

RESEARCH ARTICLE

Open Access



# Characterization of the transcriptional divergence between the subspecies of cultivated rice (*Oryza sativa*)

Malachy T. Campbell<sup>1,2\*</sup> , Qian Du<sup>3</sup>, Kan Liu<sup>3</sup>, Sandeep Sharma<sup>1,4</sup>, Chi Zhang<sup>3</sup> and Harkamal Walia<sup>1</sup> 

## Abstract

**Background:** Cultivated rice consists of two subspecies, *Indica* and *Japonica*, that exhibit well-characterized differences at the morphological and genetic levels. However, the differences between these subspecies at the transcriptome level remains largely unexamined. Here, we provide a comprehensive characterization of transcriptome divergence and cis-regulatory variation within rice using transcriptome data from 91 accessions from a rice diversity panel (RDP1).

**Results:** The transcriptomes of the two subspecies of rice are highly divergent. *Japonica* have significantly lower expression and genetic diversity relative to *Indica*, which is likely a consequence of a population bottleneck during *Japonica* domestication. We leveraged high-density genotypic data and transcript levels to identify cis-regulatory variants that may explain the genetic divergence between the subspecies. We identified significantly more eQTL that were specific to the *Indica* subspecies compared to *Japonica*, suggesting that the observed differences in expression and genetic variability also extends to cis-regulatory variation.

**Conclusions:** Using RNA sequencing data for 91 diverse rice accessions and high-density genotypic data, we show that the two species are highly divergent with respect to gene expression levels, as well as the genetic regulation of expression. The data generated by this study provide, to date, the largest collection of genome-wide transcriptional levels for rice, and provides a community resource to accelerate functional genomic studies in rice.

**Keywords:** RNA sequencing, *Oryza sativa*, Population genetics, Regulatory variation, Expression quantitative trait loci, Gene expression, Natural variation

## Background

Cultivated rice consists of two subspecies: *Indica* and *Japonica*. *Indica* varieties are cultivated throughout the tropics, and account for the majority of rice production worldwide. *Japonica* varieties, on the other hand, are grown in both tropical and temperate environments, and only account for approximately 20% of rice production.

Although the domestication history of rice remains a contested topic, the current research collectively suggests that rice was domesticated at least twice from two geographically and ecologically distinct subpopulations of *Oryza rufipogon*. The unique environmental pressures in these distinct regions, as well as preferences by early farmers for grain characteristics has resulted in large morphological and physiological differences between the two subspecies. These differences have been recognized for centuries, as evidenced by references of Keng and Hsein types of rice found in records from the Han Dynasty in China [1].

\*Correspondence: [campbell.malachy@gmail.com](mailto:campbell.malachy@gmail.com)

<sup>1</sup>Department of Agronomy and Horticulture, University of Nebraska Lincoln, 1825 N 38th St., 68583 Lincoln, NE, USA

<sup>2</sup>Department of Animal and Poultry Sciences, Virginia Polytechnic Institute and State University, 175 West Campus Drive, 24060 Blacksburg, VA, USA

Full list of author information is available at the end of the article



© The Author(s). 2020 **Open Access** This article 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 licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence 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 licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

The unique natural and agronomic selection pressures placed on the wild progenitors and early proto-domesticates resulted in drastic changes at the genetic level. Work by Huang et al. [2] showed considerable reduction in genetic diversity in *Indica* and *Japonica* compared with *O. rufipogon*. Such drastic reductions in genetic diversity are common following domestication. Moreover, the transition from an out-crossing/heterogamous nature of *O. rufipogon* to the autogamous breeding system of cultivated rice likely led to greater partitioning of genetic diversity among the two subspecies, and further differentiation of the two groups. These large genetic differences have been recognized for nearly a century as hybrids between *Indica* and *Japonica* exhibit low fertility [3]. More recently, these genetic differences have been realized with the availability of high density molecular markers and full genome sequences for both *Indica* and *Japonica* [2, 4–12]. For instance, Ding et al. [4] showed that approximately 10% of the genes in the *Indica* and *Japonica* genomes showed evidence of presence-absence variation or asymmetrical genomic locations. Several other studies have highlighted genetic differences between the subspecies as structural variants differences, gene acquisition and loss, transposable element insertion and single nucleotide polymorphisms [2, 5–12].

While the morphological and genetic differences of *Indica* and *Japonica* have received considerable attention, few studies have investigated the divergence between the two subspecies at transcriptome level [13–15]. Walia et al. [13] utilized genome-wide expression profiling to characterize the transcriptional responses for two *Indica* and *Japonica* cultivars to salinity. This study was performed to elucidate the mechanisms underlying the contrasting responses to stress exhibited by the cultivars, rather than examine the transcriptional difference between the subspecies. Moreover, separating genotypic differences from subspecies differences is not feasible with the low number of cultivars used in these studies. Lu et al. [14] compared transcriptional profiles of two *Indica* accessions and a single *Japonica* accessions and identified many novel transcribed regions, highlighted alternative splicing differences, and differentially expressed genes between accessions. Although these studies provided insights into the transcriptional differences between *Indica* and *Japonica*, given the small sample size, the scope for extending conclusions to a population level is limited. Jung et al. [15] leveraged the large number of public microarray databases to compare transcriptional diversity between the two subspecies. The 983 publicly available Affymetrix microarrays were classified into *Indica* and *Japonica* subspecies based on the cultivar name. This study showed that considerable differences in expression levels were evident between the two subspecies. However, large proportion of information is likely lost due to the heterogeneity in

sample types (e.g. tissue, developmental stage) and varying growth conditions. Thus, a more highly controlled study that utilized a larger panel with genotypic information would provide greater insight into the differences in expression levels, as well as provide a mechanism for connecting transcriptional differences between the two subspecies with genetic variation.

The objective of this study is to examine the genetic basis of the transcriptional variation at a population level within the *O. sativa* species. By combining population and quantitative genetics approaches, we aim to elucidate the genetic basis of transcriptional divergence between the two subspecies. To this end, we generated transcriptome data using RNA sequencing on shoot tissue for a panel of 91 diverse rice accession selected from the Rice Diversity Panel1 (RDP1) [16–18]. Here, we show that transcriptional diversity between *Indica* and *Japonica* subspecies is consistent with diversity at the genetic level. Moreover, we connect transcriptional differences between the two subspecies with divergent patterns of *cis*-regulatory variation. This study is the first to document the transcriptional divergence between the major subspecies of cultivated rice at a population level, and provides insight into the genetic mechanisms that have shaped this transcriptional divergence.

## Results

We selected 91 accessions to represent the genetic diversity within Rice Diversity Panel 1 (RPD1). Using the sub-population assignment described by Zhao et al. [16] and Famoso et al. [17], shoot transcriptome data was generated for 23 *tropical japonica*, 23 *indica*, 21 *temperate japonica*, 13 *admixed*, 9 *aus*, and 2 *aromatic* accessions. Genes with low variance or expression within the expression set were filtered out, as these genes are uninformative for downstream analyses focused on natural variation in gene expression. A total of 25,732 genes were found to be expressed (>10 read counts) in at least one or more of the 91 accessions. This equates to about 46% of the genes present in the rice genome (total of 55,986 genes in MSUv7 build).

### Divergence between the *Indica* and *Japonica* subspecies are evident at the genetic and transcriptional levels

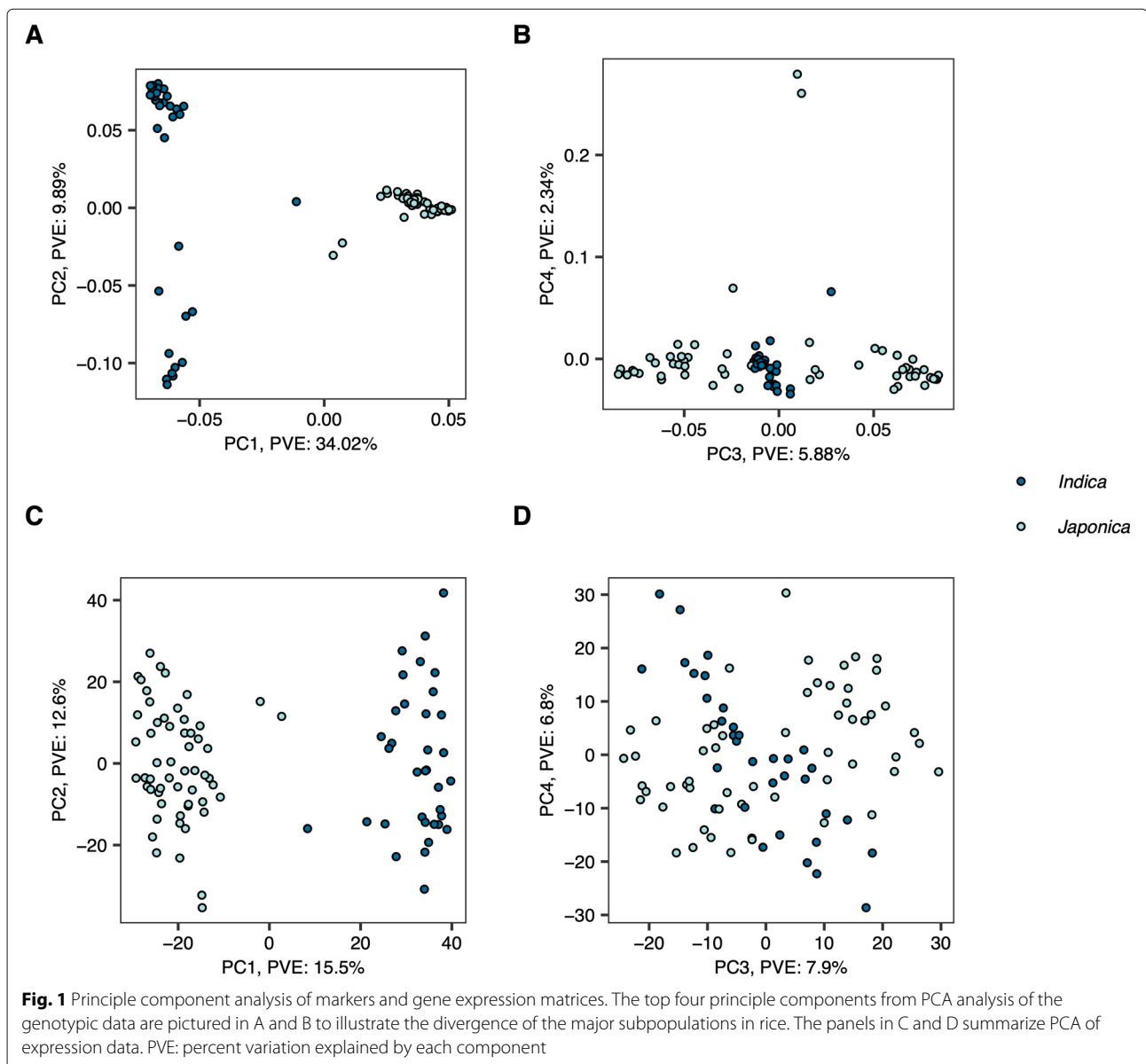
To examine patterns of variation within the transcriptomics data, we performed principle component analysis (PCA) of transcript levels for the 91 accessions. Prior to PCA, lowly expressed genes were removed if they were not expressed (<10 reads) in at least 20% of the samples. This filtering removed approximately 33,311 genes, resulting in a total of 22,675 genes that were used for the principal component analysis based on the normalized read counts. For the genetic analysis, we used 32,849 SNPs. PCA analysis of the expression matrix resulted in a clear

separation between the two subspecies along PC1, suggesting a significant transcriptional divergence between *Indica* and *Japonica* (Fig. 1). The first PC accounted for approximately 26.8% of the variation in gene expression. While PC1 was able to differentiate between the two subspecies at the transcriptional level, no clear clustering of accessions was observed along other PCs (Fig. 1). These results suggest that the two subspecies of cultivated rice have divergent transcriptomes, but the transcriptomes of the subpopulations are more similar. Consistent with these results, differentiation between the subspecies was clearly evident along PC1 using the genetic (SNP) data alone (Fig. 1a,b). The clustering of accessions along PCs 2-4 for the SNP data was consistent with those described by

