., Lu, H., Li, Y., Wang, A., Tian, Q., Zhan, Q., Lu, Y., Zhang, L., Huang, T. and
1040        Wang, Y., 2018. Pan-genome analysis highlights the extent of genomic variation in cultivated and
1041        wild rice. Nature genetics, 50(2), pp.278-284.

1043    Zhou, Y., Miao, J., Gu, H., Peng, X., Leburu, M., Yuan, F., Gu, H., Gao, Y., Tao, Y., Zhu, J. and
1044        Gong, Z., 2015. Natural variations in SLG7 regulate grain shape in rice. Genetics, 201(4),
1045        pp.1591-1599.

1047    Zhu, X., Liang, W., Cui, X., Chen, M., Yin, C., Luo, Z., Zhu, J., Lucas, W.J., Wang, Z. and Zhang, D.,
1048        2015. Brassinosteroids promote development of rice pollen grains and seeds by triggering
1049        expression of Carbon Starved Anther, a MYB domain protein. The Plant Journal, 82(4), pp.570-
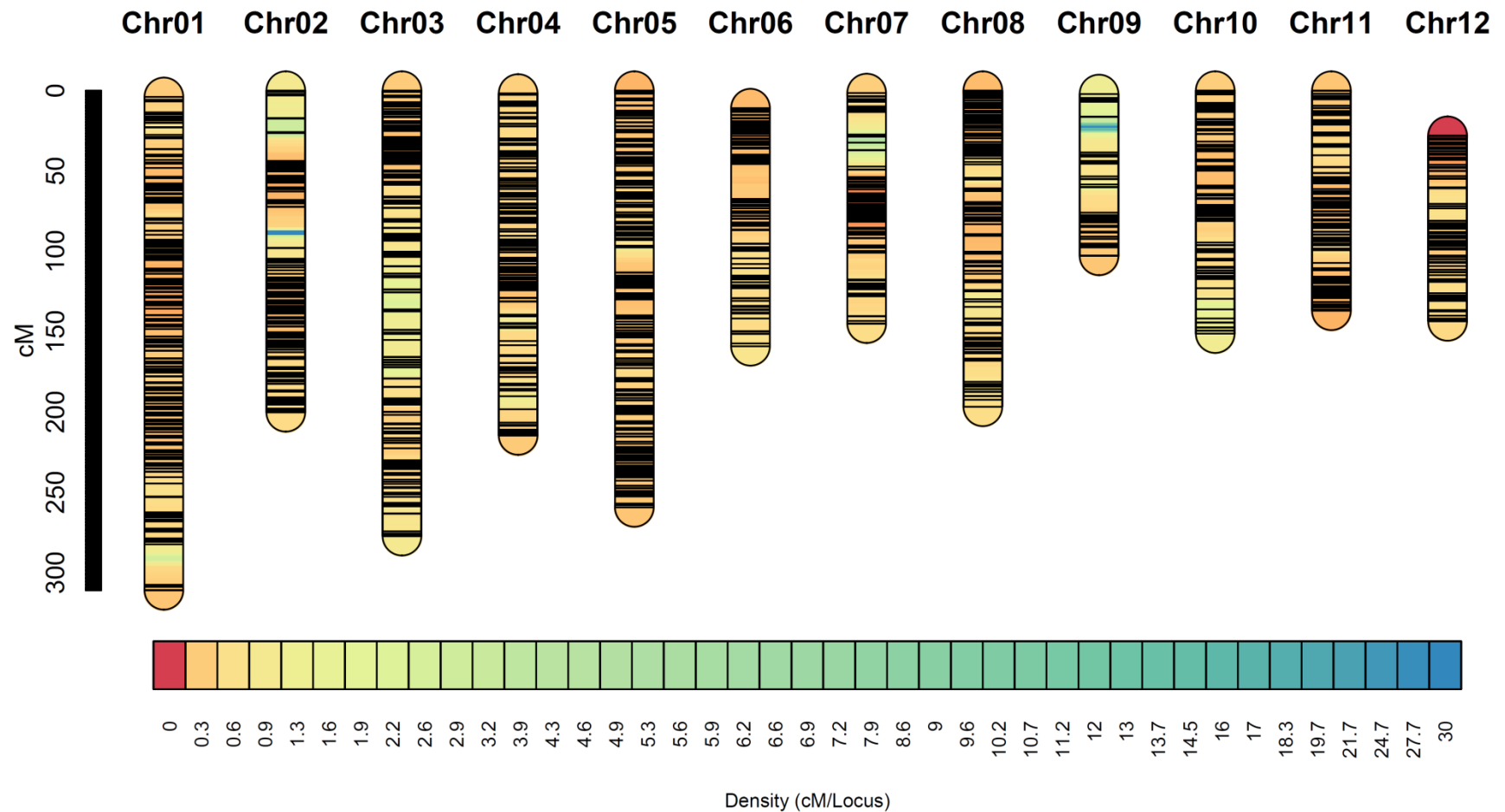1050        581.

