

## Peer Review File

---

Pan-genome inversion index reveals evolutionary insights into the subpopulation structure of Asian rice (*Oryza sativa*)



**Open Access** This file is licensed under a Creative Commons Attribution 4.0

International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to

the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. In the cases where the authors are anonymous, such as is the case for the reports of anonymous peer reviewers, author attribution should be to 'Anonymous Referee' followed by a clear attribution to the source work. The images or other third party material in this file are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

## Reviewers' Comments:

### Reviewer #1:

#### Remarks to the Author:

This manuscript presented 16 rice high quality assemblies, including cultivated and two wild species, and focused on the analysis regarding inversions in terms of evolution, effects on gene regulation, recombination rate et al. I have several major concerns as follow:

1: The authors have tried to resolve inversions in population scale by combining 18 rice high-quality assemblies. However, to my knowledge, up to 60 rice high quality assemblies have been available so far, it's much better to combine all available assemblies together to investigate inversions at population level. Moreover, the authors investigated inversion distribution in 5 additional wild species with high-quality genomes (line 246-250), unfortunately, they were noted "unpublished data", those wild species will enhance the quality of this study.

2: The authors should compare the inversions in this study with inversions which had been published, such as that data in Qin et al. Cell, 2021, wherein they also identified tons of inversions. Because the authors did not perform the comparison, so the author can determine how many new inversions have been identified, and therefore, the description at line 236 is not accurate: "two of which (INV060390 and INV080710) were previously reported". Moreover, the descriptions in the section of "five largest inversions" are boring, it was just a list of the distribution of these inversion, and some distribution information have been revealed before.

3: In the next four sections after "five largest inversions", I did not see any attractive and new contents which will significantly promote rice biological and evolutionary studies. For example: the conclusion in the section "Characterization of Transposable Element Content within Inversions and Breakpoints" (our results reveal an enrichment of TE related sequences both within inversions and at their breakpoints), have been reported in several papers. Same thing for the conclusion "a marked suppression of genetic recombination is associated with inversions" of the section "Recombination Rate and Genomic Inversions".

The investigation about the effect of inversions on the genes located within inverted regions, and their expression wasn't appropriate, the effect of other sequence variations, such as SNP and SV, were not considered, these variations were supposed to have much bigger and direct impacts on gene expression than inversions, whereas, they were overlooked in this investigation, and without discussions. Additionally, regarding the section "Phenotypic consequences of inversions: Inversion Cluster 92", based on the description and my understanding, it was much better to say the consequence of SNP.

### Reviewer #2:

#### Remarks to the Author:

This paper details the comparison of 18 long-read based genomes of rice. Of the 18, 16 had been published previously, one was updated and one was new. The goal of the study was to use high-quality genomes to infer the position and prevalence of inversions, which are likely to be an important class of structural variant (SV). There have been many studies of rice genomes recently, so the novel contribution of this one is its focus on inversions events. As set forth at the end of the Introduction, the paper makes five main claims – i.e., a catalog of 1054 inversions, an inversion rate, and some biological effects (lower LD and at least one inversion that contributes to delayed flowering).

I have many comments on the paper, and list them according to my progress through the paper. Many are admittedly quite minor and are offered in the spirit of helpful criticism. However, others are, I believe, more substantive and require substantial revision and perhaps more thought. Overall, the

data presented in the paper did not, in my mind, match the claims of the paper.

- Given that inversions are the main thrust of the paper, the paper is narrowly written, with a focus only on rice. There is lots of previous evidence – although mostly not discussed in this manuscript – that inversions affect phenotypes, gene expression, are mediated TEs, etc. From my perspective, it's a missed opportunity to not put this rice work into the context of the wider plant literature (maize, tomato, grapes, evolution, etc. etc.). The only real attempt at generality is lines 90 to 93, and a bit in the first paragraph of the Discussion.

- The new or updated genomes from rufi and punct are diploid, I think – i.e., from not a naturally selfing lineage, like rice. How were diploid genomes treated? Details on phasing, haplotypes and heterozygosity are lacking, but probably important (particularly if diploid chromosomes have SVs, so that their treatment is key to inferences)

- This is admittedly a stylistic thing, but I find it awkward to list the main results at the end of the Discussion. I'd rather know what questions and going to be asked and why. Moreover, on careful reading, I feel as if at least #2 and #3 were overstated.

- I struggled to follow the sampling and the nomenclature (e.g., Xi-adm MH63). It'd be nice to have a sampling table with acronyms, taxon of origin, etc.

- I'm confused by a great deal about the pan-genome (1st two sections of results), because I am not sure how genomes were combined to create a pan-genome. And later, the paper extols the virtue of using a pan-genome free approach. I think (?) that the term pan-genome was used mostly to say "we looked at all of the genomes", but if so that usage is a bit confusing given that pan-genome has gained a more specific meaning.

- p. 175 – I appreciate the detail given to the various workflows in the M&M to estimate pairwise inversions. It gives the impression of great care!

- Line 181 – is this statement relative to IRGSP or to Oruf and Opunct? Generally lines 181 to 190 were pretty tough to follow, given lack of knowledge about sampled genomes. Again, a table would be nice or more explanation.

- Line 200 – I desperately would have liked to see the information in this paragraph summarized on a phylogeny, showing species and group-specific number of inversions on nodes. Given that this paper repeatedly touts a strong phylogenetic basis as the rationale for sampling, the lack of a phylogeny is a somewhat glaring omission.

- I was a bit confused by Figure 1. It's nice in the sense that it shows all of the data, but I think there are no species-specific inversions in rice (vs. other species), right? And what does each blue bar represent? I assume each one is an inversion. It seems to me that there are many inversions shared among the Asian Rice group XI on the right hand side of the graph, but not circled as group specific. Perhaps it's a question of the level of detail, because it is very hard to show all of the inversions, but unfortunately I did not find the figure particularly useful. (I'd love to see inversions in the context of chromosomes, centromeres, etc.)

- Line 208 – are subpopulations groups? (as in Figure 1?)

- Line 216 – if I followed correctly, the 3K-RGP is a short-read dataset. It's not clear to me how that could be used to test/confirm inversions, particularly since later the claim is made that short read data are not useful for inferring inversions after comparing the results of this paper to the Fuertes paper. Moreover, details of methods, numbers used, 'high coverage', etc., are lacking, such that is hard to follow the basis for conclusions.

- Line 256 - TE information about breakpoints is interesting. By "analyzing TE content across the inversion index", is this all 1054 inversions?
- Figure 3C. I'm confused and need some explanation in the legend. The lines in the middle seem to be inverted, but the arrows go in the same direction. The TE (0025) seems to be an LTR that is split in both cases. What is INT?
- How many inversions were validated with bioNano, as in Figure 2? Overall, I'm not convinced that all 1054 inversions were independently validated (as claimed in point #2 in the Intro). I do believe bioNano, but details are lacking (how low is the resolution? How many could be confirmed?) I believe that short-reads are useful, but details are lacking here, too, but (again) its hard to claim they validate inversions on the one hand but are not useful at all on the other (e.g., Fuertes).
- Lines 285 - I'm assuming that 10.9% and 7.3% is much higher than the genome average, but it'd be nice to have an explicit comparison to the genome average to drive this point home.
- About the inversion rate estimates in the Discussion (line 411 and following). It may be that I'm misreading things, but I think they are generally wrong and perhaps horribly so. Here's why. There may be 22 post-inversion events in rice, but (if one thinks in phylogenetic terms) there are many many more years accumulated across the rice lineages on the tree than 10,000. As a brief example, let's assume (for simplicity) that 17 rice genomes diverged 10,000 years ago. If that were true, then the numerator in the rate calculation should be  $17 \times 10,000$  years, not 10,000 years, so that estimate is inflated about 17-fold. Of course, we don't know exactly when each of the separate rice genomes diverged from one another, so the estimate of  $17 \times 10,000$  years is too many. But hopefully the point is made that the calculations reported in this section may be way off and that the problem may require some consideration of population genetics given the sample.
- Line 432 - the comparison to Fuertes. I'd have to read that paper carefully to see how it was done, but the title implies there were 3,000 individuals. That suggest, I think, that the entire dataset was used as evidence to support or not the inference of inversions. I'd expect them to have found many more inversions, but that the inversions in this paper would be a subset of their total set. I think that'd be a more fair comparison. That said, it is indeed puzzling that there is only 194 out of 1054 that overlapped! While I agree with the authors that short read data is certainly less accurate than long-read data, it does make one wonder about the accuracy of the 18 assemblies and whether there is no only errors in Fuertes but also in assemblies that mislead inversion inference.
- One time consuming but convincing way to validate inversions is to find long reads that span the inferred junctions. I don't think that was done her, but it would certainly go a long way to confirming the inversion inferences more convincingly.

Reviewer #3:

Remarks to the Author:

Zhou et al. present a pan-genome analysis of the major sub-populations of asian rice and two wild rice species represented by a set of 18 whole genome assemblies. the two wild rice assemblies are new and are provided with this manuscript. The other assemblies have been published before. The authors main point of analysis is the cataloging of inversions larger than 100 bp and contextualizing these with data on recombination, LD, selection and gene expression. As the authors point out, the detection of inversions (also larger) is not novel for plant genomes (maybe more references to recent pan-genome studies in cereals would be justified here?), however, here maybe a first comprehensive catalog for representatives of a crop species' subpopulations is provided.

The study reads well and is well presented. The data displays sometimes are rather basic - often just direct exports out of commercial analysis and visualization software? Intuitivity of the figure displays could be improved - maybe?

The study addresses an obvious point in comparative genomics as more genomes of the same species become available. The authors are well aware of artifacts introduced into analyses when different assembly qualities affect the analysed assemblies. Therefore, I was surprised, given the today's costs for making a HiFi assembly for rice (haploid assembly consumable costs in the few hundred dollars range!) that the authors did not make an effort here to have really absolutely comparable datasets. This may only incrementally change the presented results, however, it is a weakness of the study that could have been avoided with modest effort and investment. The same argument applies to the annotation which was done with the same pipeline only for the Asian rice.

The spectrum of inversion distribution: the authors say they are randomly distributed - Ext. figure 3 should support this claim. While the figure legend is not conclusive and the figure itself is basic, I think the authors should have made an effort define size classes and redo the genome distribution scan. Is it true that all sizes of inversions are evenly distributed along the chromosomes? This is counterintuitive.