Zhao et al. [16] (Fig. 1), and were effective in discerning the two subpopulations in rice. These results collectively suggest that the two subspecies are highly divergent at the genetic and transcriptional levels.

#### Differential expression analysis reveals contrasting expression between subspecies

To further explore the differences and identify genes that display divergent expression between the two subspecies, the 91 accessions were first classified into *Indica* and *Japonica*-like groups, using the program STRUCTURE with the assumption of two groups and no admixture [19]. A total of 35 accessions were assigned to the *Indica* subspecies, while 56 were assigned to the *Japonica* sub-



species. Next, a linear mixed model was fit for each of the 26,675 genes, where subspecies was considered a fixed effect and accession as a random effect. A total of 7,417 genes were found to exhibit contrasting expression between the two subspecies ( $FDR \leq 0.001$ , Additional file 3). Of these genes, 4,210 (57%) showed significantly higher expression in *Japonica* relative to *Indica*, while 3,207 (43%) showed higher expression in *Indica* relative to *Japonica*.

This divergent expression levels observed between the two subspecies could be the result of the presence or absence of genes within the subspecies. To this end, we sought to identify genes showing a presence-absence expression variation (PAV). Genes with a read count greater than 10 were considered as expressed and coded as 1 while those with read counts less than 10 were coded as 0. These genes were further filtered, so that genes that were expressed in at least 20%, but no more than 80% of the samples were retained for downstream analyses. A logistic mixed effects model was fit for the 4,163 genes meeting this criteria. In total, 1,980 genes showed evidence of PAV between the two subspecies ( $FDR < 0.001$ ; Additional file 3). This analysis, enriched for genes that were expressed at higher frequency in *Japonica* rice compared to *Indica*. For instance, 1,435 genes were found to be expressed at a significantly greater frequency in *Japonica* relative to *Indica*, while only 545 were found to be expressed predominately in *Indica*. Moreover, we detected significant enrichment for GO terms associated stress response (GO:00006950) and response to biotic stress (GO:0009607), as well genes with kinase activity (GO:0016301). Within *Indica*-specific genes, only a single GO category was enriched for oxygen binding activity (GO:0019825; Table 1). Moreover, 173 were identified with no evidence of expression in *Indica* while only 18 were identified in *Japonica*. Collectively, these results suggest that the divergence between *Indica* and *Japonica*

subspecies may be due, in part, to differences in mean expression levels as well as presence-absence expression variation.

### **Japonica subspecies exhibits reduced genetic and transcriptional diversity**

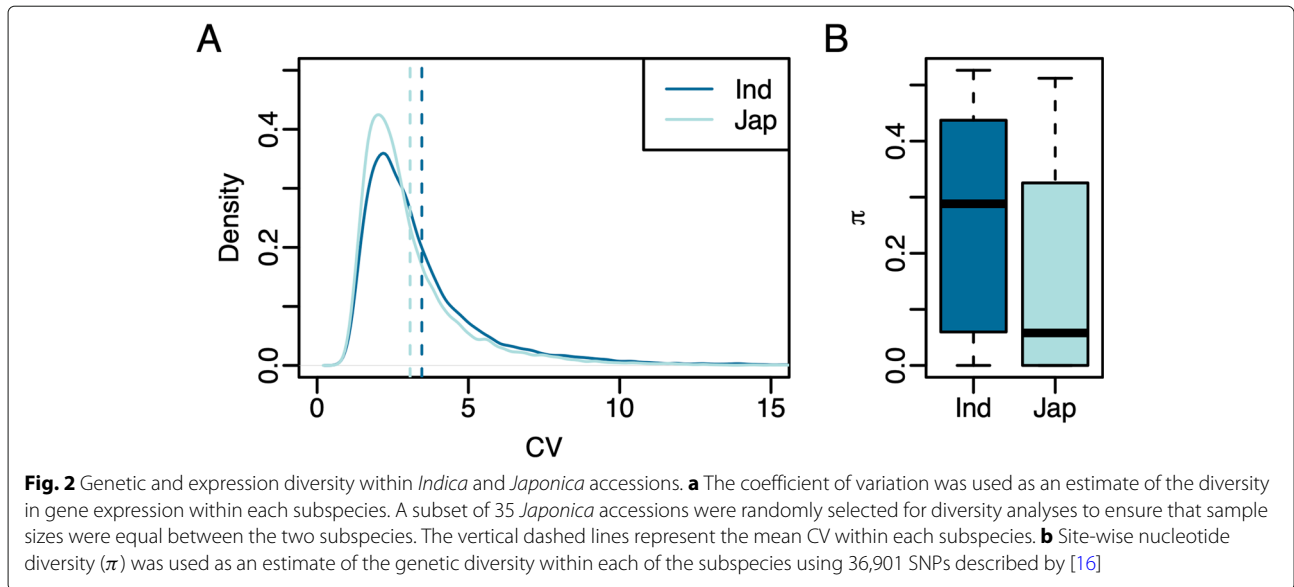
Several studies have shown that the unique domestication history of the two subspecies has resulted in large differences in the overall genetic diversity between the two subspecies, with *Indica* being more genetically diverse than *Japonica* [2, 20–22]. We next explored the variation in gene expression within each subspecies. Two metrics were used to examine the differences in diversity at both the genetic and transcriptional levels within each subspecies: nucleotide diversity ( $\pi$ ) and the coefficient of variation (CV). Diversity analyses within each subspecies may be influenced by differences in sample size. Since the number of *Japonica* accessions were greater than *Indica*, a subset of 35 *Japonica* accessions were randomly selected for diversity analyses. The results for the full set of 56 *Japonica* accessions are provided as Fig. S1.

Expression diversity was estimated using the coefficient of variation (CV) for 22,675 genes. CV was significantly different between the two subspecies (Wilcoxon rank sum test,  $p < 0.0001$ ; Fig. 2). The *Indica* subspecies exhibited approximately 12.6% higher expression diversity compared to *Japonica*. On average, CV in the *Indica* subspecies was 3.46, while in the *Japonica* subspecies the mean CV was 3.07. These results suggest that the transcriptional diversity is lower in the *Japonica* subspecies compared to *Indica*. CV estimates using the complete set of *Japonica* accession were similar (CV: 3.46 and 3.10 for *Indica* and *Japonica*, respectively; Fig. S1).

Genetic diversity within each subspecies was estimated using  $\pi$  for 33,543 SNPs in randomly selected 35 *Indica* and 35 *Japonica* accessions. Similar differences were observed for  $\pi$  as CV, however the differences between

**Table 1** Gene ontology (GO) enrichment analysis for genes exhibiting significant presence-absence expression variation (PAV) ( $FDR < 0.001$ ). GO enrichment was conducted using AgriGO using the MSU V7 genome build without transposable elements as a background. GO enrichment was conducted separately for genes expressed predominately in each subspecies

Subspecies	Ont. Cat.	GO Description	No. in input	No. in background	$p$ -value	$FDR$
<i>Japonica</i>	P	response to stress	137	4660	$1.5 \times 10^{-10}$	$5.2 \times 10^{-8}$
	P	response to stimulus	172	6928	$1.0 \times 10^{-7}$	$1.7 \times 10^{-5}$
	P	response to biotic stim.	43	1404	$2.4 \times 10^{-4}$	$2.7 \times 10^{-2}$
	F	oxygen binding	25	390	$5.0 \times 10^{-8}$	$4.5 \times 10^{-6}$
	F	nucleotide binding	92	3490	$2.4 \times 10^{-5}$	$1.1 \times 10^{-3}$
	F	transferase activity	120	5200	$3.6 \times 10^{-4}$	$9.6 \times 10^{-3}$
	F	catalytic activity	271	13508	$4.2 \times 10^{-4}$	$9.6 \times 10^{-3}$
	F	kinase activity	69	2699	$6.4 \times 10^{-4}$	$9.6 \times 10^{-3}$
<i>Indica</i>	F	oxygen binding	13	390	$1.5 \times 10^{-4}$	$8.8 \times 10^{-3}$