1051

1052

1053

1054

1055

1056

1057

**Figure 1. Pan-genome-wide distribution and structural annotation of SNPs tiled on rice pan-genome array. a)** The density of SNPs in 100 kb bins across twelve chromosomes from Nipponbare reference genome. **b)** The density of SNPs in 100 kb bins across sub-population/varietal group-specific contigs of 3K rice pan-genome. **c)** Proportionate distribution of SNPs across genic and intergenic regions of the Nipponbare genome as well as across sub-population/varietal group-specific sequences (SS/VGS) identified in 3K pan-genome. (Adm: Admixed Group; AROG: Aromatic Group; AUSG: *Aus* group; IG: *Indica* Group; JG: *Japonica* Group)
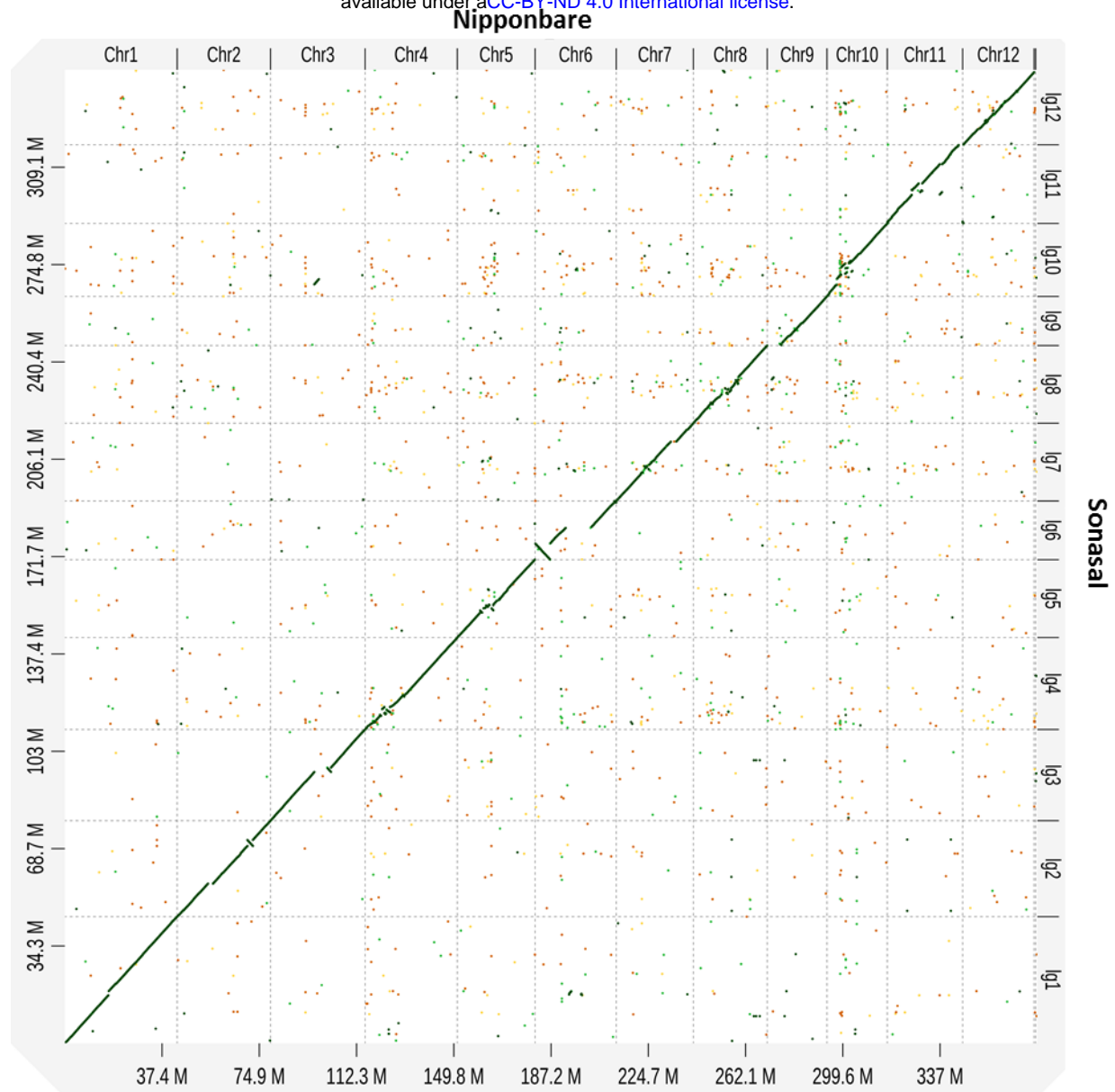
**Figure 2. Representative images of six different SNP classes identified using diversity panel genotyping data generated with rice pan-genome genotyping array (RPGA).