The authors used the 3K rice data to genotype for the presence of inversions, however, they performed this only on a selected set of genotypes representing the 15 sub-populations. Is it really computationally so intense to include all 3000? Or is this a problem of sequence coverage in a certain proportion of the 3K dataset? I would love to see the analysis on all 3000! I wonder whether the authors tried to use the inversion catalog as a proxy to model and detect additional inversions in the full 3K dataset as the 15 genotypes sequenced will for sure not give the full pan-genome inversion spectrum. The authors outline the technical feasibility in the discussion. Furthermore, it would be very instructive to give an estimate for WGS coverage to robustly scan for the presence of inversions at population scale?

TE landscape was analysed in context of inversions and it was detected that inversions are enriched with TE content in rice. The analysis of inversions could go further here. Have you systematically assessed the TE at the Inversion junctions? This could reveal patterns of the mechanistic involvement for the occurrence of inversions in the rice genome.

Gene expression and inversion: the authors showed examples where inversions and expression levels of genes in inverted regions and their orthologs in "non-inverted" haplotypes. The data is interpreted as inversions are causal, which is not unlikely, however, the authors do not show any functional validation (which is not trivial) and analysis and they also do not provide complementary datasets (methylation, HiC, ATAC etc.) which would support their hypothesis and interpretations. I would also recommend to expand the analysis to genes adjacent or in neighborhood but not directly affected by the inversions - regulatory sequence context can and will be affected also for such genes and one should see a gradient perhaps?

Similarly, the interpretation of positive selection in inverted or non-inverted haplotypes should maybe presented and discussed with more caution.

The authors make a point in the discussion that not merging sequences into a pan-genome graph was an advantage here. Well isn't this obvious? And isn't the need for graph development depending on the analytical goals - I am not sure if I can follow the argument here.

minor issues:

data is all available to what I could track, however, one has to dive into supplements. A clear data availability statement with instructions where to find the details is missing and needs to be added.

persephone visualization is challenging. Maybe it is comprehensive but I don't think it is necessarily intuitive. Have you tried the recent development from John Lovell at Hudson Alpha (GENESPACE - look at BiorXiv)?

I appreciate the authors use ref 46 for wheat in wheat as this is part of their own work, however, in the context still the IWGSC 2018 ref in Science is probably more appropriate? - unless you want to make a point out of impact of Hifi sequencing in wheat - which is not the case at current.

abstract: "effects on gene regulation" - no! you only report correlations between datasets, no functional proof.

intro: "almost" 10 billion by (exactly) 2064 - I find this mix of approximation and exactness curious - sure you are citing here, but ...

intro: why sequence diversity is a natural variation "tool box"??

1 Point-by-point response to REVIEWERS' COMMENTS

2

3 **We really appreciate the reviewers' comments and feel the manuscript is much better**  
4 **and more comprehensive. Following their comments and the editor's recommendation,**  
5 **we have written a major revision that addresses the reviewer's comments below:**

6

7 **Reviewer #1**

8 **Question 1:** The authors have tried to resolve inversions in population scale by combining 18  
9 rice high-quality assemblies. However, to my knowledge, up to 60 rice high quality  
10 assemblies have been available so far, it's much better to combine all available assemblies  
11 together to investigate inversions at population level.

12 Moreover, the authors investigated inversion distribution in 5 additional wild species with  
13 high-quality genomes (See line 246-250), unfortunately, they were noted "unpublished data",  
14 those wild species will enhance the quality of this study.

15

16 **Author response:** Thanks for your suggestions. This is a very good and important comment  
17 that enhances our work.

18 Indeed, including our platinum standard pan-genome resources, there are three recent  
19 papers that released designated "high-quality" genomes, *i.e.*, Qin et al., 2021, Cell; Zhang et  
20 al., 2022, Genome Research, and Shang et al., 2022, Cell Research. The first of the three  
21 papers used PacBio sequencing technology, while the remaining two used Oxford Nanopore  
22 Technology (ONT). Following the reviewer's comments, we re-called inversions by adding  
23 several newly sequenced high-quality genome from both Qin's and Zhang's studies (n=65).  
24 However, the genome sequences from the Shang et al., 2022 paper: "*Genome sequencing*  
25 *data of 251 accessions in this study have been deposited in the NCBI Sequence Read Archive*  
26 *(<https://www.ncbi.nlm.nih.gov/sra>) under BioProjects PRJNA656318 and PRJNA692836*",  
27 have yet to be publicly released (a screenshot as below).

28 Prior to re-calling inversions, we assessed the genome quality of the 33 genomes from  
29 Qin's study, and the 65 genomes from Zhang's study (See line 192 - 208). First, we called  
30 N50s on the contigs (as the scaffold N50 was reported in some cases) to give us a more  
31 accurate representation of a contiguous assembly and we only kept genomes with Contig  
32 N50s > 3Mbp. Secondly, we validated genome assembly correctness by genome-wide dot-  
33 plots and found that some assemblies had scaffolding, or possibly assembly errors. These two  
34 parameters can indicate the quality of genome assemblies and is essential when running a

35 meta-analysis of INVs between genomes, as poor-quality genomes could lower our ability to  
36 detect inversions, or falsely introduce INVs when in-fact it is a genome assembly error.  
37 Ultimately, we added 29 (out of 33) newly sequenced PacBio genomes from Qin’s study, and  
38 28 (out of 65) newly sequenced ONT genomes from Zhang’s study into our 18-genome data  
39 package, for a total of 75 high quality genomes.

40

SRA	PRJNA692836	SRA	PRJNA656318
<a href="#">Create alert</a>	<a href="#">Advanced</a>	<a href="#">Create alert</a>	<a href="#">Advanced</a>

⚠ The following term was not found in SRA: PRJNA692836.

🔍 No items found.

⚠ The following term was not found in SRA: PRJNA656318.

🔍 No items found.

41

42 In addition, we do agree that the 5 additional wild *Oryza* high-quality genomes (*i.e.* *O.*  
43 *nivara* [AA], *O. glaberrima* [AA], *O. barthii* [AA], *O. coarctata* [KKLL] and *O. alta*  
44 [CCDD]), will enhance our study, however we only used these genomes to validate species  
45 specific inversions. In the future, we plan to use these genomes to investigate abiotic, biotic  
46 resistance, and neo-domestication. Thus, for the present paper they were removed.

47

48 References:

49 Qin P, Lu H, Du H, et al. Pan-genome analysis of 33 genetically diverse rice accessions  
50 reveals hidden genomic variations[J]. Cell, 2021, 184(13): 3542-3558. e16.

51

52 Zhang F, Xue H, Dong X, et al. Long-read sequencing of 111 rice genomes reveals  
53 significantly larger pan-genomes[J]. Genome Research, 2022, 32(5): 853-863.

54 Shang L, Li X, He H, et al. A super pan-genomic landscape of rice[J]. Cell Research, 2022:  
55 1-19.

56

57 Shang L, Li X, He H, et al. A super pan-genomic landscape of rice[J]. Cell Research, 2022,  
58 32(10): 878-896.

59

60 **Question 2:** The authors should compare the inversions in this study with inversions which  
61 had been published, such as that data in Qin et al. Cell, 2021, wherein they also identified  
62 tons of inversions. Because the authors did not perform the comparison, so the author can  
63 determine how many new inversions have been identified, and therefore, the description at  
64 line 236 is not accurate: “two of which (INV060390 and INV080710) were previously  
65 reported”.



66        Moreover, the descriptions in the section of “five largest inversions” are boring, it was  
67 just a list of the distribution of these inversion, and some distribution information have been  
68 revealed before.

69

70 **Author response:** We thank the reviewer’s comments. We revisited Qin’s paper in 2021, and  
71 also communicated with Qin. They did not make their data available publicly, but the authors  
72 kindly shared all their inversion data with us. In Qin’s work, they identified 718 inversions,  
73 of which 553 were overlapped with our inversion index. Since the two studies used different  
74 pipelines, we discussed with them to include all of their genomes in our study and their *de*  
75 *nov*o called inversions.

76        In addition, we do agree with reviewer 1 that that the “five largest inversion” section is  
77 boring, so thus we removed that section from the paper and simply described the key  
78 messages (See line 230).

79

80 **Question 3:** In the next four sections after “five largest inversions”, I did not see any  
81 attractive and new contents which will significantly promote rice biological and evolutionary  
82 studies. For example: the conclusion in the section “Characterization of Transposable  
83 Element Content within Inversions and Breakpoints” (our results reveal an enrichment of TE  
84 related sequences both within inversions and at their breakpoints), have been reported in  
85 several papers. Same thing for the conclusion “a marked suppression of genetic  
86 recombination is associated with inversions” of the section “Recombination Rate and  
87 Genomic Inversions”.

88        The investigation about the effect of inversions on the genes located within inverted  
89 regions, and their expression wasn’t appropriate, the effect of other sequence variations, such  
90 as SNP and SV, were not considered, these variations were supposed to have much bigger  
91 and direct impacts on gene expression than inversions, whereas, they were overlooked in this  
92 investigation, and without discussions.

93        Additionally, regarding the section “Phenotypic consequences of inversions: Inversion  
94 Cluster 92”, based on the description and my understanding, it was much better to say the  
95 consequence of SNP.

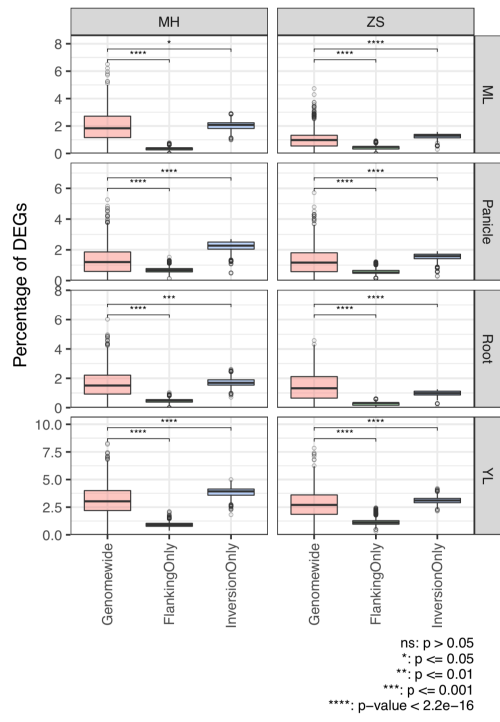
96

97 **Author response:** Yes, we agree with the reviewer, inversions and other structural variants  
98 having an effect on phenotypes has been previously demonstrated. The novelty in this paper  
99 is actually organismal (rice), on a global (pan-genome created for the sub-populations of

100 *Oryza sativa*) and population scale (incorporating 3K data set, that is a compilation of  
101 populations representing known subpopulations), and explored the effect of inversions (a  
102 subtype of large structural variants) on TE and gene composition. We demonstrated that TEs  
103 are enriched at breakpoints which supports TEs as a mechanism to produce inversions, and  
104 the frequency that they may be at a population scale.