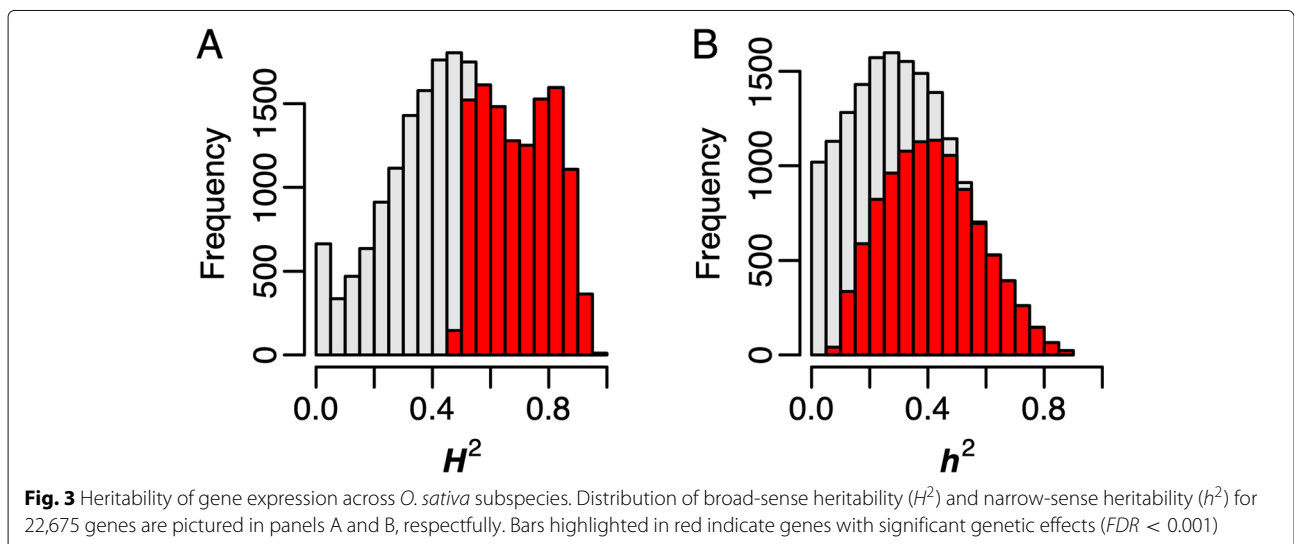


subspecies was much greater (Wilcoxon rank sum test,  $p < 0.0001$ ; Fig. 2). The *Indica* subspecies showed a 64.7% higher nucleotide diversity ( $\pi$ ) compared to *Japonica*. On average,  $\pi$  estimates were 0.26 for *Indica* and 0.17 for *Japonica*. These results are consistent with reports by Huang et al. [2] and Garris et al. [23], and are in agreement with the expression diversity reported above. Together these data suggest that the *Japonica* subspecies exhibits less genetic and transcriptional diversity compared to *Indica*.

#### Gene expression is heritable in cultivated rice

The above analyses shows a strong differentiation between the subspecies at transcriptional and genetic levels, and presents a possible linkage between expression and

genetic diversity. However, the extent of variation in gene expression that can be accounted by genetic variation is not yet determined. To estimate the extent to which variation in gene expression is under genetic control, a mixed model was fit to the expression of each of the 22,675 genes and the variance between accessions was estimated. The significance of the random *between* – *accession* term was determined using a likelihood-ratio test. The broad-sense heritability ( $H^2$ ) was estimated as the proportion of the total variance explained by between-accession variance to total variance. A total of 11,895 genes showed a significant *between* – *accession* variance ( $FDR < 0.001$ ;  $H^2 \geq 0.47$ ), which accounts for approximately 53% of the genes expressed in at least 20% of the samples (Fig. 3a; Additional file 4).  $H^2$  ranged from 0.97 to 0.47, with 4,606



genes showing highly heritable expression ( $H^2 > 0.75$ ), 7,145 showing moderate  $H^2$  ( $0.5 < H^2 \leq 0.75$ ), and the remaining 146 showing low  $H^2$ .

To determine the extent to which additive genetic effects could explain variance in gene expression, a genomic relationship matrix was constructed using 32,849 SNPs following VanRaden [24] and variance components were estimated using a mixed linear model for each gene. A total of 10,125 genes were identified with significant  $h^2$  (Additional file 4). Of these, 234 genes had highly heritable expression ( $h^2 \geq 0.75$ ), while 2,750 genes showed moderate heritability ( $0.5 \leq h^2 < 0.75$ ) (Fig. 3b). An additional 7,141 genes showed low narrow sense heritability ( $h^2 < 0.5$ ). Collectively, these results indicate that many genes in the rice transcriptome are under genetic control.

#### Genetic variability of gene expression is considerably different between subspecies

The analyses above indicate that the two subpopulations differ at the transcriptional and genetic levels, and that for many genes, variation in expression can be explained by genetic effects. We next asked whether the heritability of gene expression is different between the two subspecies. To this end, the expression dataset was partitioned into *Indica* and *Japonica* subsets and genes with low expression in each subspecies were removed (expressed in less than 20% of the samples). Since the number of accessions for the two subspecies are unequal, 35 *Japonica* accessions were randomly sampled to ensure equal sample size, and the number of genes that were expressed in each subspecies were counted. Here, a gene was considered as expressed if 10 or more reads mapped to the gene in 20% or more of the samples. A total of 22,444 genes were found to be expressed in at least 20% of the samples for the *Japonica* subspecies, while 22,068 were found to be expressed in the *Indica* subspecies. A large number of genes were common to both subspecies (21,166 genes). A total of 1,278 genes were found to be uniquely expressed in *Japonica*, and 902 were found to be uniquely expressed in *Indica*.

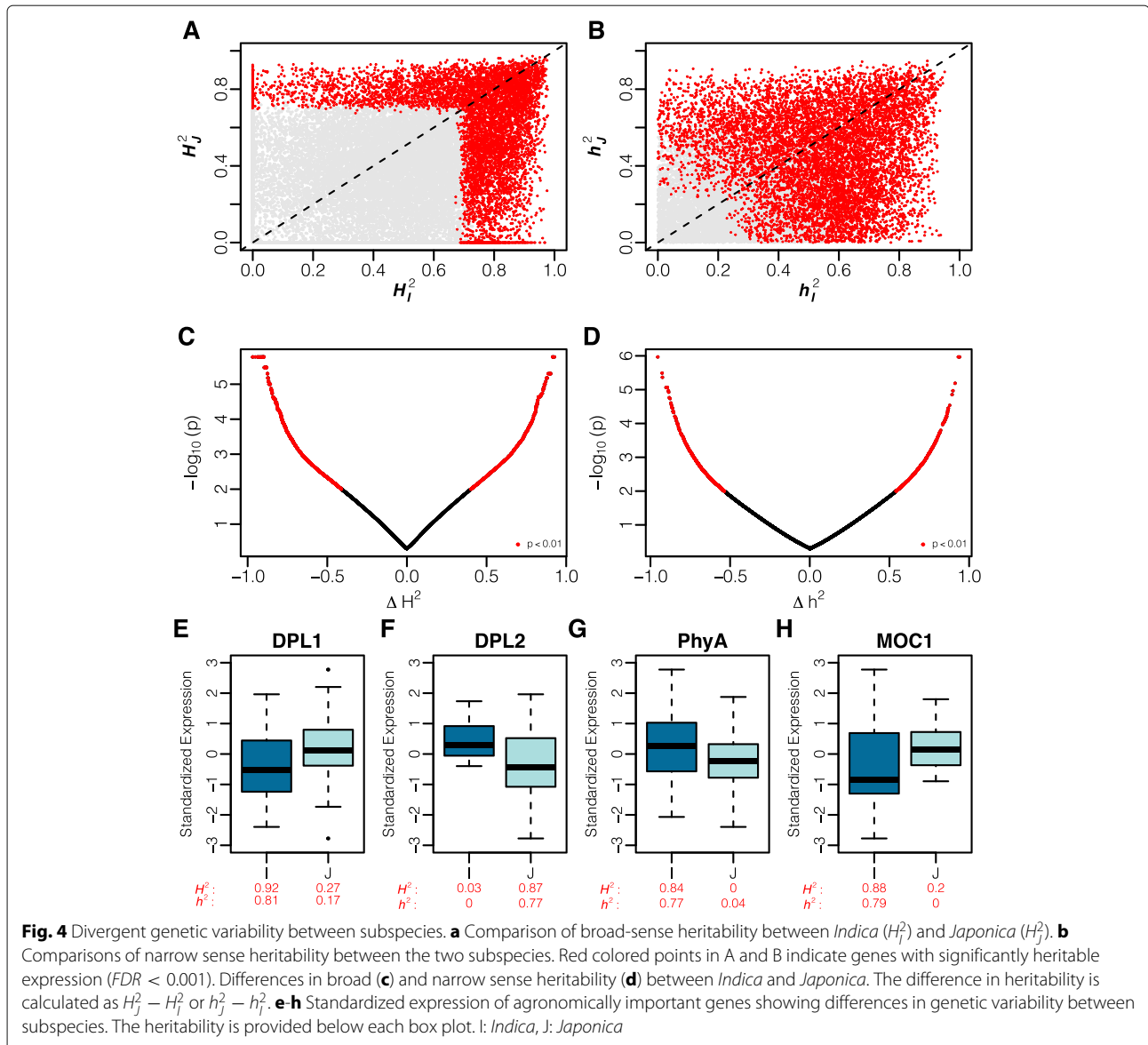
A total of 5,005 genes exhibited significant  $H^2$  in *Indica* and 3,338 genes in *Japonica* ( $FDR < 0.001$ ; Additional file 5). For these genes,  $H^2$  ranged from 0.67 to 0.98 in *Indica* and 0.67 to 0.97 in *Japonica*. A larger number of genes were identified with significant additive genetic variance, with 6,804 identified in *Indica* and 5,103 found in *Japonica*. For these genes, narrow-sense heritability ranged from 0.201 to 0.953 in *Indica* and 0.220 to 0.948 in *Japonica*. Interestingly, few genes showed significant heritable expression in both subspecies. For instance, only 1,681 and 2,644 genes were found to have significant  $H^2$  and  $h^2$ , respectively, in both *Indica* and *Japonica*. Moreover, a comparison of  $H^2$  and  $h^2$  between subspecies showed that for many genes, heritability estimates

were considerably different between *Indica* and *Japonica* (Fig. 4).

To systematically identify genes showing significant differences in  $H^2$  or  $h^2$  ( $\Delta H^2$  and  $\Delta h^2$ , respectively) between subspecies, accessions were randomly partitioned into two groups of equal size and the difference in heritability was estimated between groups. The resampling approach was repeated 100 times. A total of 1,860 genes showed significant differences in  $H^2$  ( $p < 0.01$ ) between the two subspecies, with a minimum absolute difference in  $H^2$  of 0.40. Fewer genes were identified with a significant difference in  $h^2$  between *Japonica* and *Indica* (Additional file 6). Only 1,325 genes were found with significant differences in  $h^2$  between *Indica* and *Japonica*, and the absolute difference in  $h^2$  ranged from 0.54 to 0.95 (Fig. 4).

These differences in heritability may be due to insufficient phenotypic variation (e.g. lack of expression diversity), or changes in the genetic or environmental factors that contribute to phenotypic variation. Thus, to further examine the potential causes of the observed differences in heritability, we quantified the expression diversity (CV), genetic variation and environmental variation within each subspecies for genes exhibiting  $\Delta H^2$  and  $\Delta h^2$ , as well as those with shared heritable variation. For genes exhibiting subspecies-specific genetic variability, the loss of heritability was largely due to an increase in environmental effects on phenotypic variation in the subspecies lacking heritability rather than loss of phenotypic variation. This is clearly evident in Additional file 2. The mean CV for  $\Delta H^2$  genes decreased slightly in subspecies lacking genetic variability. However, for these same genes the proportion of phenotypic variation that was explained by environmental effects increased significantly in subspecies lacking genetic variability. Collectively, these results suggest that the differences in heritability exhibited between the subspecies is driven largely by loss of genetic variability and an increase in environmental effects rather than a loss of phenotypic variation.