**
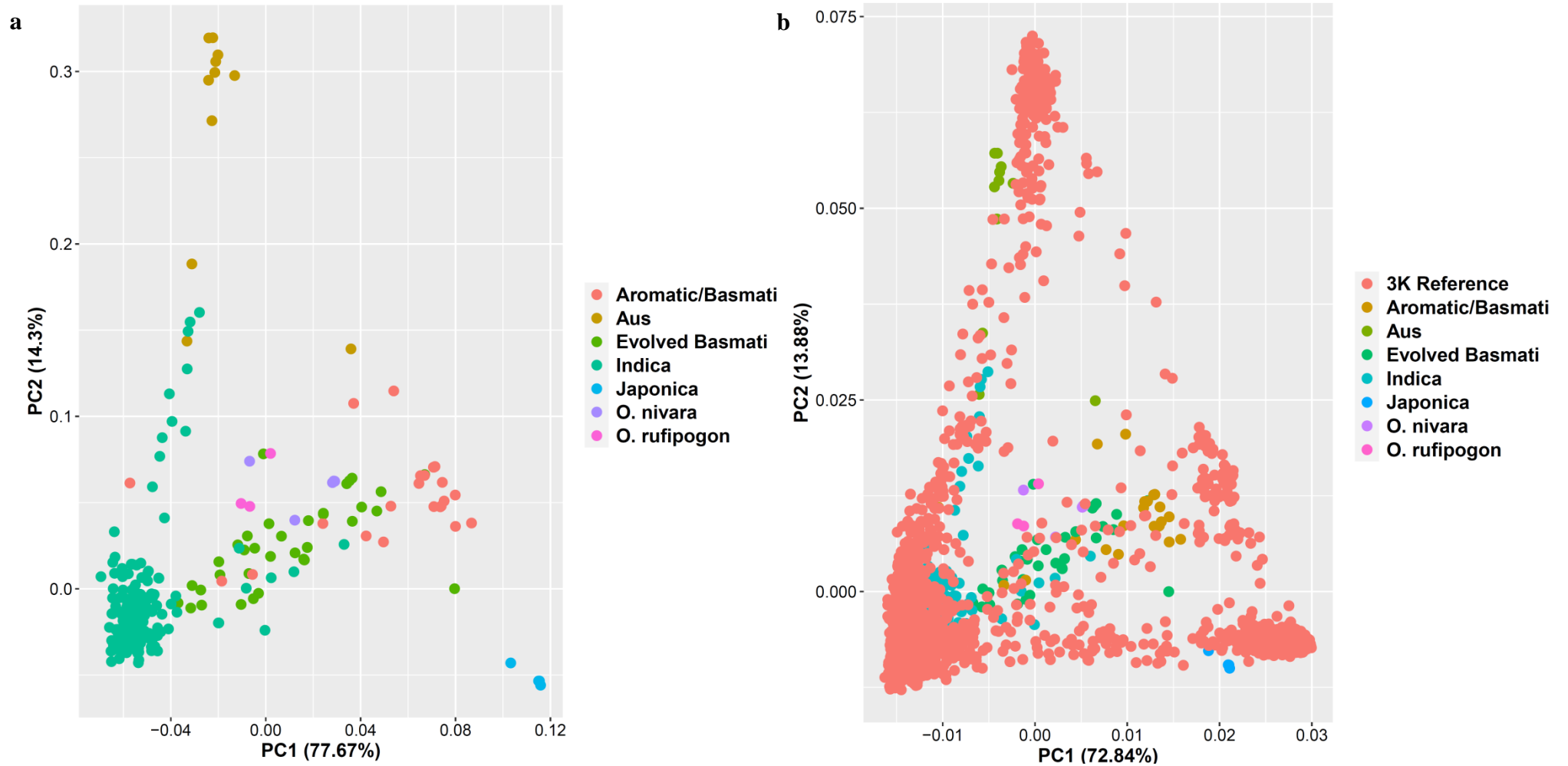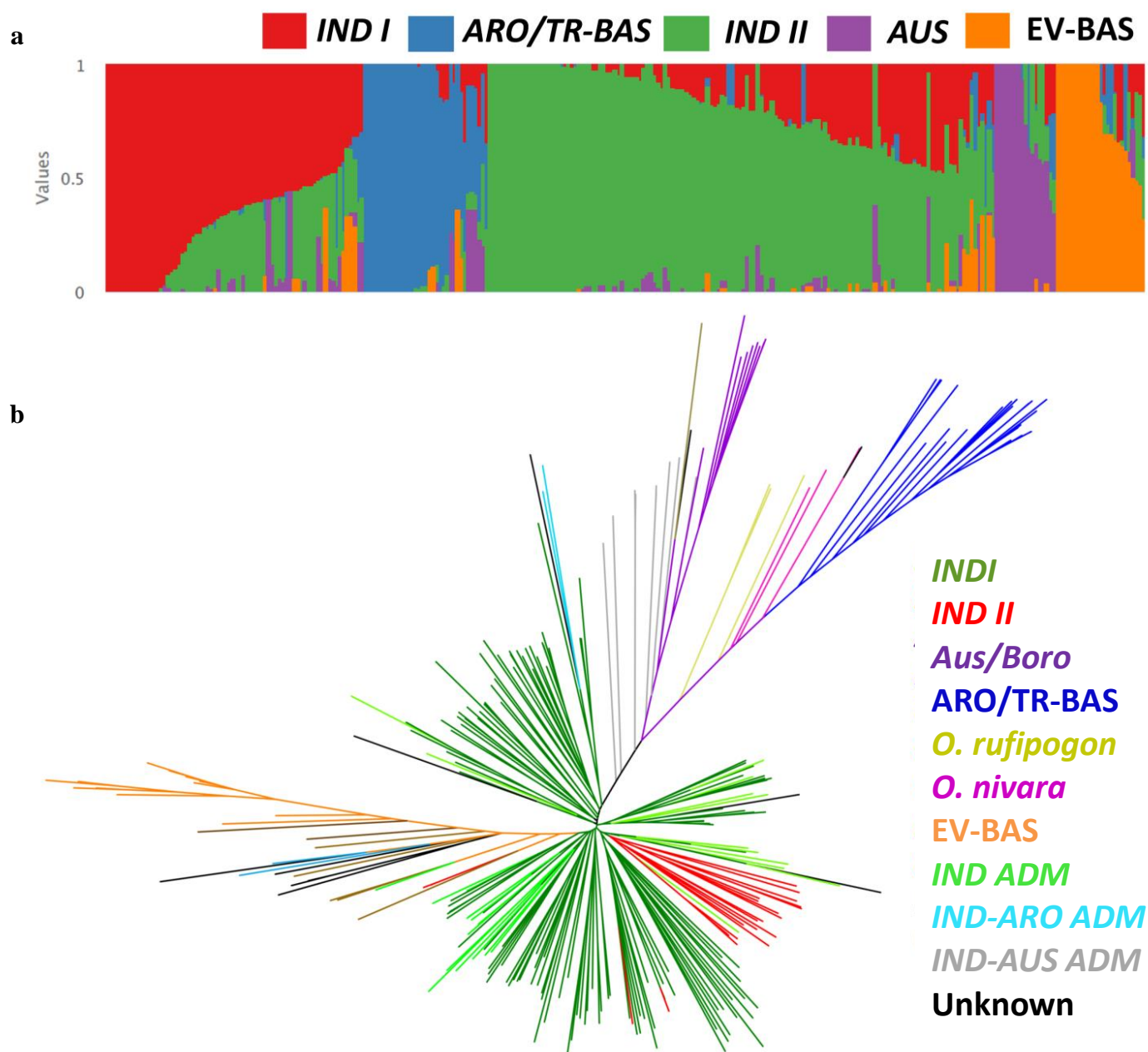
**Figure 3. Ultra high-density genetic linkage map of a RIL mapping population (Sonasal × Pusa Basmati 1121) developed using genotype data generated with rice pan-genome genotyping array (RPGA).** SNPs are represented as black horizontal lines on each linkage group. The vertical scale at the right depicts