105 We agree with the reviewer that a gene's expression is impacted by multiple sequence  
106 variation, *e.g.* SNPs, InDels. As one type of structure variations, inversions could impact  
107 gene expression in multiple ways, including both indirectly (*e.g.*, increase changes in  
108 promotor regions or change binding sites that are affected by 3D structure of sequences), and  
109 directly (via gene disruption). In this case, we are looking for clues of the effect of inversions  
110 on gene expression. We compared portions of DEGs located in genome-wide regions, within  
111 inversion regions, and 20 Kb flanking regions by a permutation test ( $n = 1000$ ). We carried  
112 out comparisons from four tissues, *i.e.* panicle, mature leaves, young leaves and root. The  
113 results showed that the portion of DEGs within inversion were significantly higher than  
114 genome-wide regions, and the portion of DEGs in flanking regions were significantly lower  
115 than genome-wide regions (see a figure as shown below). We do agree with the reviewer that  
116 a gene's expression might be overlooked here, and it might be worth an independent research  
117 study, but here it provide us with clues as to the effects of inversions on gene expression.

118 Lastly, we agree with the reviewer that it was better to say the consequence of SNPs in  
119 Inversion Cluster 92, so we removed "Phenotypic consequences of inversions: Inversion  
120 Cluster 92".



122 **Reviewer #2**

123

124 **Question 1:** Given that inversions are the main thrust of the paper, the paper is narrowly  
125 written, with a focus only on rice. There is lots of previous evidence – although mostly not  
126 discussed in this manuscript – that inversions affect phenotypes, gene expression, are  
127 mediated TEs, etc. From my perspective, it's a missed opportunity to not put this rice work  
128 into the context of the wider plant literature (maize, tomato, grapes, evolution, etc. etc.). The  
129 only real attempt at generality is lines 90 to 93, and a bit in the first paragraph of the  
130 Discussion.

131 **Author response:** Thanks for your comments, and we agree with that. Previous studies in  
132 many species have shown evidence for the biological consequences of inversions. In this  
133 study, we focused our study on inversions with a large number of high-quality genomes in  
134 rice. To generally compare our study with other species, we added a new paragraph in the  
135 introduction section of the manuscript (See line 89-101, and they were highlighted in yellow).

136

137 **Question 2:** The new or updated genomes from rufi and punct are diploid, I think – i.e., from  
138 not a naturally selfing lineage, like rice. How were diploid genomes treated? Details on  
139 phasing, haplotypes and heterozygosity are lacking, but probably important (particularly if  
140 diploid chromosomes have SVs, so that their treatment is key to inferences)

141 **Author response:** Thanks for your comment. Both *O. rufipogon* and *O. punctata* are  
142 autogamous. We believe the degree of outcrossing can be ~30%. The *O. rufipogon* and *O.*  
143 *punctata* samples used in this study were from a single single-seed decent plant. Since the  
144 expected fixation is about 0.7, phasing should not be such an issue for the interpretation of  
145 our data.

146

147 **Question 3:** This is admittedly a stylistic thing, but I find it awkward to list the main results at  
148 the end of the Discussion. I'd rather know what questions and going to be asked and why.  
149 Moreover, on careful reading, I feel as if at least #2 and #3 were overstated.

150

151 **Author response:** Thanks for your comment. You are correct as this is a stylish thing. Our  
152 group has written a few major papers like this in the past and feel it is important to get the  
153 message out there at the beginning. The paper now has 7 points and we kept them brief, to  
154 the point and tried to avoid overstatement (See line 135-162).

155

156 **Question 4:** - I struggled to follow the sampling and the nomenclature (e.g., Xi-adm MH63).  
157 It'd be nice to have a sampling table with acronyms, taxon of origin, etc.

158  
159 **Author response:** Thanks for your comment. We have updated the acronyms, taxon of  
160 origin, and full name of accessions in **Table 1** and **Supplementary Table 2**, and the  
161 acronyms were applied though the manuscript, tables and figures.

162  
163 **Question 5:** - p. 175 – I appreciate the detail given to the various workflows in the M&M to  
164 estimate pairwise inversions. It gives the impression of great care!

165  
166 **Author response:** Thanks for your comments.

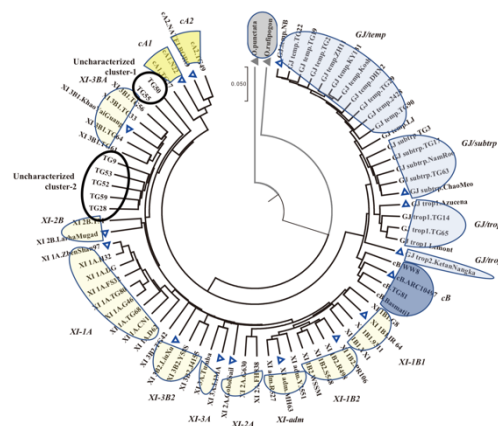
167  
168 **Question 6:** - Line 181 – is this statement relative to IRGSP or to Oruf and Opunct?  
169 Generally lines 181 to 190 were pretty tough to follow, given lack of knowledge about  
170 sampled genomes. Again, a table would be nice or more explanation.

171  
172 **Author response:** Thanks for your comment. We have slightly updated the description (See  
173 line 216-225) and there is a table to help with explanation ([Supplementary Table 3](#)).

174  
175 **Question 7:** - Line 200 – I desperately would have liked to see the information in this  
176 paragraph summarized on a phylogeny, showing species and group-specific number of  
177 inversions on nodes. Given that this paper repeatedly touts a strong phylogenetic basis as the  
178 rationale for sampling, the lack of a phylogeny is a somewhat glaring omission.

179  
180 **Author response:** We agree with reviewer's  
181 comment. To address the suggestion and  
182 improve the study, we added a section (See  
183 line 260) on the analysis of a phylogenetic tree  
184 (see Figure 3). We hope this could help readers  
185 to better understand our study.

186 **Figure 3**



187  
188 **Question 8:** - I was a bit confused by Figure 1. It's nice in the sense that it shows all of the  
189 data, but I think there are no species-specific inversions in rice (vs. other species), right? And

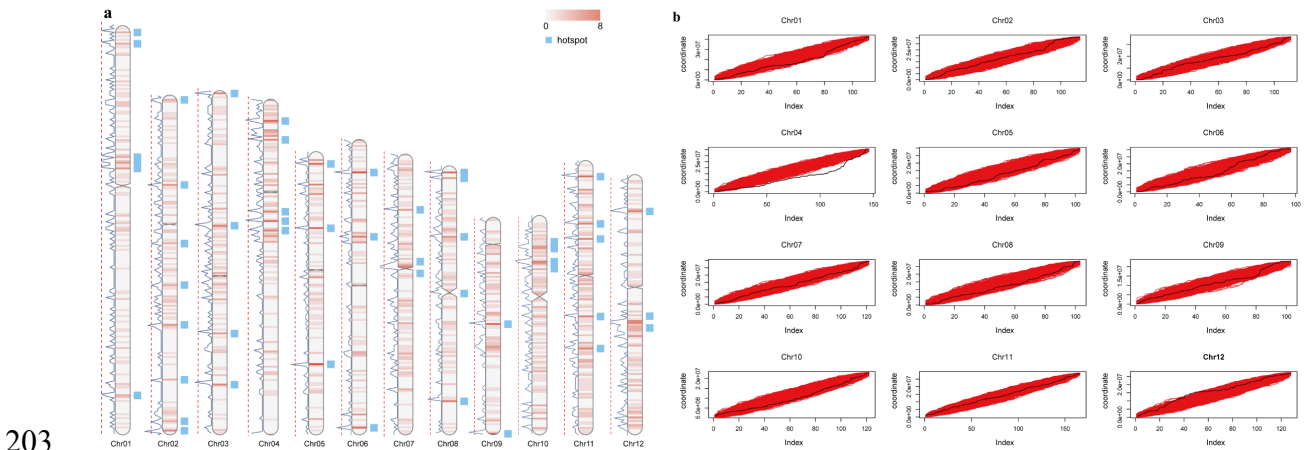
190 what does each blue bar represent? I assume each one is an inversion. It seems to me that  
191 there are many inversions shared among the Asian Rice group XI on the right hand side of  
192 the graph, but not circled as group specific. Perhaps it's a question of the level of detail,  
193 because it is very hard to show all of the inversions, but unfortunately I did not find the figure  
194 particularly useful. (I'd love to see inversions in the context of chromosomes, centromeres,  
195 etc.)

196

197 **Author response:** Thanks for the comments. We agree with the reviewer. We have removed  
198 figure 1 since the figure didn't help as suggested.

199

200 We re-called all inversions based on 75 genomes and following this suggestion, we added to  
201 the study inversion distribution, hotspots, and their overlap with centromeres (See line 241-  
202 258 and [Figure 2](#)).



203

204

205 **Question 9:** - Line 208 – are subpopulations groups? (as in Figure 1?)

206 **Author response:** In the previous version, Line 208 are describing subpopulation specific or  
207 shared inversions. In the current version (See line 301-307), we have classified them into four  
208 different groups: Inversions segregating in *O. sativa* (S), Inversions segregating in both *O.*  
209 *sativa* and *O. rufipogon* (SR), *O. rufipogon* specific (R), and *O. punctata* specific or AA-fixed  
210 (i.e., ancestral state not clear) (P) (Figure 4).