Interestingly, several genes that have been reported to have divergent genetic variants between *Indica* and *Japonica* were found within  $\Delta H^2$  and  $\Delta h^2$  genes. For instance, *DOPPELGANGER1* (*DPL1*) showed significantly higher  $H^2$  and  $h^2$  in *Indica* relative to *Japonica* ( $H^2$ : 0.92 and 0.27, respectively,  $p_{\Delta H^2} = 0.011$ ;  $h^2$ : 0.81 and 0.17,  $p_{\Delta h^2} = 0.004$ ; Fig. 4e). However for *DOPPELGANGER2*, the converse was true. Significantly higher  $H^2$  and  $h^2$  was observed in *Japonica* relative to *Indica* ( $H^2$ : 0.87 and 0.03,  $p_{\Delta H^2} < 0.001$ ;  $h^2$ : 0.77 and 0, respectively,  $p_{\Delta h^2} = 0.005$ ; Fig. 4f). Mizuta et al. [25] showed that *DPL1* and *DPL2* are important regulators of *Indica*-*Japonica* hybrid incompatibility, and non-functional alleles arose independently for *DPL1* and *DPL2* within the *Indica* and *Japonica* subspecies respectively. Thus the results reported by Mizuta



et al. [25] are consistent with the divergent genetic variability in expression observed in our study. In addition to *DPL1* and *DPL2*, a gene that is important for the regulation of shoot growth/ architecture, *MOC1*, also displayed divergent genetic variability between subspecies. *MOC1* showed significant differences in both  $H^2$  and  $h^2$  (Fig. 4h). Collectively, these results show that the two subspecies are divergent at the transcriptional and genetic levels. Moreover, many genes exhibit large differences in genetic variability between the *Indica* and *Japonica*, suggesting that these genes may be regulated by divergent genetic mechanisms.

#### Joint eQTL analysis assesses cis-regulatory divergence between subspecies

The differences in the narrow-sense heritability between

subspecies observed for some genes suggest a divergence in the genetic regulation of these genes. Using the transcriptional and genotypic data for this population, we next sought to identify genetic variants that can explain this divergent genetic regulation. To this end, a joint eQTL analysis was conducted across subspecies using the eQTL Bayesian model averaging (BMA) approach described by Flutre et al. [26]. With this approach, the posterior probability of specific configurations can be formally tested; in other words, the probability that an eQTL is present/active in both the *Indica* and *Japonica* subspecies or unique to a given subspecies can be determined. The 91 accessions were classified into *Indica* and *Japonica* subspecies using STRUCTURE as described earlier, yielding 35 *Indica*-type and 56 *Japonica*-type accessions. eQTLs were modeled for genes showing significant  $H^2$  in

at least one subspecies (6,307 genes) and 274,499 SNPs. For each gene, associations were tested for SNPs within 100kb of the transcription start site. A total of 5,097 genes were detected with one or more eQTL at an FDR of 0.05 (Additional file 7). This equates to approximately 81% of the genes displaying heritable expression, and indicates that a large portion of genes with heritable expression are regulated by variants in close proximity to the gene.

To identify eQTL genes that were specific to a given subspecies, the SNP with the highest probability of being the eQTL was selected for each gene, and the posterior probability for all three configurations (*Indica*-specific, *Japonica*-specific, and across subspecies) was compared. Of the 5,097 eQTL genes detected, 80% (4,077 genes; 3,826 unique SNPs) were detected across subspecies, 18% (914 genes; 880 unique SNPs) were detected for *Indica* accessions, and only 2% (106 genes; 103 unique SNPs) were detected in *Japonica* accessions. These results indicate that while a large portion of *cis*-eQTLs are shared across the two subspecies of cultivated rice, many genes are regulated by unique *cis* regulatory mechanisms that are specific to the *Indica* subspecies.

#### Signatures of selection are evident among subspecies specific eQTL

The presence or absence of *cis*-regulatory variants within a given subspecies may be the result of the unique domestication histories that have shaped *Indica* and *Japonica*, and/or driven by environmental adaptation of the wild progenitors from which they were derived. The absence of variation at the eQTL SNP could be due to sampling during differentiation of the wild progenitors or during domestication (e.g. lost purely by chance), or due to selective pressures imposed by the environment or humans. In the case of selection, we expect to see reduced genetic diversity around the eQTL compared to the rest of the genome. To determine whether the absence of subspecies-specific eQTL are the result of selection, we calculated the average nucleotide diversity ( $\pi$ ) in 100 Kb windows around significant subspecies-specific eQTL within each subspecies and compared these values to the overall average  $\pi$  for 100 Kb windows across the genome within each subspecies using a two-sided *t*-test. Comparisons within each subspecies of  $\pi$  for eQTLs and the genome-wide average should account for the inherent differences in  $\pi$  between the two subspecies.

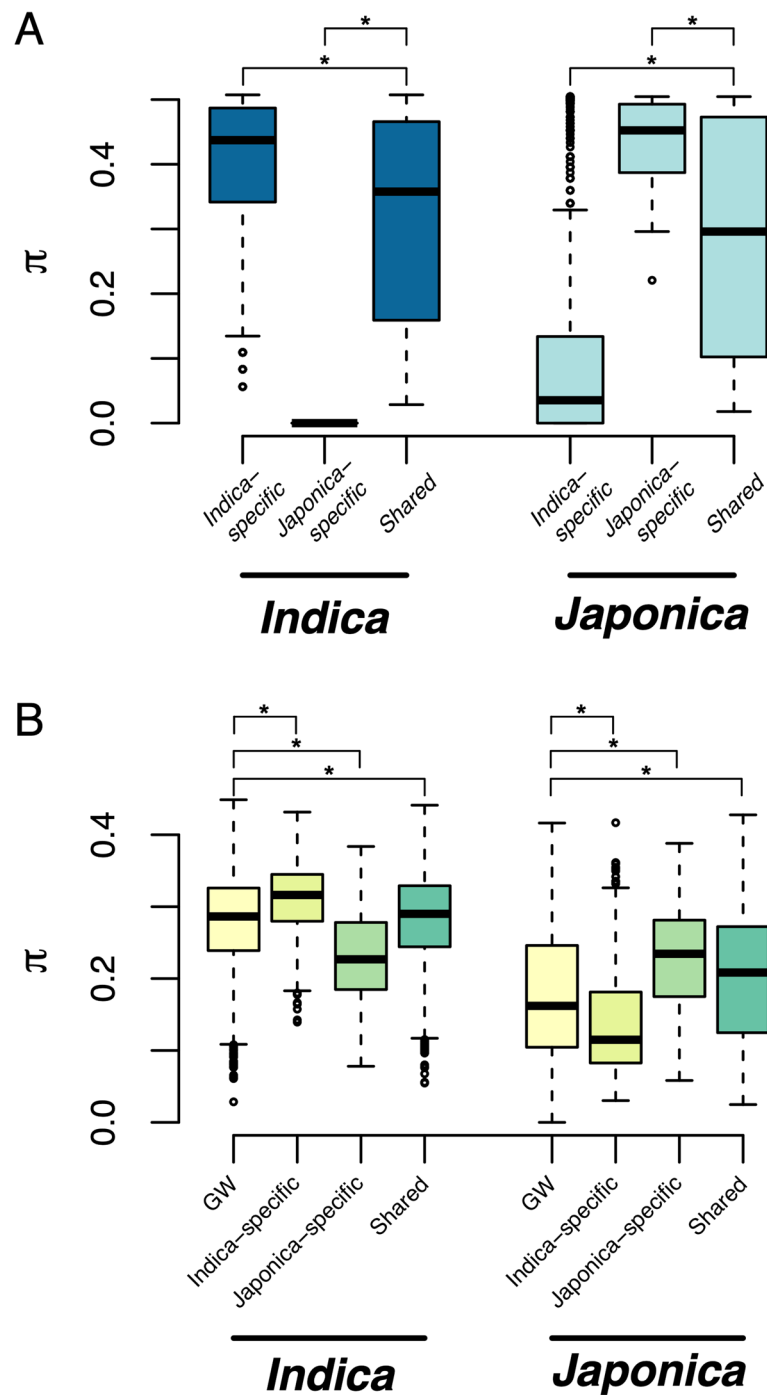
Consistent with what would be expected under selection, a significant reduction in nucleotide diversity was observed for eQTL SNPs that were absent in a subspecies, as well as for regions around subspecies-specific eQTL (Fig. 5). For instance, for *Indica*-specific eQTL, the average  $\pi$  in *Japonica* was approximately 22% lower than the genome-wide average (0.138 and 0.176, respec-

tively;  $p < 1 \times 10^{-15}$ ). Similarly, the average  $\pi$  in *Indica* for *Japonica*-specific eQTL was about 16% lower than the genome-wide average (0.235 and 0.279, respectively;  $p = 3.85 \times 10^{-10}$ ). Interestingly, slightly higher nucleotide diversity was observed for regions around subspecies-specific eQTL in subspecies in which they were detected compared to genome-wide nucleotide diversity, as well as for shared eQTL when compared to genome-wide nucleotide diversity. Collectively, these results indicate that the absence of eQTL within a given subspecies may be the result of selective pressures that reduced genetic diversity within the eQTL regions.

Given the small sample size in the current study ( $n = 91$ ) we sought to confirm these results using resequencing data for a larger population of 3,024 diverse rice accessions [27–30]. To this end, we extracted SNP information for 3,024 rice accessions in the same 100 Kb window surrounding eQTL, and examined  $\pi$  within each subpopulation for these regions. As above,  $\pi$  within these regions were compared with genome-wide averages for 100 kb windows. The 3,024 rice accessions are classified into 12 subpopulations: *admix* (103 accessions), *aromatic* (76 accessions), *aus* (201 accessions), *indica1A* (209 accessions), *indica1B* (205 accessions), *indica2* (285 accessions), *indica3* (475 accessions), *indica-X* (615 accessions), *japonica-X* (83 accessions), *subtropical japonica* (112 accessions), *temperate japonica* (288 accessions), and *tropical japonica* (372 accessions). The *Indica* subspecies are represented by *indica1A*, *indica1B*, *indica2*, *indica3*, and *indica-X*; while the *Japonica* subspecies consists of the *japonica-X*, *subtropical japonica*, *temperate japonica*, and *tropical japonica* subpopulations.