the genetic distance in cM. The color scale at the bottom represents the genetic distance represented by each locus across linkage groups in the form of density (cM/Locus).
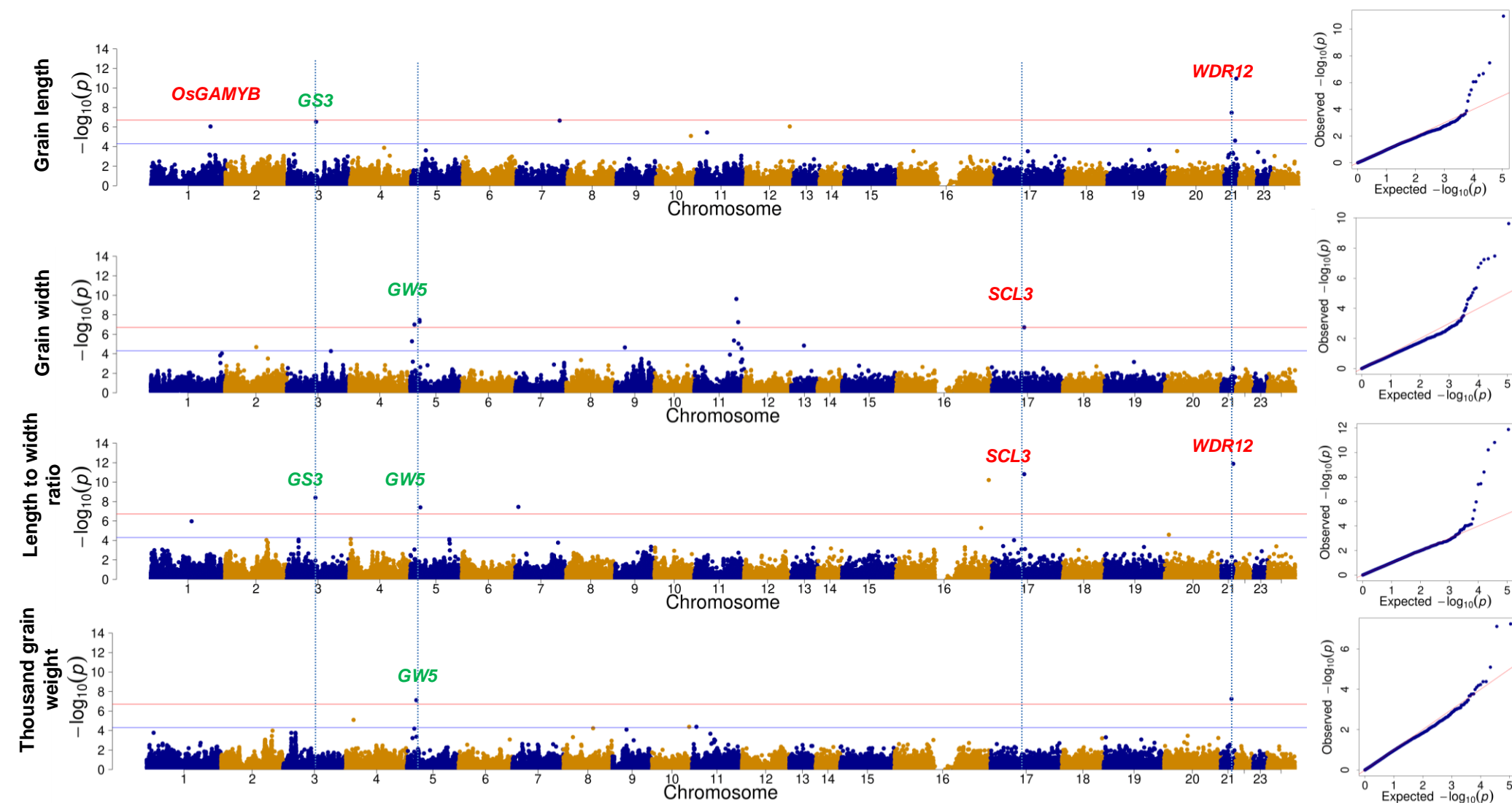
**Figure 4. Dot plot illustrating syntenic relationship between the Sonasal (aromatic) draft genome and Nipponbare (*japonica*) reference genome**. The 12 LGs (linkage groups) correspond to 12 Sonasal chromosomes.