211

212 **Question 10:** - Line 216 – if I followed correctly, the 3K-RGP is a short-read dataset. It's not  
213 clear to me how that could be used to test/confirm inversions, particularly since later the  
214 claim is made that short read data are not useful for inferring inversions after comparing the

215 results of this paper to the Fuertes paper. Moreover, details of methods, numbers used, ‘high  
216 coverage’, etc., are lacking, such that is hard to follow the basis for conclusions.

217

218 **Author response:** We agree with reviewer’s comments. We compared the inversions from  
219 genome assembly and short reads by using overlapping genomes. In doing so, we discovered  
220 a very high frequency of false positives when using short read data (Fuertes et al., 2019) (See  
221 line 373-378, [Supplementary Table 14](#)). In our case, we did not use short reads to call  
222 inversion directly. We mapped the 3K-RGP Illumina reads to the reference genome. We took  
223 a detailed look into the alignment in the genome browser (IGV), to validate the mapping of  
224 short reads that span the previously confirmed inversions to show a clear breakpoint (See line  
225 317 – 330) between samples. We initially did this manually but now also applied a machine  
226 learning approach to identify inversion events across the whole 3K-RGP data set. We added  
227 this analysis with details into the manuscript ([Supplementary Note 4](#)).

228

229

230 **Question 11:** - Line 256 - TE information about breakpoints is interesting. By “analyzing TE  
231 content across the inversion index”, is this all 1054 inversions?

232

233 **Author response:** Thanks for your question. Yes, it was for all 1,054 inversions. In this  
234 update, we applied our analysis to the full 75 genome data set and identified 1,769 inversions.

235

236 **Question 12:** - Figure 3C. I’m confused and need some explanation in the legend. The lines in  
237 the middle seem to be inverted, but the arrows go in the same direction. The TE (0025) seems  
238 to be an LTR that is split in both cases. What is INT?

239

240 **Author response:** We apologize for the confusion. We have corrected the second arrow,  
241 which shows the direction in accession Minghui 63. The transposable element (ID: Os0025)  
242 is a long terminal repeat (LTR) retrotransposon, which include two parts that were shown in  
243 Fig 5C (the 3C in previous version), i.e. long-terminal repeats (Os0025\_LTR) and integrase  
244 gene (Os0025\_INT), which is located in one of part in internal domain region of the  
245 retrotransposon. We updated this explanation in the Figure 5c legend.

246

247 INT refers to integrase gene in Internal Domain, please see the structure of LTRs that cited  
248 from Alzohairy et al., 2014 (<https://doi.org/10.1071/FP13339>)

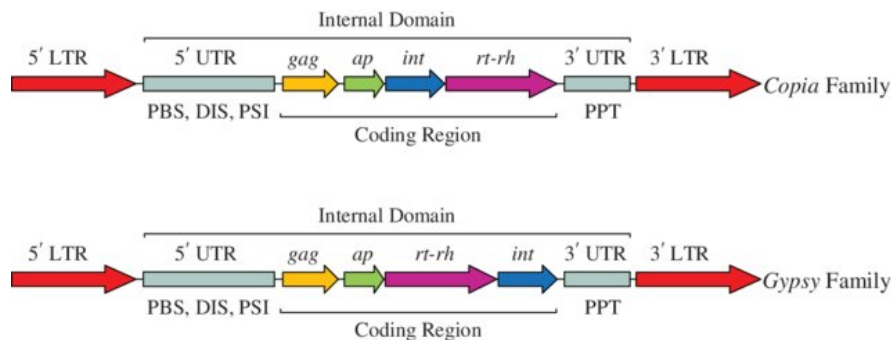


Fig. 1. Schematic structure differences between long-terminal repeat (LTR) retrotransposons (RTs) of Copia and Gypsy families. 5' gag, group-specific antigen or capsid protein gene; ap, aspartic protease gene; int, integrase gene; rt, reverse transcriptase gene; rh, ribonuclease-H gene; 3' UTR, 3' untranslated region; PBS, primer binding site; DIS, dimerisation signal; PSI, packaging signal; PPT, polypurine tract.

249  
250  
251  
252  
253  
254  
255

256 **Question 13:** - How many inversions were validated with bioNano, as in Figure 2? Overall,  
257 I'm not convinced that all 1054 inversions were independently validated (as claimed in point  
258 #2 in the Intro). I do believe bioNano, but details are lacking (how low is the resolution? How  
259 many could be confirmed?) I believe that short-reads are useful, but details are lacking here,  
260 too, but (again) it's hard to claim they validate inversions on the one hand but are not useful  
261 at all on the other (e.g., Fuertes).

262

263 **Author response:** We apologize for the confusion. We did not validate all inversions based  
264 on bionano, but only for inversions greater than 1 Mb and corrected our claims in the  
265 introduction. Of note, we do have Bionano data for 12 genomes in the 18-genome data  
266 package.

267 We clarified how we applied the Illumina reads for cross checking in [Supplementary Fig.](#)  
268 [3](#). Compared to Fuertes's work, we did not call inversions based on short read mapping, but  
269 only applied the alignment patterns on the inversion breakpoints. Fuertes's work, used short  
270 read data to call INVs, which is the biggest difference between our two studies.

271

272 **Question 14:** - Lines 285 – I'm assuming that 10.9% and 7.3% is much higher than the  
273 genome average, but it'd be nice to have an explicit comparison to the genome average to  
274 drive this point home.

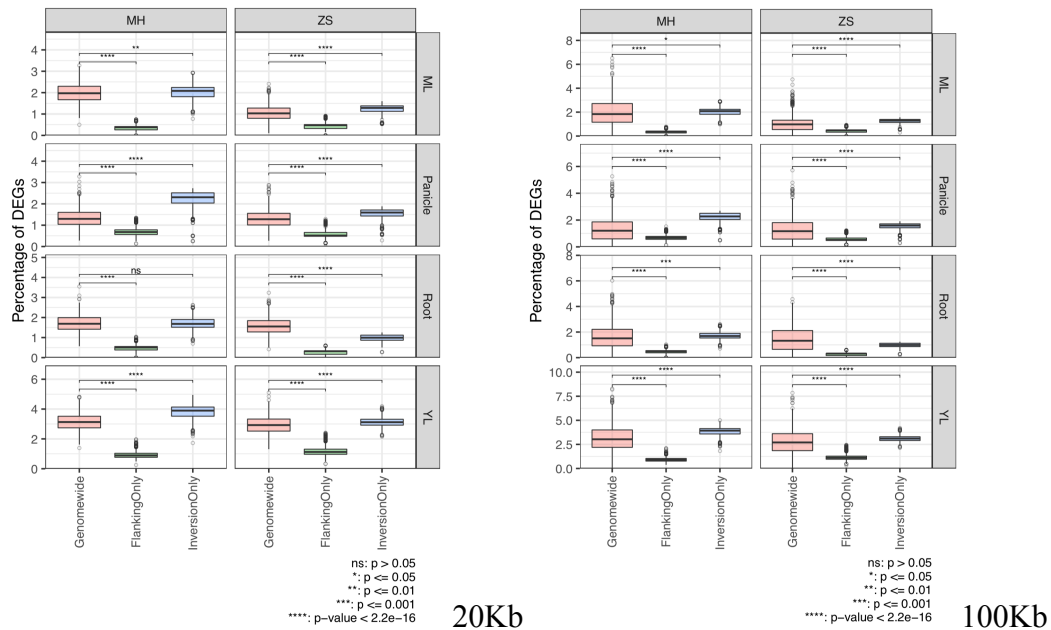
275

276 **Author response:** Yes, we agree with the reviewer. Following your suggestion, we  
277 performed the DEG analysis in 4 independent tissues (panicle, root, young leaf and mature  
278 leaf), and compared MH63 and ZS97 to the IRGSP RefSeq. We carried out a resampling test  
279 (n=1000) by comparing the portion of DEGs in genome-wide regions (20 Kb and 100 Kb



280 window), within inversion regions, and flanking regions (20 Kb). From the results, we  
 281 observed that the portion of DEGs within inversions are significantly higher than that found  
 282 in genome-wide regions.

283



**Question 15:** - About the inversion rate estimates in the Discussion (line 411 and following).

It may be that I'm misreading things, but I think they are generally wrong and perhaps horribly so. Here's why. There may be 22 post-inversion events in rice, but (if one thinks in phylogenetic terms) there are many many more years accumulated across the rice lineages on the tree than 10,000. As a brief example, let's assume (for simplicity) that 17 rice genomes diverged 10,000 years ago. If that were true, then the numerator in the rate calculation should be 17\*10,000 years, not 10,000 years, so that estimate is inflated about 17-fold. Of course, we don't know exactly when each of the separate rice genomes diverged from one another, so the estimate of 17\*10,000 years is too many. But hopefully the point is made that the calculations reported in this section may be way off and that the problem may require some consideration of population genetics given the sample.

**Author response:** We agree that this part of discussion was lacking clarity and regret that it was interpreted not as we intended. Indeed, if the estimate in question was based on 17 rice genomes, it would be wrong. What we did was based on only two genomes (Nipponbare (IRGSP) and KetanNangka (KN), both *Geng/Japonica (GJ)* genomes), and the number presented (22) is the number of inversions between these two genomes excluding any

303 inversion that could be found in any other from 15 genomes, *i.e.*, we are only using  
304 inversions private to the two genomes used. This is to avoid counting inversions in regions  
305 that could be introgressions from *Xian/Indica (XI)*, as the time to the most recent common  
306 ancestor (MRCA) between *Xian/Indica* and *Geng/Japonica* goes much farther than 10k years.  
307 This approach was in fact over-cautious because it is reasonable that *KetanNangka* would  
308 share some inversions with other japonica genomes.

309 Our inversion rate estimates were calculated as the (number of inversions between 2  
310 genomes) / (2\*estimated time to coalescence *i.e.* 10kya). Since this *O. sativa*-only estimate  
311 was much higher than the two preceding estimates, in the revised version we are giving a  
312 more detailed analysis, by

- 313 1) using a larger collection of *GJ* genomes (14) to base our estimates on a larger set of  
314 comparisons (out of 75 genomes),
- 315 2) measuring divergence in the inversion regions using SNP data in order to test whether  
316 any of the remaining inversions could have been introgressed from *XI* populations in the past  
317 but did not make it into the representative genomes.
- 318 3) Adjusting divergence time used based on recent work (Gutaker et al., 2020) and  
319 theoretical considerations, using a more conservative estimate of 14,200 years to allow time  
320 for coalescence within the ancestral population.