Consistent with the results derived from the 91 accessions,  $\pi$  within subspecies-specific eQTL was lower in subpopulations lacking the eQTL (Fig. 6). For instance, for the *Japonica* subpopulations (*japonica-x*, *subtropical japonica*, *temperate japonica*, and *tropical japonica*)  $\pi$  estimates for *Indica*-specific eQTL were considerably lower than those for *Indica* subpopulations (*indica-1A*, *indica-1B*, *indica-2*, *indica-3*, and *indica-x*). The converse was true for *Japonica*-specific eQTL, with lower  $\pi$  observed in *Indica* subpopulations relative to *Japonica*. However for the shared eQTL,  $\pi$  estimates were higher than the genome-wide averages, suggesting that genetic diversity within regions that regulate gene expression is maintained. To identify specific loci that may have been targeted by selection, we selected eQTL regions with an average  $\pi$  within a 100 Kb window that was below the 5% quantile for genome-wide average for a given subspecies. Consistent with the results above, we observed a greater frequency of low diversity eQTL regions in subspecies lacking the subspecies-specific eQTL. For instance, approximately 11% of the 880 *Indica*-specific

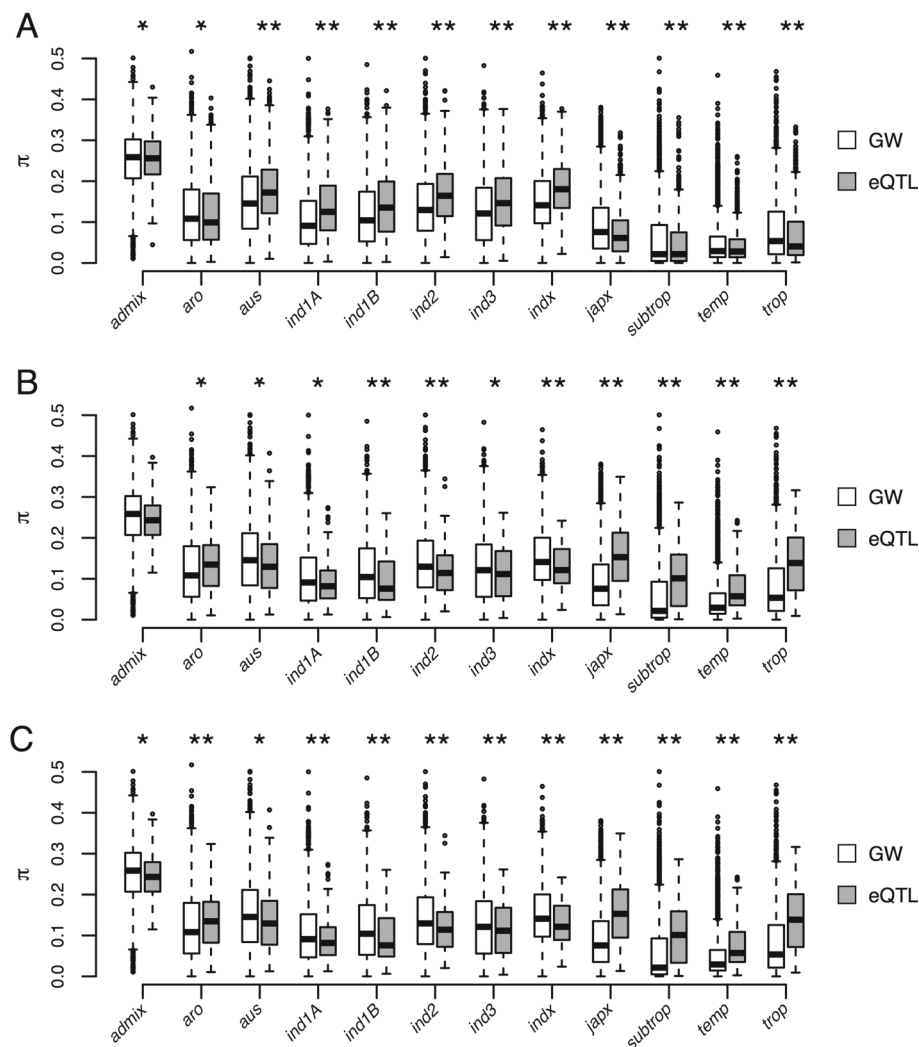




**Fig. 5** Nucleotide diversity at cis-eQTL. **a** Nucleotide diversity ( $\pi$ ) for the most significant SNP for each cis-eQTL. The distribution of  $\pi$  is pictured from each subspecies and each eQTL type. **b** Distribution of  $\pi$  for 100 kb windows around the most significant SNP for each cis-eQTL. Genome-wide (GW)  $\pi$  was determined by randomly selecting X SNPs that were more than 100 kb from a cis-eQTL and low diversity SNPs (MAF < 0.1 in both subspecies) were removed prior to analyses. Asterisks indicate a significant differences determined via Tukey’s test between eQTL types ( $p < 1 \times 10^{-8}$ )

eQTL were found in regions of low diversity in *Japonica* ( $\pi_{Jap} \leq 0.0645$ ). While for *Japonica*-specific eQTL, 14% (14 of the 103) eQTL regions were lying in regions of low diversity in *Indica* ( $\pi_{Ind} \leq 0.1617$ ). However, for

shared eQTL and for subspecies in which the subspecies-specific eQTL was detected, the converse was true. Only a small percentage of eQTL regions were found within regions of low diversity. For instance, approximately 3.5%



**Fig. 6** Nucleotide diversity at cis-eQTL within subpopulations for 3,053 rice accessions. Average nucleotide diversity ( $\pi$ ) for 100 kb regions surrounding *Indica*-specific, *Japonica*-specific, and shared eQTL are pictured in panels A, B, and C, respectively. For each subpopulation and class of eQTL (e.g. *Indica*-specific, *Japonica*-specific, and shared)  $\pi$  was calculated for each SNP within 100 kb of the most significant eQTL SNP.  $\pi$  for the eQTL windows were compared to a genome wide (GW) average in which regions with eQTL and site with low diversity (MAF  $< 0.01$  in 10 of 12 subpopulations) were excluded. Asterisks indicate significant differences between GW and eQTL regions determined using a two-sided Student's *t*-test (\*  $p < 0.05$ ; \*\*  $p < 0.001$ ). Subpopulations are named following [27] (aro: aromatic; ind1A: indica-1A; ind1B: indica-1B; ind2: indica-2; indx: indica-X; japx: japonica-X; subtrop: subtropical japonica; temp: temperate japonica; trop: tropical japonica)

of shared eQTL were found in regions of low diversity in both *Indica* and *Japonica*, and less than 1% of subspecies eQTL were found in regions of low diversity in the subspecies in which they were detected. Collectively these results suggest that selective pressures may have shaped the cis-regulatory divergence of the *Indica* and *Japonica* subspecies.

## Discussion

The differentiation between the *Indica* and *Japonica* subspecies of cultivated rice has been intensively characterized at the morphological, biochemical, and genetic levels

[2, 3, 5–12, 31–35]. However, the divergence at the transcriptional levels remains understudied. Here, we provide a comprehensive analysis of the transcriptional and cis-regulatory divergence between the major subspecies of rice, and show that the presence or absence of cis regulatory variants within the subspecies is a component of this divergence.

The transcriptional divergence is most evident in the large number of expressed genes showing differences in the magnitude or frequency of expression. Of the 25,732 genes showing evidence of expression in the current study, approximately 29% showed significant differences

in expression levels between the two subspecies. Moreover, approximately 8% of expressed genes showed evidence of presence-absence expression variation. While few studies have examined the differences in expression levels between diverse populations of *Indica* and *Japonica*, recent studies have utilized whole genome sequencing to shed light on the genetic differentiation between the subspecies of cultivated rice [2, 27]. In a recent study, Wang et al. [27] found that on average approximately 15% of all genes showed evidence of PAV between the genomes of *Indica* and *Japonica* accessions, further indicating that PAV is pervasive between the subspecies of cultivated rice. While the number of PAV reported by Wang et al. [27] are nearly two fold higher than those reported in the current study, it is important to note that only a single tissue was sampled for 91 accessions at a single time point. Therefore, while the expression data provides considerable insight into transcriptional variation in cultivated rice, it likely captures only a portion of the total transcriptome given the lack of temporal and spatial resolution. Moreover, Wang et al. [27] captured PAV using 3,010 resequenced rice genomes, while the current study utilized only a fraction of the variation of Wang et al. [27] with RNA sequencing of 91 accessions. Thus, increased sample size via larger populations and more sampling within tissue and developmental context may lead to a better agreement between PAV at the genome and transcriptional levels.

One major challenge for genomic studies that utilize both *Indica* and *Japonica* accessions is choosing an appropriate genome. While cultivated rice consists of two subspecies, many studies that have used accessions from both subspecies often map sequences to the Nipponbare reference genome [16, 17, 27, 30, 36]. Several studies have highlighted structural variation both within and between subspecies of cultivated rice. Thus, some genomic features may not be shared between diverse accessions and Nipponbare [10, 37]. The current study highlights many differences between the *Indica* and *Japonica* subspecies, but does so under the assumption that the genomes of the two subspecies should not be too different. The overall high colinearity of the genomes of the two subspecies and the ability to recover fertile F1 individuals from *Indica*-*Japonica* hybrids suggests that this is a reasonable assumption.

#### **Potential causes of transcriptional divergence between *Indica* and *Japonica***

Lower mean expression values or absence of expression in a given subspecies may be the result of both heritable and non-heritable effects. The availability of high density SNP information for RDP1 allowed us to begin to elucidate the genetic basis of the observed transcriptional divergence between the subspecies of cultivated

rice. A notable portion of genes with evidence of PAV or DE also showed differences in genetic variability between the subspecies (13% and 9% of DE genes showed differences in  $H^2$  and  $h^2$ , respectively, and 20% and 15% of PAV genes showed differences in  $H^2$  and  $h^2$ ), indicating that for many genes, the genetic mechanisms that regulate expression may be different between the two subspecies. However, many genes that display divergent expression patterns have non-significant differences in genetic variability. There are several explanations for this. For one, the thresholds used to identify genetically divergent genes were quite stringent. For instance, genes must have a difference in genetic variability in either the broad sense greater than 0.4022 between subspecies to be labeled as statistically significant, and in the narrow sense 0.5364. Therefore, it is possible that many more DE or PAV genes have different genetic architectures in the two subspecies, but were missed because of the stringency of statistical threshold. A second possibility is that many of the genes showed divergent expression are influenced greatly by the environment, and thus have low heritability. Thus, these genes would be filtered out in these genetic analyses.

The heritable transcriptional divergence may be due to genetic variants that influence gene expression and are divergent between *Indica* and *Japonica*. These include large structural variants (e.g. deletions, insertions, inversions, and/or duplications), or SNPs that may act in cis or trans to influence gene regulation. While high density SNP information is available for this population and can be leveraged to identify SNPs that regulate expression and are divergent between the subspecies, the identification of larger structural variants that influence expression is only attainable through full genome sequencing, which is not currently available for RDP1. As more genetic resources become available for RDP1 this would be a promising future direction to resolve the causal basis of these transcriptional differences.