**Figure 5. Principal component analysis (PCA)-based molecular diversity analysis of a diversity panel with rice accessions representing different rice sub-populations (*indica, japonica, aus,* and aromatic/Basmati). a)** PCA of 271 rice accessions, where, principal component 1 (PC1) and PC2 explains 77.67% and 14.43% genetic variation, respectively, in the accessions analyzed. **b)** PCA of 271 Indian rice accessions along with 3K rice genome reference panel, where, PC1 and PC2 explains 72.84% and 13.88% genetic variation, respectively, in the accessions analyzed.
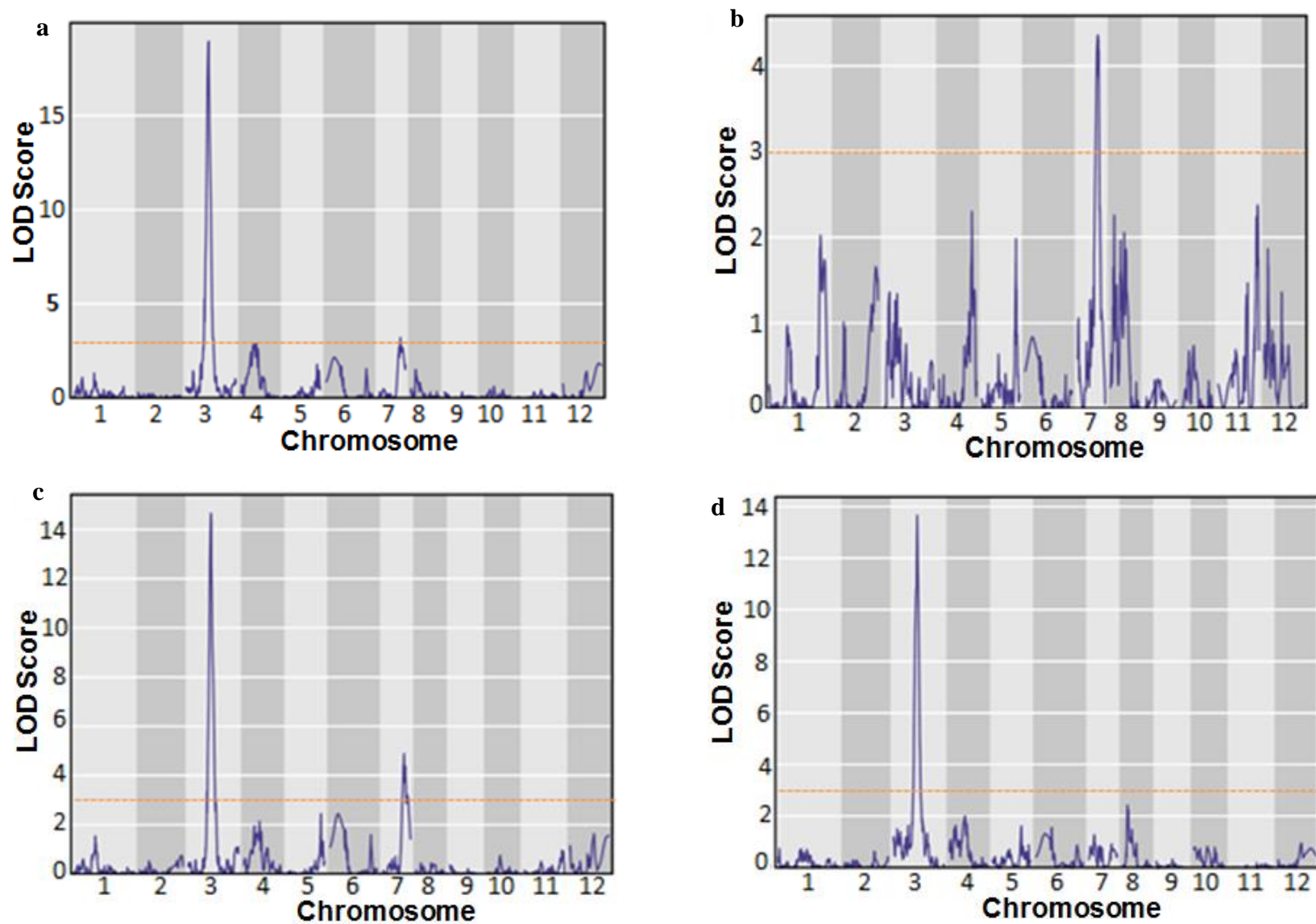
**Figure 6. Population structure and neighbor-joining tree-based molecular diversity in a diversity panel of Indian rice accessions. a**) Population structure of 265 diverse rice accessions with K=5, Inferred populations are color-coded with five different colors. **b**) Neighbor-joining tree of 271 diverse rice accessions belonging to three different cultivated and wild rice species viz. *O. sativa, O. nivara* and *O. rufipogon*. The branches of