321

322 This leads to a more conservative estimate, which is closer to the one based on cross-species  
323 comparison, however is still higher than cross-species one. We want to note however, that  
324 estimates based on inter-species comparisons are likely serious underestimates as detecting  
325 inversions becomes more difficult with higher divergence - which can be already seen  
326 between diverged populations of *O. sativa*.

327 Additionally, under larger divergence times, there may be under-counted due to the  
328 recurrence of inversions (inversion can occur in same region). In a recent study in human  
329 (Porubsky et al., 2022), the authors estimated a recurrence rate of 2.7 inversions per 10k  
330 generations per locus (at certain loci), which is slightly less than half of our genome-wide  
331 rate.

332 We thank the reviewer for making this point and have updated the methods, results and  
333 discussion accordingly (See line 269 - 298).

334

- 335 1. Porubsky D, Höps W, Ashraf H, et al. Recurrent inversion polymorphisms in humans  
336 associate with genetic instability and genomic disorders[J]. Cell, 2022, 185(11): 1986-2005.  
337 e26.
- 338 2. Gutaker R M, Groen S C, Bellis E S, et al. Genomic history and ecology of the  
339 geographic spread of rice[J]. Nature plants, 2020, 6(5): 492-502.

340

341

342 **Question 16:** - Line 432 – the comparison to Fuertes. I'd have to read that paper carefully to  
343 see how it was done, but the title implies there were 3,000 individuals. That suggest, I think,  
344 that the entire dataset was used as evidence to support or not the inference of inversions. I'd  
345 expect them to have found many more inversions, but that the inversions in this paper would  
346 be a subset of their total set. I think that'd be a more fair comparison. That said, it is indeed  
347 puzzling that there is only 194 out of 1054 that overlapped! While I agree with the authors  
348 that short read data is certainly less accurate than long-read data, it does make one wonder  
349 about the accuracy of the 18 assemblies and whether there is no only errors in Fuertes but  
350 also in assemblies that mislead inversion inference.

351

352 **Author response:** Indeed, the 3K-RGP study (Fuertes et al., 2019) listed 1,255,033  
353 inversions, but they were identified along with other structure variations, and were listed  
354 separately for each accession, i.e., with redundancy. Due to the limitation of detecting  
355 inversions with short-reads, the inversions were false positive errors (Type I error) and were  
356 not validated. Taking advantage of the overlaps with our dataset, we determined that about  
357 90% of inversion reported by Fuentes et al. (2019) could be false. In our study, we used  
358 dotplots to validate and long reads for correcting the inversion identification. Of note, we also  
359 found that even with whole a genome alignment strategy, that raw inversions from any caller  
360 should be validated to obtain precise inversions (~75%).

361 Following your suggestion, we compared our 1,769 inversions with the all inversions  
362 (entire dataset) that generated by 3K-RGP (Fuentes et al. 2019), and we found 293 were  
363 overlapped with previous report. The details are reported in [Supplementary Table 4](#).

364

365 **Question 17:** - One time consuming but convincing way to validate inversions is to find long  
366 reads that span the inferred junctions. I don't think that was done her, but it would certainly  
367 go a long way to confirming the inversion inferences more convincingly.

368 **Author response:** Thanks for your comments. Following the reviewer's suggestion, we  
369 validated a subset of 264 random inversion with long-read data and found that 97.73% of the  
370 inversions could be supported at both the left and right breakpoints in the reference and  
371 queried genomes. We updated this in the manuscript (See line 197-201) and with the  
372 following table ([Supplementary Table 5](#)).

373

Genome	Number of INVs	Number of INVs could be supported by PacBio LongReads
LaMu	78	75
NaBo	74	74
CMeo	35	35
MH63	77	74
Total	264	258 (97.73%)

374

375

376 **Reviewer #3**

377 Zhou et al. present a pan-genome analysis of the major sub-populations of asian rice and two  
378 wild rice species represented by a set of 18 whole genome assemblies. the two wild rice  
379 assemblies are new and are provided with this manuscript. The other assemblies have been  
380 published before.

381

382 **Question 1:** The authors main point of analysis is the cataloging of inversions larger than 100  
383 bp and contextualizing these with data on recombination, LD, selection and gene expression.  
384 As the authors point out, the detection of inversions (also larger) is not novel for plant  
385 genomes (maybe more references to recent pan-genome studies in cereals would be justified  
386 here?), however, here maybe a first comprehensive catalog for representatives of a crop  
387 species' subpopulations is provided.

388

389 **Author response:** Yes, mostly, the inversions were seldom studied in detail. Our study is a  
390 first comprehensive identification, validation, and detailed study in crop species at  
391 subpopulation level. We also added more references on past studies of INVs and their  
392 downstream effects to comprehensively survey the literature and results of previous work, in  
393 rice, cereals and beyond.

394

395 **Question 2:** The study reads well and is well presented. The data displays sometimes are  
396 rather basic - often just direct exports out of commercial analysis and vizualisation software?  
397 Intuitivity of the figure displays could be improved - maybe?

398

399 **Author response:** Thanks for the reviewer's support and suggestions. We have improved our  
400 image designs, and we limited visualization outputted commercial software. The only  
401 commercial software used are ones that do not have an equivalent open-source visualization  
402 *i.e.*, bionano maps.

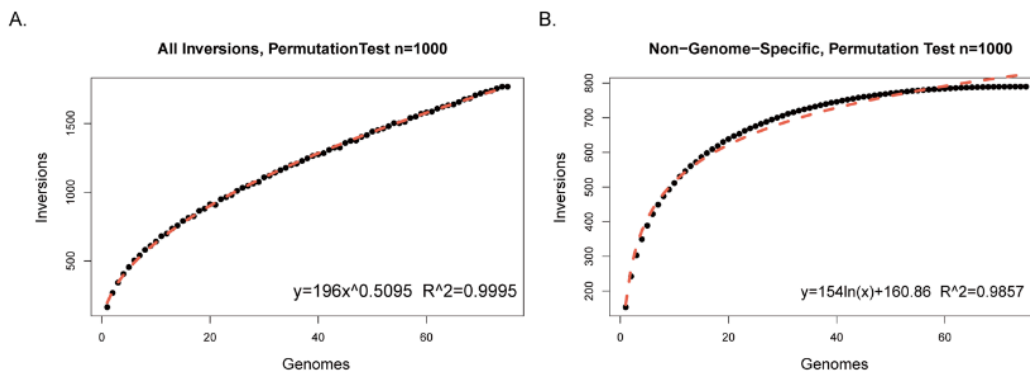
403

404 **Question 3:** The study adresses an obvious point in comparative genomics as more genomes  
405 of the same species become available. The authors are well aware of artifacts introduced into  
406 analyses when different assembly qualities affect the analysed assemblies. Therefore, I was  
407 surprised, given the today's costs for making a Hifi assembly for rice (haploid assembly  
408 consumable costs in the few hundred dollars range!) that the authors did not make an effort  
409 here to have really absolutely comparable datasets. This may only incrementally change the

410 presented results, however, it is a weakness of the study that could have been avoided with  
411 modest effort and investment. The same argument applies to the annotation which was done  
412 with the same pipeline only for the Asian rice.

413

414 **Author response:** Thanks for these great comments. Yes, indeed, Hifi assemblies have  
415 greatly enhanced the number and quality of genome for rice and many species. As a genome  
416 size with 400 Mb, rice could be the model crop for those species as well, which can also be  
417 cost effective. Luckily, the rice community has contributed a large number of assembled  
418 genomes during the last two years, which helped us to obtain a comparable dataset. In this  
419 case, we collected all publicly available high-quality genomes and performed a permutation  
420 simulation (n = 1000) to determine the number of genomes we needed to interrogate to reach  
421 near inversion saturation. In doing so, we found that 60 high-quality genomes would be  
422 required to have a 99% chance of identify the majority inversions with allele frequencies of 2  
423 or greater (Figure 1).



424

425

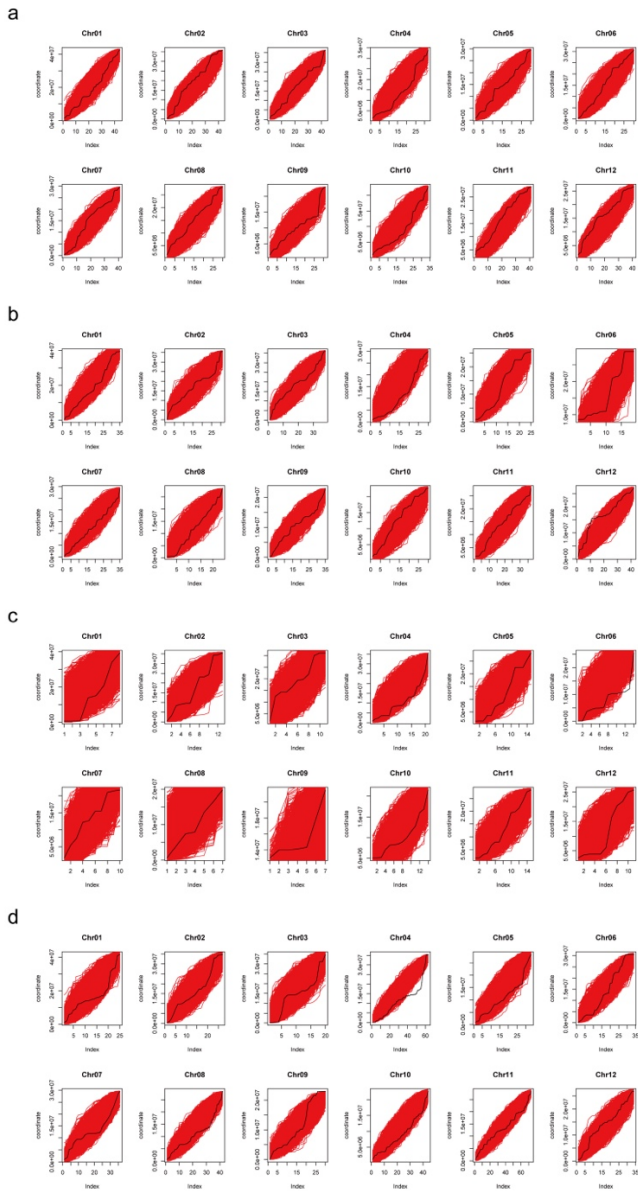
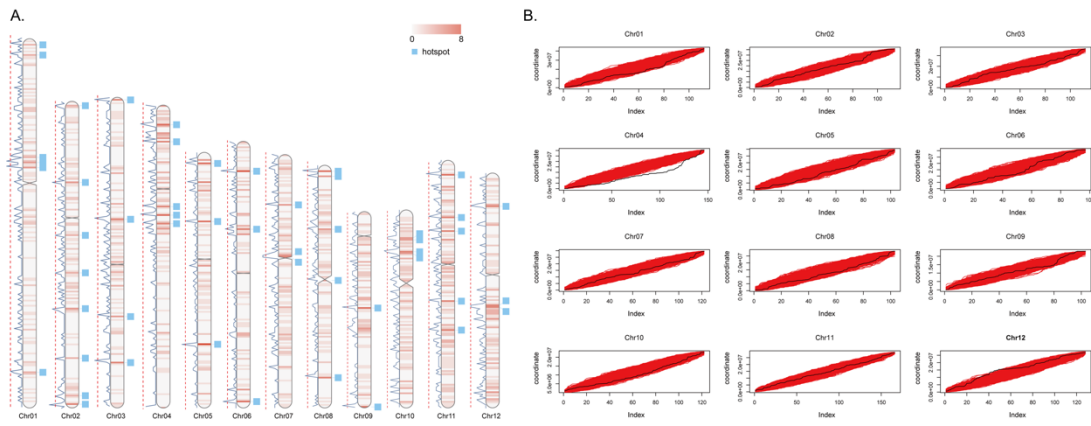
426 **Question 4:** The spectrum of inversion distribution: the authors say they are randomly  
427 distributed - Ext. figure 3 should support this claim. While the figure legend is not conclusive  
428 and the figure itself is basic, I think the authors should have made an effort define size classes  
429 and redo the genome distribution scan. Is it true that all sizes of inversions are evenly  
430 distributed along the chromosomes? This is counterintuitive.