The availability of high density SNP information for RDP1 allowed us to begin to elucidate the genetic basis of the observed transcriptional divergence between the subspecies of cultivated rice, and classify genetic effects into those that are common between subspecies, or unique to a given subspecies. While the eQTL-BMA approach has proven to be a powerful framework for assessing the specificity of eQTL for a given tissue or population, one potential limitation of eQTL-BMA is that the framework only allows us to model cis-eQTL. Trans-eQTLs are often difficult to detect due the penalties associated with the large number of statistical tests performed, and because trans-eQTL often have small effect sizes and thus require larger dataset for detection. Several studies in humans have shown that cis-eQTL typically only explain 30-40% of genetic variation in expression [38–40]. Thus, the divergent regulatory variants captured in the current study only

reflect a portion of the differences in genetic variation between the two subspecies. Further studies are necessary to shed light on the contribution of trans-regulatory variants on the genetic differentiation between *Indica* and *Japonica* transcriptomes.

The joint eQTL analysis facilitated the identification of 5,097 genes associated with one or more SNP in *cis*. For most of these genes (81%), the *cis*-regulatory variant was shared between both subspecies, indicating that much of the *cis*-regulatory variation is common between the two subspecies. This high degree of overlap is somewhat expected. For one, both *Indica* and *Japonica* originate from populations of the same species, *Oryza rufipogon*. Moreover, crosses between *Indica* and *Japonica* often produce viable offspring, indicating a high degree of colinearity and functional similarity between the genomes. Thus, while considerable differentiation between founder *Oryza rufipogon* populations has been reported and further divergence has likely occurred since domestication, the common origin and inter-specific comparability suggests that the transcriptional regulation and genome structure is similar [2].

#### Functional significance of transcriptional divergence

The current study elucidates the transcription divergence between the major subspecies of cultivated rice. Many of these genes with divergent expression, genetic variability, or regulatory variation have been reported to be underlying important agronomic traits, such as photoperiod adaptation and development. Therefore these observed differences may have potential agronomic significance. However, further studies are necessary to determine whether these expression patterns are conserved in other tissues or developmental stages.

Among these divergent genes, we identified three genes (*OsPhyA*, *OsPhyC*, and *OsCO3*), that have been reported to be associated with the timing of reproductive development in response to day length that had significant heritability in *Indica* only. The two phytochrome genes, *OsPhyA* and *OsPhyC* are activated under long-day conditions and repress flowering time through *OsGhd7* [41, 42]. Although no studies have shown whether *OsCO3* participates directly in the pathway involving *OsPhy* genes, disruption of *OsCO3* interferes with photoperiod sensitivity and/or flowering time [43]. For instance, Kim et al. [43] showed that the overexpression of *OsCO3* delayed flowering under short-day conditions. In most rice varieties, short-days promote the transition from vegetative to reproductive growth [44]. However, *temperate japonica* rice varieties, which are adapted to higher latitudes have been selected to initiate flowering in long-days to escape the negative impact of low temperatures in autumn on pollen fertility [45–47]. All genes showed heritable expression only in the *Indica* subspecies, indicating that in the

*Japonica* subspecies expression variation may be driven largely by non-genetic effects. Moreover, the patterns of genetic variability for these genes are consistent with their potential role in the adaptation of flowering in different environments for *Indica* and *Japonica*.

In addition to genes regulating rice phenology, several genes were identified that have been reported to play important roles in the regulation of shoot architecture (*D18*, *MT2b*, and *MOC1*). For instance, two genes *dwarf18* (*D18*) and *Metallothionein2b* (*MT2b*) have been reported to regulate plant height [48, 49]. *D18* encodes a GA- $\beta$  hydroxylase and is involved with GA biosynthesis. Loss of function mutants exhibit a severe dwarf phenotype [48]. Interestingly, *D18* was found have an *Indica*-specific eQTL, but did not exhibit a difference in  $H^2$  or  $h^2$  between the two subspecies ( $p = 0.046$  and  $p = 0.19$ , respectively), indicating that genetic differences may be confined to local regions around *D18*.

#### Conclusion

The morphological and genetic differences between subspecies of cultivated rice have been studied extensively, however the divergence of *Indica* and *Japonica* at the transcriptional and regulatory levels is largely unresolved. Here, we provide, to date, the first detailed population-level characterization of transcriptional diversity within cultivated rice, and assess the divergence in transcriptomes and expression variation between *Indica* and *Japonica*. We find that many agronomically important genes exhibit differences in expression levels, and/or *cis*-regulatory variation between the subspecies. These resources provided by this study can serve as a foundation for future functional genomics studies in rice and other crops, and can be further utilized to connect gene function with natural variation in gene expression.

#### Methods

##### Plant materials and growth conditions

This study used 91 diverse accessions from the Rice Diversity Panel1 (RDP1) [16–18]. Seeds were obtained from the USDA-ARS Dale Bumpers Rice Research Center. The 91 accessions consisted of 13 *admixed*, 2 *aromatic*, 9 *aus*, 23 *indica*, 21 *temperate japonica*, and 23 *tropical japonica* accessions.

Seeds were dehusked manually and germinated in the dark for two days at 28°C in a growth cabinet (Percival Scientific), and were exposed to light (120  $\mu\text{mol m}^{-2}\text{s}^{-1}$ ) twelve hours before transplanting to acclimate them to the conditions in the growth chamber. The seeds were transplanted to 3.25" x 3.25" x 5" pots filled with Turface MVP (Profile Products) in a walk-in controlled environment growth chamber (Convion). The plants were cultivated in the absence of intentional stress conditions. The pots were placed in 36" x 24" x 8" tubs, that were filled with

tap water. Four days after transplanting the tap water was replaced with half-strength Yoshida solution [50] (pH 5.8). The pH of the solution was monitored twice daily and was recirculated from a reservoir beneath the tubs to the growth tubs. The temperatures were maintained at 28°C and 25°C in day and night respectively and 60% relative humidity. Lighting was maintained at  $800 \mu\text{mol m}^{-2}\text{s}^{-1}$  using high-pressure sodium lights (Phillips).

### RNA extraction and sequencing

Ten days after transplant, aerial parts of the seedlings were excised from the roots and frozen immediately in liquid nitrogen. The samples were ground with TissueLyser II (Invitrogen) and total RNA was isolated with RNAeasy isolation kit (Qiagen) according to manufacturer's instructions. On-column DNase treatment was performed to remove genomic DNA contamination (Qiagen). Sequencing was performed using Illumina HiSeq 2500. Sixteen RNA samples were combined in each lane. Two biological replicates were used for each accession.

### Sequence alignment, expression quantification, and differential expression analysis

Quality control for raw reads was performed using the package FastQC [51]. The Illumina 101-bp single-end reads were screened and trimmed using Trimmomatic to ensure each read has average quality score larger than 30 and longer than 15 bp, and were aligned to the rice genome (*Oryza sativa* MSU Release 6.0) using TopHat (v.2.0.10), allowing up to two base mismatches per read. Reads mapped to multiple locations were discarded [52, 53]. The number of reads for each gene sequence was counted using the HTSeq-count tool with the "union" resolution mode [54]. For down-stream genetic analyses, a variance stabilizing transformation was performed on normalized read counts to provide approximately homoskedastic values in DESeq2 [55].

To identify genes that exhibited differential expression between the two subspecies, a mixed linear model was fit that included subspecies as the main fixed effect and accession as a random effect in lme4 [56]. This full model was compared to a reduced model that lacked subspecies as a fixed effect using a likelihood-ratio test. Prior to differential expression analysis, expression levels were quantile normalized to ensure a Gaussian distribution. Benjamini and Hochberg's method was used to control the false discovery rate, and genes with an FDR  $\leq 0.001$  were considered differentially expressed [57].

Genes showing differences in presence-absence expression variation (PAV) was determined using a mixed-effects logistic regression model. Briefly, for each sample the expressed genes (number of reads >10) were assigned 1, while those with 10 or less reads were assigned a 0. A logistic regression model was fit using the 'glmer' func-

tion in 'lme4' and included subspecies as a fixed effect and accession as random [56]. The significance of the fixed effect of subspecies was determined by comparing the full model above with a reduced model that lacked subspecies using a likelihood-ratio test. Benjamini and Hochberg's method was used to control the false discovery rate, and genes with an FDR  $\leq 0.001$  were considered as having presence-absence expression variation [57].

### Subspecies classification

The 91 accessions were classified into two subspecies using the software STRUCTURE [19]. Briefly, the software was run using the 44k SNP data, assuming two subpopulations (K=2), with 20000 MCMC replicates and a burn-in of 10000 MCMC replicates.

### Expression and genetic diversity analyses

Principle component analysis of gene expression was conducted for the 91 accessions using 22,675 genes after variance stabilizing transformation. For PCA of SNP data the 44k dataset described by Zhao et al. [16] was used. SNPs with a MAF <0.10 were removed prior to PCA analysis.

The coefficient of variation (CV) was used to estimate the diversity in gene expression within the *Indica* and *Japonica* subspecies. Prior to estimating CV genes with low expression (i.e. those with read counts of  $\leq 10$  in  $\geq 20\%$  of the samples) were removed, leaving a total of 22,503 genes in *Japonica* and 21,719 genes in *Indica*. For the estimation of  $\pi$ , SNPs were extracted for each subspecies and SNPs with MAF < 0.05 were excluded. In total 201,891 SNPs were retained for *Indica* and 161,715 for *Japonica*.  $\pi$  was estimated at each SNP using the site-pi function in VCFtools [58].

### Heritability estimates

Heritability, both in the broad ( $H^2$ ) and narrow sense ( $h^2$ ), was estimated across subspecies for 22,675 genes that were expressed in both *Indica* and *Japonica*. To estimate  $H^2$  a mixed model was fit using lme4 where accession was considered a random effect, and significance of  $H^2$  was assessed using a restricted likelihood-ratio test in the RLRTsim package [56, 59]. Benjamini and Hochberg's method was used to control the false discovery rate, and genes with an FDR  $\leq 0.001$  were considered to have significant genetic variability [57]. To assess heritability in the narrow sense ( $h^2$ ) a mixed model was fit in asreml-R [60]. Briefly, a genomic relationship matrix (G) was estimated according to [24] using the approximately 36,901 SNPs described by Zhao et al. [16]. G is estimated as  $G = \frac{Z_{cs}Z_{cs}'}{m}$ , where  $Z_{cs}$  is the centered and scaled marker matrix and  $m$  is the number of markers. A likelihood-ratio test was used to assess significance and Benjamini and Hochberg's method was used to control the false discovery rate. Genes

with an  $FDR \leq 0.001$  were considered to have significant genetic variability [57].