phylogenetic are colored to depict their respective genetic ancestry. ADM: Admixture, ARO: Aromatic, *Aus/Boro*: *Aus/Boro* ecotypes, EV-BAS: Evolved Basmati, *IND*: *Indica,* TR-BAS: Traditional Basmati.
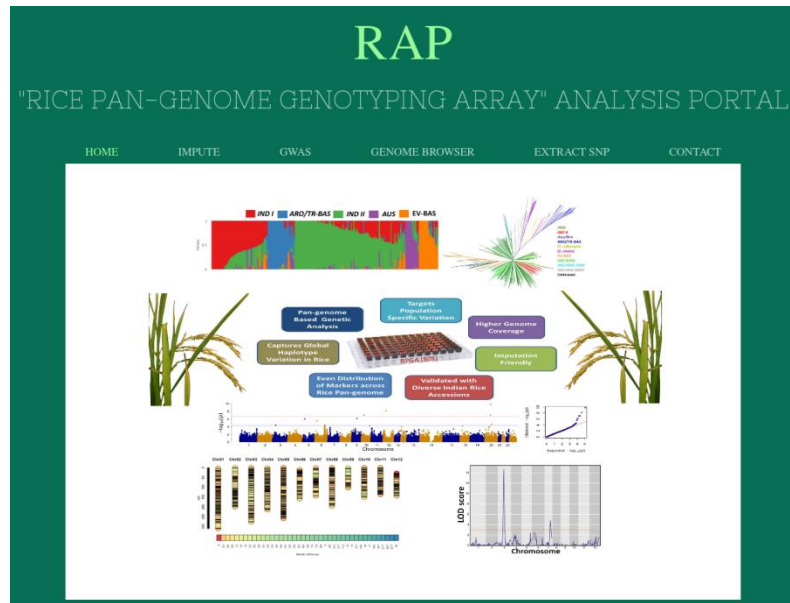
**Figure 7. Manhattan plots and QQ plots depicting results of genome-wide association study (GWAS) for grain length, grain width, length-to-width ratio, thousand-grain weight trait**. Chromosomes 1 to 12 represent twelve Nipponbare chromosomes, whereas, chromosomes 13 to 24 represent twelve sub-population group specific pseudo-chromosomes (Chr13:Adm, Chr14:AROG11, Chr15:AUSG6, Chr16:IG1, Chr17:IG2, Chr18:IG3, Chr19:IG4, Chr20:IG5, Chr21:JG10, Chr22:JG7, Chr23:JG8, Chr24:JG9). Trait associated loci coinciding with previously known grain weight genes are marked in green. The blue horizontal line represents a stringent Benjamini-Hochberg threshold whereas the red horizontal line represents a less stringent threshold. (Adm: Admixed Group; AROG: Aromatic Group; AUSG: *Aus* group; IG: *Indica* Group; JG: *Japonica* Group).