431

432 **Author response:** We have re-analyzed the spectrum of inversion distributions and adjusted  
433 the manuscript accordingly (See line 246-252). We found that inversions were evenly  
434 distributed except on all rice chromosomes (n=12) with the exception of chromosome 4  
435 (Figure 2). Following your suggestion, we split divided the inversions into 4 groups, i.e. <1

436 Kb, 1-5 Kb, 5-10 Kb, and >10 Kb, and found that only >10 Kb inversions on chromosome 4  
 437 contributed to uneven distributed (Supplementary Fig. 2).

438



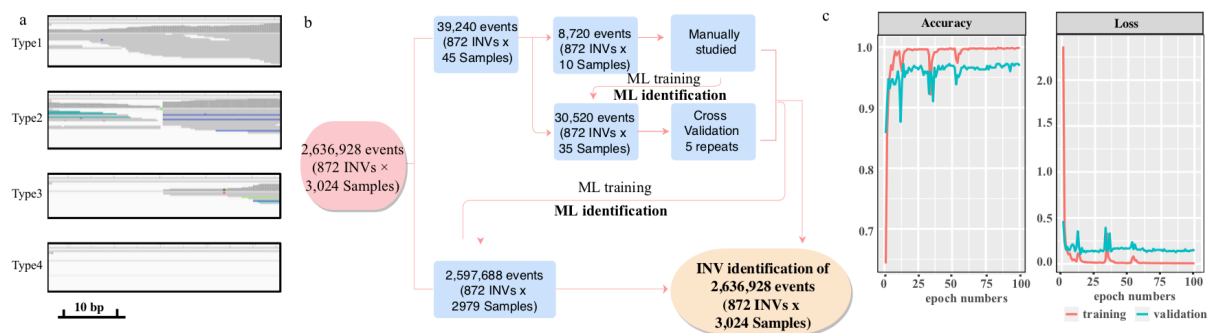
439

440 a. < 1 Kb, b. 1 – 5 Kb, c. 5 – 10 Kb, d. > 10 Kb.

441 **Question 5:** The authors used the 3K rice data to genotype for the presence of inversions,  
 442 however, they performed this only on a selected set of genotypes representing the 15 sub-  
 443 populations. Is it really computationally so intense to include all 3000? Or is this a problem  
 444 of sequence coverage in a certain proportion of the 3K dataset? I would love to see the  
 445 analysis on all 3000! I wonder whether the authors tried to use the inversion catalog as a  
 446 proxy to model and detect additional inversions in the full 3K dataset as the 15 genotypes  
 447 sequenced will for sure not give the full pan-genome inversion spectrum. The authors outline  
 448 the technical feasibility in the discussion. Furthermore, it would be very instructive to give an  
 449 estimate for WGS coverage to robustly scan for the presence of inversions at population  
 450 scale?

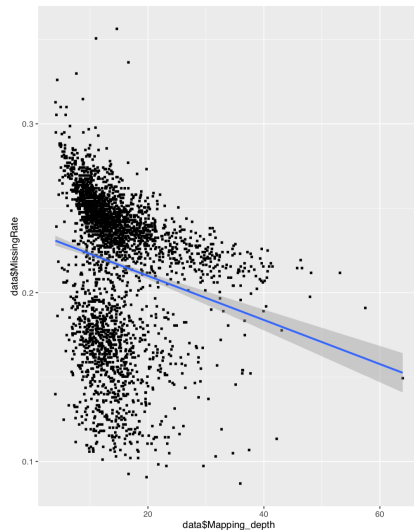
451  
 452 **Author response:** We agree with reviewer, and yes it would be nice if we could scan the  
 453 entire 3K-RGP dataset. As discussed above, when compared the inversion results that used  
 454 both short read and long read data, 90% of the short-read inversion call appeared to be false  
 455 positives. However, we also found that the alignment in IGV could hint at the presence of  
 456 inversions if we know the coordinates.

457 Since the reviewer suggested that it would be very instructive to give an estimate for  
 458 coverage to robustly scan for the presence of inversions at population scale, and to see the  
 459 analysis on all 3K-RGP samples, we applied a machine learning approach. To do this, we  
 460 started by manually curating 872 inversions (*O. sativa* specific) across 45 accessions  
 461 (overlapped with short reads and long reads), and used these 39,240 inversion events to train  
 462 a machine learning model (online method). This model was used to study the 872 inversions  
 463 across the remaining 2,979 samples, with 5-fold cross validation. In doing so, we were able to  
 464 assess the presence or absence of all 872 inversion at the 3K-RGP population level  
 465 (Supplementary Note 4).



466  
 467 From the estimation, we found only slightly negative correlate ( $df = 3020$ ,  $p$ -value  $< 2.2e$ -  
 468 16,  $cor = -0.1955148$ ) with WGS sequence coverage and inversion validation at the  
 469 population level (as below).





470

471

472 **Question 6:** TE landscape was analysed in context of inversions and it was detected that  
473 inversions are enriched with TE content in rice. The analysis of inversions could go  
474 further here. Have you systematically assessed the TE at the Inversion junctions? This  
475 could reveal patterns of the mechanistic involvement for the occurrence of inversions in  
476 the rice genome.

477

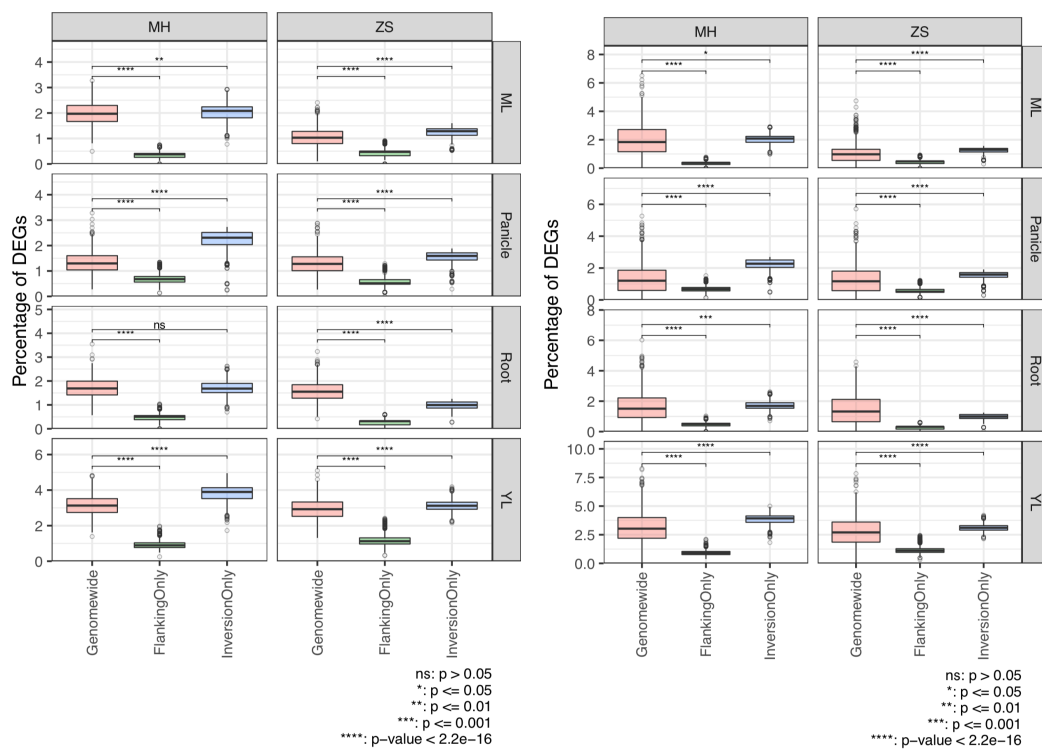
478 **Author response:** We agree with reviewer. We updated the TE related analysis, and add the  
479 results and discussion of patterns of mechanisms for the occurrence of inversions (See  
480 line 367-378 and [Supplementary Table 10](#)).