Heritability was assessed within subspecies using the same approaches as described above. However, due to the unequal sample size for the *Indica* and *Japonica* subspecies, a random set of 35 *Japonica* accessions were selected. Genes showing low expression ( $< 10$  reads in  $< 20\%$  of samples) in either subspecies were removed prior analysis, leaving 22,444 genes in *Japonica* and 22,068 genes in *Indica*.

#### Assessing differences in genetic variability between subspecies

To identify genes showing significant differences in genetic variability ( $H^2$  or  $h^2$ ) between subspecies, a permutation approach was used. Here, the 91 accessions were randomly partitioned into two groups of equal size (35 accessions each). Heritability was estimated as described above and the difference in heritability between each group was calculated. The resampling approach was repeated 100 times for both  $H^2$  and  $h^2$ . This process effectively estimated a null distribution of  $\Delta H^2$  and  $\Delta h^2$  values. The heritability estimates for each subspecies was used to calculate the differences in  $H^2$  and  $h^2$  between the two subspecies as  $\Delta H^2 = H_j^2 - H_l^2$  or  $\Delta h^2 = h_j^2 - h_l^2$ . These values were compared with the null distribution to assess significance.

#### Joint cis-eQTL analysis

eQTLs were jointly detected using the eQTL-BMA (Bayesian model averaging) described by Flutre et al. [26] for 26,675 genes and 274,499 SNPs ( $MAF > 0.10$ ) [61]. Prior to eQTL mapping BLUPs for each gene was calculated and the gene expression level of each gene was transformed into the quantiles of a standard Normal distribution with ties broken randomly. To control for the effects of population structure the first four PCs derived from PCA analysis of 44k SNP dataset were included in the linear model. Briefly, to identify eQTL and control false discovery rate (FDR) a gene-level permutation approach was used within the eQTL-BMA software. Using the `eqtlbma_bf` program, 10,000 permutations were performed with the following settings: `-maf 0.1`, `-nperm 10000`, `-trick 1`, `-tricut 10` and `-error uvlr`. Genes were considered to have an eQTL if the  $FDR < 0.05$ . These permutations were used to estimate  $\pi_0$ , the probability for a gene to have no eQTL in any subspecies. Here, expression from both *Japonica* and *Indica* samples were analyzed together with the option `-error uvlr` specified. Next, a hierarchical model with an expectation-maximization algorithm was used to estimate hyper-parameters and configuration probabilities using the `eqtlbma_hm` program. These configurations were *Indica*-specific, *Japonica*-specific, and present in both subspecies. Lastly, the

`eqtlbma_avg_bfs` program was run to obtain (i) the posterior probability (PP) of a gene to have an eQTL in at least one subspecies, (ii) PP for a SNP to be the causal SNP for the eQTL, (iii) PP for the SNP to be an eQTL, (iv) PP for the eQTL to be present in one subspecies, and (v) PP for the eQTL to be present for a specific configuration. SNP-gene pairs were determined to be specific to a given subspecies or shared if the  $PP > 0.5$  for a given configuration.

#### Detecting evidence of selection at cis-eQTL

To determine whether the absence of an eQTL was due to of selection, first SNPs from the HDRA dataset within 100kb of each significant eQTL were extracted for the 91 accessions [61]. For each SNP, nucleotide diversity was determined using the `site-pi` function in `VCftools` and was averaged across the 100kb window [58]. Secondly, a genome-wide diversity level was determined for each subspecies. Here, SNPs that were within 100kb of an eQTL were excluded, as well as those that exhibited low diversity in both subspecies ( $MAF < 0.1$  in both *Indica* and *Japonica*). Nucleotide diversity was determined as described above for each SNP, and the average was taken for 100kb windows. For each class of eQTL (e.g. *Indica*-specific, *Japonica*-specific, and shared), a two-sided Student's *t*-test was performed to assess whether the mean  $\pi$  was different from the genome-wide average for each subspecies and class of eQTL.

A similar approach was taken for the 3kg data [30]. For each eQTL SNP, all SNPs within 100kb of the eQTL SNP was extracted from the 4.8M core SNP data. The  $MAF$  was determined for each of the 12 subpopulations in the 3kg data, and SNPs that had low diversity ( $MAF < 0.01$ ) in 10 of the 12 subpopulations were excluded from further analyses. As above,  $\pi$  was calculated for each site. An average  $\pi$  was determined for each subpopulation at each eQTL by taking the average  $\pi$  across the 100kb window. To obtain a genome wide average, eQTL regions were excluded and  $\pi$  was obtained for each subpopulation by averaging  $\pi$  across the 100kb region. Finally, as above a two-sided Student's *t*-test was performed to assess whether the mean  $\pi$  was different from the genome-wide average for each subpopulation and class of eQTL.

#### Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s12864-020-06786-6>.

**Additional file 1:** Genetic and expression diversity within *Indica* and *Japonica* accessions.

**Additional file 2:** Assessing phenotypic variation and environmental effects for genes exhibiting genetic variability within each subspecies.

**Additional file 3:** Results of differential expression and presence-absence variation analyses.

**Additional file 4:** Heritability of gene expression for full population (*Indica* and *Japonica*).

**Additional file 5:** Results of heritability analyses within each subspecies.

**Additional file 6:** Genes showing differences in heritability between subspecies.

**Additional file 7:** cis-eQTL results.

### Abbreviations

MAF: minor allele frequency;  $\pi$ : nucleotide diversity; SNP: single nucleotide polymorphism; eQTL: expression quantitative trait loci; Ind: *Indica*; Jap: *Japonica*; BMA: Bayesian model averaging; DE: differentially expressed; PAV: presence-absence expression variation; FDR: false-discovery rate; CV: coefficient of variation; GO: gene ontology; PCA: principle component analysis; PC: principle component; RDP1: rice diversity panel 1

### Acknowledgements

Not applicable.

### Authors' contributions

Study was conceived by MTC under the supervision of HW; MTC designed the experiments; MTC and SS performed the experiments; MTC and SS collected and processed all plant materials; RNA-seq mapping and read quantification was done by QD and KL under the supervising of CZ; differential expression and genetic analyses were performed by MTC; MTC wrote the manuscript with input from HW; HW and CZ edited the manuscript. All authors have read and approved the manuscript.

### Funding

Funding for this research was provided by the National Science Foundation (United States) through Awards 1238125 and 1736192 to Harkamal Walia. Funds were used for the design of the study, as well as the collection and analyses of the data.

### Availability of data and materials

All sequencing data is available via NCBI Gene Expression Omnibus under the accession number GSE98455.

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interests

The authors declare that they have no competing interests.

### Author details

<sup>1</sup>Department of Agronomy and Horticulture, University of Nebraska Lincoln, 1825 N 38th St., 68583 Lincoln, NE, USA. <sup>2</sup>Department of Animal and Poultry Sciences, Virginia Polytechnic Institute and State University, 175 West Campus Drive, 24060 Blacksburg, VA, USA. <sup>3</sup>School of Biological Sciences, University of Nebraska Lincoln, 1901 Vine St., 68503 Lincoln, NE, USA. <sup>4</sup>Marine Biotechnology and Ecology Division, CSIR-CSMCRI, Bhavnagar, Gujarat, India.