**Figure 8. Molecular mapping of QTLs regulating grain size/weight traits in rice. a)** Thousand-grain weight, **b)** grain length, **c)** grain width, and **d)** length-to-width ratio. The orange color horizontal lines represent the significance threshold calculated using 1000 permutations and red asterisk sign represents the significant QTL peaks.

**Figure 9**. **Interface for Rice Pan-genome Genotyping Array Analysis Portal (RAP). a)** RAP home page, **b)** Impute function available in RAP which enable imputation of user uploaded RPGA genotype data, **c)** GWAS utility available in RAP which enable user to conduct pan-genome-based GWAS, and **d)** RAP enable extraction of SNPs from Locus ID and genomic coordinates.

# Rice Pan-genome Genotyping Array (RPGA)

**"Rice Pan-genome Genotyping Array (RPGA)"** is a first-ever pan-genome-based SNP genotyping assay developed for crop plants. It enables researchers to target population-specific genomic variation by assaying SNP markers unique to different rice populations (*indica*, aromatic, *aus* and *japonica*, etc.), in addition to markers from the *japonica* Nipponbare reference genome. Therefore, RPGA provides higher genomic coverage across diverse rice populations compared to other genotyping arrays available for rice. This makes RPGA as an ideal SNP genotyping solution for pan-genome based genetic studies including complex trait dissection (GWAS, eQTL and bi-parental QTL mapping, etc.), population and evolutionary genetics, and molecular breeding applications (marker-assisted selection and genomic selection, etc.) to drive genetic improvement of rice.
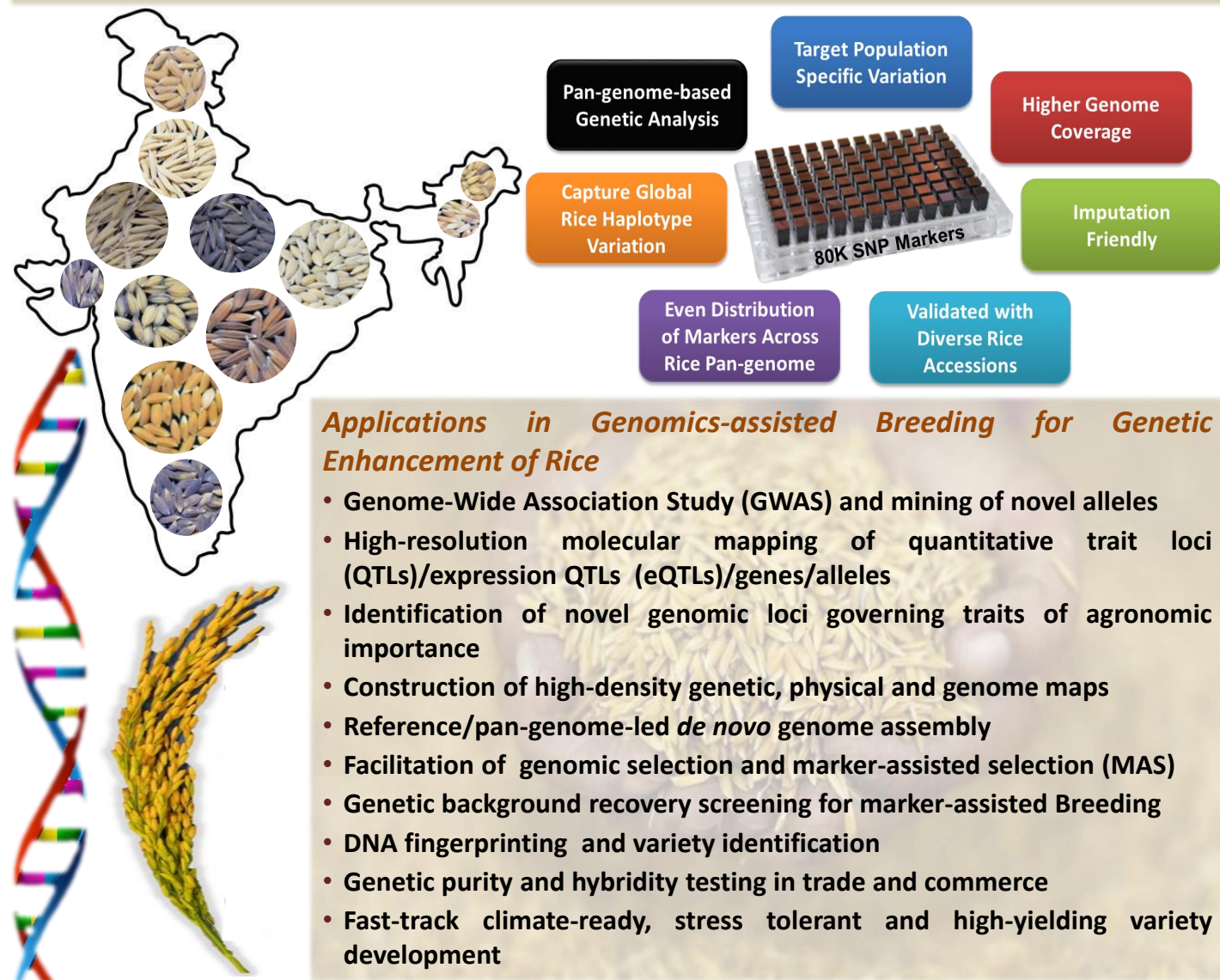


Pan-genome-based Genetic Analysis

Target Population Specific Variation

Higher Genome Coverage

Capture Global Rice Haplotype Variation

80K SNP Markers

Imputation Friendly

Even Distribution of Markers Across Rice Pan-genome

Validated with Diverse Rice Accessions

## *Applications in Genomics-assisted Breeding for Genetic Enhancement of Rice*

- **Genome-Wide Association Study (GWAS) and mining of novel alleles**
- **High-resolution molecular mapping of quantitative trait loci (QTLs)/expression QTLs (eQTLs)/genes/alleles**
- **Identification of novel genomic loci governing traits of agronomic importance**
- **Construction of high-density genetic, physical and genome maps**
- **Reference/pan-genome-led *de novo* genome assembly**
- **Facilitation of genomic selection and marker-assisted selection (MAS)**
- **Genetic background recovery screening for marker-assisted Breeding**
- **DNA fingerprinting and variety identification**
- **Genetic purity and hybridity testing in trade and commerce**
- **Fast-track climate-ready, stress tolerant and high-yielding variety development**

Figure 10. Overview of "Rice Pan-genome Genotyping Array (RPGA)"