481

482 **Question 7:** Gene expression and inversion: the authors showed examples where inversions  
483 and expression levels of genes in inverted regions and their orthologs in "non-inverted"  
484 haplotypes. The data is interpreted as inversions are causal, which is not unlikely,  
485 however, the authors do not show any functional validation (which is not trivial) and  
486 analysis and they also do not provide complementary datasets (methylation, HiC,  
487 ATAC etc.) which would support their hypothesis and interpretations. I would also  
488 recommend to expand the analysis to genes adjacent or in neighborhood but not  
489 directly affected by the inversions - regulatory sequence context can and will be  
490 affected also for such genes and one should see a gradient perhaps?

491 **Author response:** Thanks for these comments. We agree with the reviewer that that  
492 functional validation could support the hypotheses. In this paper, we are aiming to show clue  
493 that inversions could affect gene expression.

494 Following your suggestion, we explored the gene expression neighborhoods, *i.e.*, flanking  
 495 regions of inversions. We studied the portion of DEGs in flanking regions (as shown below,  
 496 within inversions, and genome-wide regions (See the two charts below showing gene  
 497 expression patters in 20 Kb (left) and 100 Kb (right) flanking regions. To compare the  
 498 portion of DEGs in these regions, we performed a permutation test (n = 1000) in four tissues  
 499 (panicle, mature leaf, young leaf and root). From the charts below, we observed that the  
 500 portion of DEGs within inversions was scientifically higher than genome-wide regions, and  
 501 the portion of DEGs in flanking regions was significantly lower than in genome-wide  
 502 regions. These results demonstrate that gene expression in flanking regions of inversions  
 503 could be affected directly or indirectly by inversions themselves, as the reviewer commented.  
 504



505  
 506  
 507 **Question 8:** Similarly, the interpretation of positive selection in inverted or non-inverted  
 508 haplotypes should maybe presented and discussed with more caution.

509  
 510 **Author response:** Since we removed the analysis of cluster inversion 92, we deleted this part  
 511 of the manuscript.

512  
 513 **Question 9:** The authors make a point in the discussion that not merging sequences into a  
 514 pan-genome graph was an advantage here. Well isn't this obvious? And isn't the need for

515 graph development depending on the analytical goals - I am not sure if I can follow the  
516 argument here.

517

518 **Author response:** Sorry for the confusion here. Yes, this is obvious. Since the pangenome  
519 idea is one of the most popular research topics nowadays, and there are a lot of tools or  
520 approaches that have been developed for building pangenomes. However, there are no  
521 computation tools that cover all structure variations, especially inversions. If we want retain  
522 all genetic diversity, we recommend the avoidance of computational tools that collapse  
523 genomes for building pangenomes.

524

525 **Question 10:** data is all available to what I could track, however, one has to dive into  
526 supplements. A clear data availability statement with instructions where to find the details is  
527 missing and needs to be added.

528

529 **Author response:** Thanks for this comment. We made a website to describe the motivation  
530 of the project and the related datasets. Please see the following link which was updated in  
531 the manuscript <https://yongzhou2019.github.io/Rice-Population-Reference-Panel/>. We hope  
532 this will help readers to easily track all available data used in this study.

533

534 **Question 11:** persephone visualization is challenging. Maybe it is comprehensive but I don't  
535 think it is necessarily intuitive. Have you tried the recent development from John  
536 Lovell at Hudson Alpha (GENESPACE - look at BiorXiv)?

537

538 **Author response:** Thanks for this comment. We tried GENESPACE, and it is a great tool  
539 indeed, especially for orthologs and synteny analysis. Normally we do static analysis of  
540 syntenic orthogroups in GENESPACE in most analyses. However, Persephone is a multi-  
541 genome browser, and this is dynamic but not static like GENESPACE. Of note, we also use  
542 Persephone as a platform to share information with the community who is interested in our  
543 data. In this case, we uploaded our 18-genome data package in Persephone including the  
544 genomes, gene annotations, TEs and structure variations, which are all available through this  
545 link <https://web.persephonesoft.com/>.

546

547 **Question 12:** I appreciate the authors use ref 46 for wheat in wheat as this is part of their own  
548 work, however, in the context still the IWGSC 2018 ref in Science is probably more

549 appropriate? - unless you want to make a point out of impact of Hifi sequencing in wheat -  
550 which is not the case at current.

551

552 **Author response:** We added two papers that reported by IWGSC who reported the large  
553 genome of wheat as well, please see below:

554

555 47. Consortium, I.W.G.S. et al. Shifting the limits in wheat research and breeding using a  
556 fully annotated reference genome. 361, eaar7191 (2018).

557 48. Consortium, I.W.G.S. et al. A chromosome-based draft sequence of the hexaploid  
558 bread wheat (*Triticum aestivum*) genome. 345, 1251788 (2014).

559

560 **Question 13:** abstract: "effects on gene regulation" - no! you only report correlations between  
561 datasets, no functional proof.

562

563 **Author response:** We agree with the comments and changed to gene expression.

564

565 **Question 14:** intro: "almost" 10 billion by (exactly) 2064 - I find this mix of approximation  
566 and exactness curious - sure you are citing here, but ...

567

568 **Author response:** We have modified the text to "Since the world population is expected to  
569 increase to approximately 10-billion by 2060-2070"

570

571 **Question 15:** intro: why sequence diversity is a natural variation "tool box"??

572

573 **Author response:** We deleted this sentence.

## Reviewers' Comments:

### Reviewer #1:

#### Remarks to the Author:

A super pan-genomic landscape of rice had been published at 12 July 2022 (Shang et al., Cell Research, 2022), leading that the novelty of this manuscript is greatly decreased.

I noted that the description regarding the biological significance (in terms of gene expression) of inversion in manuscript was different with that in rebuttal document, only the gene expression within inversion was described in manuscript, whereas, the effects on gene within, nearby, at genomic scale was described in rebuttal document. This mistake should not happened.

In the section of "Characterization of gene content within inversions and their breakpoints", the authors used expression datasets of only two indica accessions to try to show the effect of inversion on gene expression. In my view, two samples were few for the effect analysis on gene expression of inversion. Additionally, the authors detected only ~10 genes within inversion were differentially expressed in each tissue between two accessions, given that the effect of other variations on gene expression were not excluded, therefore, the inversion effect on gene expression needs to reconsidered.

### Reviewer #3:

#### Remarks to the Author:

The manuscript has been greatly improved by addressing all reviewer comments which required to include more data and re-run a substantial number of analyses. I am happy with the changes and have no further requests. Congratulations to a comprehensive and informative piece of rice genomics research.

### Reviewer #4:

#### Remarks to the Author:

The authors generated two de novo genomes for two wild rice species: *O. rufipogon* and *O. punctata*, respectively, and analyzed the data with ~70 previously published genomes to study inversion in Asian rice. I think the manuscript is clearly written. However, I don't think the manuscript reached the novelty requirements of NC. I also have the following major and minor comments for the authors.

#### Major:

First, the rice pan-genome has been assembled in at least five studies started from short or long reads with/without the use of outgroups. To be honest, I didn't get why the authors assembled another pangenome with a similar design.

At the same time, the inversions have been investigated using short reads, long reads, and assemblies. I didn't see many new insights compared with previous studies.

For the two genomes the authors generated for wild rice, I guess the wild samples have high heterozygosity. In this case, diploid assemblies are expected.

The presentation of the figures is unclear, sometimes, hard to follow.

#### Minor:

L84-88: I didn't get this. Did the authors mean SNPs are SVs?

L91: Almost all the rice pangenome papers have analyzed the inversions.

L94-101: There are lots of good examples in plants too.

L122-124: How about the genome research and cell research papers?

L140-144: I am sure this is not right. The bias in sampling leads to such biased results since rufipogon has much higher genetic diversity.

L149-151: The authors should be very careful about this estimation, more outgroups will definitely increase this estimate.

L152-157: This has been indicated in previous studies.

L158-159: This is well-known information.

L260-267: Phylogenetic analyses are expected to be done using neutral markers, I believe most of the INVs are under selection.

L269-298: The author should be very careful with such estimations, see my comments above.

L380-406: If I understand this correctly, this is novel to some extent.

Discussion: the authors should compare their results with previous publications in detail. It is hard to get what is new.

1 Point-by-point response

## 2 REVIEWERS' COMMENTS

3

4 **Reviewer #1** (Remarks to the Author):

5

6 **Question 1)** A super pan-genomic landscape of rice had been published at 12 July 2022  
7 (Shang et al., Cell Research, 2022), leading that the novelty of this manuscript is greatly  
8 decreased.

9

10 **Author response:** We address this question in the response above. In addition, the question  
11 of novelty was overruled by the Nature editors so we will not discuss the issue of novelty  
12 further.

13

14 **Question 2)** I noted that the description regarding the biological significance (in terms of  
15 **gene expression**) of inversion in manuscript was different with that in rebuttal document,  
16 only the gene expression within inversion was described in manuscript, whereas, the effects  
17 on gene within, nearby, at genomic scale was described in rebuttal document. This mistake  
18 should not happened.

19 **Author response:** We apologize for the confusion. We added the results of gene expression  
20 within inversions, flanking regions, and genomic scale (from Line 375-382) as below:

21 “Based on a comparison of transcript abundance levels between Nipponbare and MH63 and  
22 ZS97 (across four tissue types - root, panicle, young leaf, and mature leaf) we detected that 5  
23 - 12 genes from MH63 and 4 - 11 genes from ZS97 within inversions, 2 - 7 genes from  
24 MH63 and 2 - 4 genes from ZS97 in inverted flanking regions, and 19 - 42 genes from MH63  
25 and 9 - 30 genes from ZS97 that were located in non-inverted randomly resampled 20Kb  
26 regions, were differentially expressed (DEG, fold change > 2, *p* value < 0.01)

27 (Supplementary Fig. 7).”

28

29 **Question 3)** In the section of “Characterization of gene content within inversions and their  
30 breakpoints”, the authors used expression datasets of only two indica accessions to try to  
31 show the effect of inversion on gene expression. In my view, two samples were few for the  
32 effect analysis on gene expression of inversion. Additionally, the authors detected only ~10

33 genes within inversion were differentially expressed in each tissue between two accessions,  
34 given that the effect of other variations on gene expression were not excluded, therefore, the  
35 inversion effect on gene expression needs to be reconsidered.