Received: 5 December 2019 Accepted: 19 May 2020

Published online: 08 June 2020

### References

- Oka H, et al. Genetic diversity of wild and cultivated rice. *Rice Biotechnol.* 1991;55–81.
- Huang X, Kurata N, Wang Z-X, Wang A, Zhao Q, Zhao Y, Liu K, Lu H, Li W, Guo Y, et al. A map of rice genome variation reveals the origin of cultivated rice. *Nature.* 2012;490(7421):497.
- Kato A. On the affinity of rice varieties as shown by the fertility of rice plants. *Centr Agric Inst Kyushu Imp Univ.* 1928;2:241–76.
- Ding J, Araki H, Wang Q, Zhang P, Yang S, Chen J-Q, Tian D. Highly asymmetric rice genomes. *BMC Genomics.* 2007;8(1):154.
- Goff SA, Ricke D, Lan T-H, Presting G, Wang R, Dunn M, Glazebrook J, Sessions A, Oeller P, Varma H, et al. A draft sequence of the rice genome (*Oryza sativa* L. ssp. *japonica*). *Science.* 2002;296(5565):92–100.
- Yu J, Hu S, Wang J, Wong GK-S, Li S, Liu B, Deng Y, Dai L, Zhou Y, Zhang X, et al. *Science.* 2002;296(5565):79–92.
- Feltus FA, Wan J, Schulze SR, Estill JC, Jiang N, Paterson AH. An snp resource for rice genetics and breeding based on subspecies *indica* and *japonica* genome alignments. *Genome Res.* 2004;14(9):1812–9.
- Stein JC, Yu Y, Copetti D, Zwickl DJ, Zhang L, Zhang C, Chougule K, Gao D, Iwata A, Goicoechea JL, et al. Genomes of 13 domesticated and wild rice relatives highlight genetic conservation, turnover and innovation across the genus *Oryza*. *Nat Genet.* 2018;50(2):285.
- Koide Y, Ogino A, Yoshikawa T, Kitashima Y, Saito N, Kanaoka Y, Onishi K, Yoshitake Y, Tsukiyama T, Saito H, et al. Lineage-specific gene acquisition or loss is involved in interspecific hybrid sterility in rice. *Proc Natl Acad Sci.* 2018;115(9):1955–62.
- Schatz MC, Maron LG, Stein JC, Wences AH, Gurtowski J, Biggers E, Lee H, Kramer M, Antoniou E, Ghiban E, et al. Whole genome de novo assemblies of three divergent strains of rice, *Oryza sativa*, document novel gene space of *aus* and *indica*. *Genome Biol.* 2014;15(11):506.
- Huang X, Lu G, Zhao Q, Liu X, Han B. Genome-wide analysis of transposon insertion polymorphisms reveals intraspecific variation in cultivated rice. *Plant Physiol.* 2008;148(1):25–40.
- Wang X, Kudrna DA, Pan Y, Wang H, Liu L, Lin H, Zhang J, Song X, Goicoechea JL, Wing RA, et al. Global genomic diversity of *Oryza sativa* varieties revealed by comparative physical mapping. *Genetics.* 2014;196(4):937–49. <https://doi.org/10.1534/genetics.113.159970>.
- Walia H, Wilson C, Zeng L, Ismail AM, Condamine P, Close TJ. Genome-wide transcriptional analysis of salinity stressed *japonica* and *indica* rice genotypes during panicle initiation stage. *Plant Mol Biol.* 2007;63(5):609–23.
- Lu T, Lu G, Fan D, Zhu C, Li W, Zhao Q, Feng Q, Zhao Y, Guo Y, Li W, et al. Function annotation of rice transcriptome at single nucleotide resolution by rna-seq. *Genome research.* 2010;20(9):1238–49.
- Jung K-H, Gho H-J, Giong H-K, Chandran AKN, Nguyen Q-N, Choi H, Zhang T, Wang W, Kim J-H, Choi H-K, et al. Genome-wide identification and analysis of *japonica* and *indica* cultivar-preferred transcripts in rice using 983 affymetrix array data. *Rice.* 2013;6(1):19.
- Zhao K, Tung C-W, Eizenga GC, Wright MH, Ali ML, Price AH, Norton GJ, Islam MR, Reynolds A, Mezey J, et al. Genome-wide association mapping reveals a rich genetic architecture of complex traits in *Oryza sativa*. *Nat Commun.* 2011;2:467.
- Famoso AN, Zhao K, Clark RT, Tung C-W, Wright MH, Bustamante C, Kochian LV, McCouch SR. Genetic architecture of aluminum tolerance in rice (*Oryza sativa*) determined through genome-wide association analysis and QTL mapping. *PLoS Genet.* 2011;7(8):1002221.
- Eizenga GC, Ali M, Bryant RJ, Yeater KM, McClung AM, McCouch SR, et al. Registration of the rice diversity panel 1 for genomewide association studies. *J Plant Registrations.* 2014;8(1):109–16.
- Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics.* 2000;155(2):945–59.
- Caicedo AL, Williamson SH, Hernandez RD, Boyko A, Fledel-Alon A, York TL, Polato NR, Olsen KM, Nielsen R, McCouch SR, et al. Genome-wide patterns of nucleotide polymorphism in domesticated rice. *PLoS Genet.* 2007;3(9):163.
- Huang X, Sang T, Zhao Q, Feng Q, Zhao Y, Li C, Zhu C, Lu T, Zhang Z, Li M, et al. Genome-wide association studies of 14 agronomic traits in rice landraces. *Nat Genet.* 2010;42(11):961.
- Mather KA, Caicedo AL, Polato N, Olsen KM, McCouch S, Purugganan MD. The extent of linkage disequilibrium in rice (*oryza sativa* L.). *Genetics.* 2007;177(4):2223–32. <https://doi.org/10.1534/genetics.107.079616>.
- Garris AJ, Tai TH, Coburn J, Kresovich S, McCouch S. Genetic structure and diversity in *Oryza sativa* L. *Genetics.* 2005;169(3):1631–8.
- VanRaden PM. Efficient methods to compute genomic predictions. *J Dairy Sci.* 2008;91(11):4414–23.
- Mizuta Y, Harushima Y, Kurata N. Rice pollen hybrid incompatibility caused by reciprocal gene loss of duplicated genes. *Proc Natl Acad Sci.* 2010;107(47):20417–22.
- Flutre T, Wen X, Pritchard J, Stephens M. *PLoS Genet.* 2013;9(5):1003486.
- Wang W, Mauleon R, Hu Z, Chebotarov D, Tai S, Wu Z, Li M, Zheng T, Fuentes RR, Zhang F, et al. Genomic variation in 3,010 diverse accessions of asian cultivated rice. *Nature.* 2018;557(7703):43.

28. Mansueto L, Fuentes RR, Borja FN, Detras J, Abriol-Santos JM, Chebotarov D, Sanciangco M, Palis K, Copetti D, Poliakov A, et al. Rice snp-seek database update: new snps, indels, and queries. *Nucleic Acids Res.* 2016;45(D1):1075–81.
29. Mansueto L, Fuentes RR, Chebotarov D, Borja FN, Detras J, Abriol-Santos JM, Palis K, Poliakov A, Dubchak I, Solovyev V, et al. SNP-Seek II: A resource for allele mining and analysis of big genomic data in *Oryza sativa*. *Curr Plant Biol.* 2016;7:16–25.
30. Alexandrov N, Tai S, Wang W, Mansueto L, Palis K, Fuentes RR, Ulat VJ, Chebotarov D, Zhang G, Li Z, et al. Snp-seek database of snps derived from 3000 rice genomes. *Nucleic Acids Res.* 2014;43(D1):1023–7.
31. Terao H, Mizushima U. Some considerations on the classification of *Oryza sativa* L. into two subspecies, so called Japonica and Indica. *Jpn J Bot.* 1942;10:213–58.
32. Matsuo T. Genecological studies on cultivated rice. *Bull Natl Inst Agr Sci Jpn D.* 1952;3:1–111.
33. Morinaga T. Classification of rice varieties on the basis of affinity. *Jpn J Breed.* 1954;4:1–14.
34. Morishima H, Oka H-I. Phylogenetic differentiation of cultivated rice, xxii. numerical evaluation of the indica-japonica differentiation. *Jpn J Breed.* 1981;31(4):402–13.
35. Glaszmann J-C. Isozymes and classification of asian rice varieties. *Theor Appl Genet.* 1987;74(1):21–30.
36. Wang DR, Agosto-Pérez FJ, Chebotarov D, Shi Y, Marchini J, Fitzgerald M, McNally KL, Alexandrov N, McCouch SR. An imputation platform to enhance integration of rice genetic resources. *Nat Commun.* 2018;9(1):1–10.
37. Xu X, Liu X, Ge S, Jensen JD, Hu F, Li X, Dong Y, Gutenkunst RN, Fang L, Huang L, et al. Resequencing 50 accessions of cultivated and wild rice yields markers for identifying agronomically important genes. *Nat Biotechnol.* 2012;30(1):105.
38. Price AL, Helgason A, Thorleifsson G, McCarroll SA, Kong A, Stefansson K. Single-tissue and cross-tissue heritability of gene expression via identity-by-descent in related or unrelated individuals. *PLoS Genet.* 2011;7(2):1001317.
39. Grundberg E, Small KS, Hedman ÅK, Nica AC, Buil A, Keildson S, Bell JT, Yang T-P, Meduri E, Barrett A, et al. Mapping cis-and trans-regulatory effects across multiple tissues in twins. *Nat Genet.* 2012;44(10):1084.
40. Hore V, Viñuela A, Buil A, Knight J, McCarthy MI, Small K, Marchini J. Tensor decomposition for multiple-tissue gene expression experiments. *Nat Genet.* 2016;48(9):1094.
41. Takano M, Inagaki N, Xie X, Yuzurihara N, Hihara F, Ishizuka T, Yano M, Nishimura M, Miyao A, Hirochika H, et al. Distinct and cooperative functions of phytochromes a, b, and c in the control of deetiolation and flowering in rice. *Plant Cell.* 2005;17(12):3311–25.
42. Lee Y-S, Yi J, An G. OsphyA modulates rice flowering time mainly through osgi under short days and ghd7 under long days in the absence of phytochrome b. *Plant Mol Biol.* 2016;91(4-5):413–27.
43. Kim S-K, Yun C-H, Lee JH, Jang YH, Park H-Y, Kim J-K. Osco3, a constans-like gene, controls flowering by negatively regulating the expression of ft-like genes under sd conditions in rice. *Planta.* 2008;228(2):355.
44. Song YH, Shim JS, Kinmonth-Schultz HA, Imaizumi T. Photoperiodic flowering: time measurement mechanisms in leaves. *Annu Rev Plant Biol.* 2015;66:441–64.
45. Huang C-L, Hung C-Y, Chiang Y-C, Hwang C-C, Hsu T-W, Huang C-C, Hung K-H, Tsai K-C, Wang K-H, Osada N, et al. Footprints of natural and artificial selection for photoperiod pathway genes in *Oryza*. *Plant J.* 2012;70(5):769–82.
46. Itoh H, Tatsumi T, Sakamoto T, Otomo K, Toyomasu T, Kitano H, Ashikari M, Ichihara S, Matsuoka M. A rice semi-dwarf gene, tan-ginbozu (d35), encodes the gibberellin biosynthesis enzyme, ent-kaurene oxidase. *Plant Mol Biol.* 2004;54(4):533–47.
47. Naranjo L, Talón M, Domingo C. Diversity of floral regulatory genes of japonica rice cultivated at northern latitudes. *BMC Genomics.* 2014;15(1):101.
48. Itoh H, Ueguchi-Tanaka M, Sentoku N, Kitano H, Matsuoka M, Kobayashi M. Cloning and functional analysis of two gibberellin 3 $\beta$ -hydroxylase genes that are differently expressed during the growth of rice. *Proc Natl Acad Sci.* 2001;98(15):8909–14.
49. Yuan J, Chen D, Ren Y, Zhang X, Zhao J. Characteristic and expression analysis of a metallothionein gene, osmt2b, down-regulated by cytokinin suggests functions in root development and seed embryo germination of rice. *Plant Physiol.* 2008;146(4):1637–50.
50. Yoshida S, Forno D, Cock J, Gomez K. Laboratory manual for physiological studies of rice, 3rd edn Manila: International Rice Research Institute. 1976.
51. Andrews S, et al. FastQC: a quality control tool for high throughput sequence data. 2010.
52. Trapnell C, Pachter L, Salzberg SL. Tophat: discovering splice junctions with rna-seq. *Bioinformatics.* 2009;25(9):1105–11.
53. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for illumina sequence data. *Bioinformatics.* 2014;30(15):2114–20.
54. Anders S, Pyl PT, Huber W. Htseq—a python framework to work with high-throughput sequencing data. *Bioinformatics.* 2015;31(2):166–9.
55. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for ma-seq data with deseq2. *Genome Biol.* 2014;15(12):550.
56. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. *J Stat Softw.* 2015;67(1):1–48. <https://doi.org/10.18637/jss.v067.i01>.
57. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc Ser B Methodol.* 1995;57(1):289–300. <https://doi.org/10.1111/j.2517-6161.1995.tb02031.x>.
58. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G, Marth GT, Sherry ST, et al. The variant call format and vcfutils. *Bioinformatics.* 2011;27(15):2156–2158.
59. Scheipl F, Greven S, Kuechenhoff H. Size and power of tests for a zero random effect variance or polynomial regression in additive and linear mixed models. *Comput Stat Data Anal.* 2008;52(7):3283–99.
60. Butler D, Cullis BR, Gilmour A, Gogel B. *Asreml-r reference manual.* 2009. The State of Queensland, Department of Primary Industries and Fisheries, Brisbane.
61. McCouch SR, Wright MH, Tung C-W, Maron LG, McNally KL, Fitzgerald M, Singh N, DeClerck G, Agosto-Perez F, Korniliev P, et al. Open access resources for genome-wide association mapping in rice. *Nat Commun.* 2016;7:10532.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

### Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