36 **Author response:** Thanks for the comments. We do agree with the reviewer's view. For  
37 gene expression, we used data from the IRGSP RefSeq (*GJ/japonica*) and two *XI/Indica*  
38 genomes (MH63 and ZS97), and fortunately, the RNA was deep sequenced with multiple  
39 replicates. Regarding the differential transcript abundance, we do agree with the reviewer that  
40 other variations, e.g., SNPs and PAVs have effects on genes' expression and addressed this  
41 concern by stating "the possible effect of inversions" (Line 373).

42

43 **Reviewer #3** (Remarks to the Author):

44 The manuscript has been greatly improved by addressing all reviewer comments which  
45 required to include more data and re-run a substantial number of analyses. I am happy with the  
46 changes and have no further requests. Congratulations to a comprehensive and informative  
47 piece of rice genomics research.

48 **Author response:** Thank you very much for your previous advice and support.

49

50 **Reviewer #4** (Remarks to the Author):

51 The authors generated two de novo genomes for two wild rice species: *O. rufipogon* and *O.*  
52 *punctata*, respectively, and analyzed the data with ~70 previously published genomes to study  
53 inversion in Asian rice. I think the manuscript is clearly written. However, I don't think the  
54 manuscript reached the novelty requirements of NC. I also have the following major and  
55 minor comments for the authors.

56 **Author response:** Thanks for the comments. The major and minor comments are really  
57 helpful. We hope we have addressed all your questions and suggestions. As for novelty, the  
58 question was overruled by the Nature editors so we will not discuss this concern.

59

60 Major:

61 **Question 1)** First, the rice pan-genome has been assembled in at least five studies started  
62 from short or long reads with/without the use of outgroups. To be honest, I didn't get why the  
63 authors assembled another pangenome with a similar design.



64 **Author response:** As discussed in the revised discussion, a total of 5 Asian rice pan-genome  
65 studies have been published to date. Two used Illumina data and 3 PacBio, ONT or both and  
66 only 3 studies called inversions (Table 3, as shown above). The one Illumina study showed a  
67 91.5% false discovery rate and was thus discounted almost entirely. The Qin et al., 2021  
68 paper sequenced 33 high quality genomes and used the same strategy we did to call  
69 inversions. However, this paper did not interrogate the full genetic diversity of Asian rice,  
70 encompassing only 9 subpopulations, and did not validate any of their inversions. The Shang  
71 et al., 2022 paper generated 174 ONT genomes, however, all assemblies were reported as  
72 contigs and not chromosome level scaffolds which can lead to serious errors in the evaluation  
73 of inversions (please see our response to the editor above). In addition, the paper did report  
74 the presence of 2,784 raw inversions, however no validations were reported, and, more  
75 importantly, this data has not been made available (please see our response to the editor  
76 above).

77 Regarding the last comment about design – of all the pan-genome studies, our design is the  
78 most unique to all 5 in that we made no attempts to computationally collapse our 73-genome  
79 data into a graph and kept all assemblies in their native state. This permitted us to precisely  
80 identify inversions and inversion boundaries, as well as assess their frequency at the  
81 population level using machine learning.

82

83 **Question 2)** At the same time, the inversions have been investigated using short reads, long  
84 reads, and assemblies. I didn't see many new insights compared with previous studies.

85 **Author response:** Thanks for your comments.

86 Please see answer to question 1.

87 Further, for all previous pangenome reports in Asian rice, there were no attempts made to:

- 88 1. Validate all inversions at the individual genome and population level
- 89 2. Investigate the effects of inversions on recombination, LD and gene expression
- 90 3. Calculate inversion rates at the species level - i.e., *O. sativa* vs, *O. rufipogon*, vs. *O.*  
91 *punctata*.

92

93 **Question 3)** For the two genomes the authors generated for wild rice, I guess the wild  
94 samples have high heterozygosity. In this case, diploid assemblies are expected.

95

96 **Author response:** Thanks for the comments. In the case of both *O. rufipogon* and *O.*  
97 *punctata*, while they are wild species, they are predominantly selfing. Out-crossing rates for  
98 the two species are on the order of ~25%. Additionally, the *O. rufipogon* accession IRGC  
99 106523 was acquired 1991-07-30 and the *O. punctata* accession IRGC 105690 on 1987-04-  
100 28. Both were maintained over a number of generations since acquisition by the IRRI  
101 genebank obtaining seed from bagged panicles. Seed from these cycles of regeneration would  
102 be expected to be selfed and lead to more fixation within the accession. Further, DNA was  
103 obtained from a single plant of each for production of the genome builds.

104

105 **Question 3)** The presentation of the figures is unclear, sometimes, hard to follow.

106

107 **Author response:** Thanks for the comments. We shared and discussed all figures with  
108 authors and colleagues again, and made edits based on their suggestions. We hope the figures  
109 are now easier to follow.

110

111 Minor:

112

113 **Question 4)** L84-88: I didn't get this. Did the authors mean SNPs are SVs?

114

115 **Author response:** Thanks for the point. We have modified this in the manuscript (from Line  
116 84-88), and consider both SNPs and structure variations (SVs, i.e., INs/DELs, TRAs, and  
117 INVs) as "standing natural variation", as below:

118

119 “One source of the raw material required to meet this urgent demand is the standing natural  
120 variation that exists in the genomes of the more than 500,000 accessions of rice and its wild  
121 relatives deposited in germplasm banks around the world<sup>2</sup>, i.e., single nucleotide  
122 polymorphisms [SNPs], insertions/deletions [INs/DELs], translocations [TRAs], and  
123 inversions [INVs].”

124

125 **Question 5)** L91: Almost all the rice pangenome papers have analyzed the inversions.

126 **Author response:** Thanks for the comment.

127

128 Please see answer to Question 2

129

130 There are three recent papers that have reported the inversions, i.e., Fuentes et al., 2019, Qin  
131 et al., 2021 and Shang et al., 2022 (Cell Research) (Table 3, as shown above). We compared  
132 with the first of two papers in detail and discuss in the discussion, however, the Shang et al.  
133 (2022) paper has not released their data (i.e., inversion data, sequences and annotations [see  
134 our reply above], and contigs [fragmented genomes] so it was impossible to assess (please  
135 see our response to the editor, point 6).

136

137 **Question 6)** L94-101: There are lots of good examples in plants too.

138

139 **Author response:** Thanks for the great comment.

140 We modified the manuscript accordingly by adding some of important studies in plants (lines  
141 101-103), as below:

142

143 “In plants, INVs have been reported to play roles in, for example - local adaptation<sup>26,27</sup>,  
144 genome-environment associations<sup>27</sup>, gene regulation<sup>26,28,29</sup>, flowering time<sup>28</sup>, seed  
145 germination<sup>28</sup>, and fruit shape<sup>29</sup>.”

146

147 **Question 7)** L122-124: How about the genome research and cell research papers?

148

149 **Author response:** Thanks for the question. The Zhang et al., 2021 (genome research paper)  
150 didn't mention any inversion studies, and only PAV related data are available. The Shang et  
151 al., 2022 (cell research paper) mentioned inversions, but the data was not released and the  
152 quality of the genomes (in contigs only) are not well suited for inversion studies (as discussed  
153 in Question 5, and our response to the editor, point 6).

154

155

156 **Question 8)** L140-144: I am sure this is not right. The bias in sampling leads to such biased  
157 results since rufipogon has much higher genetic diversity.

158

159 **Author response:** Thanks for the comment. We did not use a diversity panel of *O. rufipogon*  
160 in this case, so we did not consider the genetic diversity of *O. rufipogon* here. We reported  
161 this and that might be true based our pan-genome index result.

162

163 **Question 9)** L149-151: The authors should be very careful about this estimation, more  
164 outgroups will definitely increase this estimate.

165 **Question 10)** L269-298: The author should be very careful with such estimations, see my  
166 comments above.

167

168 **Author response:** Thanks for these comments. Since these two comments are related, so we  
169 combined them into one reply.

170 Since we had somewhat different estimates from different outgroups, it may indeed make  
171 sense to rephrase the summary more carefully, e.g., instead of using a point estimate (700  
172 inversions/MY) we modified it as 735 – 749 inversions/MY.

173 We assume that we have already used outgroups from a “wide phylogenetic range” - from  
174 one very close to *GJ-tmp* (like *GJ-trop*) to *XI/indica* to *O. rufipogon* to *O. punctata*. One  
175 minor thing in using other *O. rufipogon* references is the common introgression from *O.*  
176 *sativa* to *O. rufipogon*. At the same time, it might be difficult to use more outgroups for  
177 genome-wide comparison. Going farther beyond *O. punctata* may result in too few  
178 alignments (due to the sequence similarity) and seriously underestimate the inversion rate.

179

180 **Question 11)** L152-157: This has been indicated in previous studies.

181 **Question 12)** L158-159: This is well-known information.

182

183 **Author response:** We thank the reviewer for these comments. Since these two comments are  
184 related we combined them into a single reply.

185

186 We agree that our results confirm prior work that TEs are associated with inversions and their  
187 breakpoints. We pointed this out in the manuscript, in results Line 336, and wrote:

188

189 “Transposable elements (TEs) are known to be associated with inversions<sup>11,44</sup>”

190

191 **Question 13)** L260-267: Phylogenic analyses are expected to be done using neutral makers, I  
192 believe most of the INVs are under selection.

193

194 **Author response:** Thanks for the comment. In this case we would like to point out that  
195 inversions could also be applied for phylogenic analysis, since the result was almost identical  
196 with our previous analysis using SNPs.

197

198 **Question 14)** L380-406: If I understand this correctly, this is novel to some extent.

199

200 **Author response:** Thanks very much for your comment.

201

202 **Question 15)** Discussion: the authors should compare their results with previous publications  
203 in detail. It is hard to get what is new.

204

205 **Author response:** We apologize the confusion. We actually compared our results with  
206 previous publications, e.g., Fuentes et al., 2019 and Qin et al., 2021, which was described in  
207 lines 199-205, as below:

208

209 “Our inversion index was compared with previous studies in rice, e.g., the 3K-RGP<sup>11</sup> study  
210 and the pan-genome analysis of 33 rice genomes<sup>12</sup>. Inversions were treated as “identical” if  
211 they matched the following two criteria: 1) the inversion length difference was smaller than  
212 200 bp, and 2) the differences between the coordinates of breakpoints of inversions were  
213 smaller than 100 bp. In doing so, we found that of the 1769 nonredundant inversions  
214 identified, 38.6% were previously identified ([Supplementary Data 4](#)).”