Journal Pre-proof

BjulR: A Multi-omics Database with Various Tools for Accelerating Functional Genomics Research in *Brassica juncea*

Linna Zhang, Jinyuan Xiao, Congyuan Liang, Yifan Chen, Changchun Yu, Xinle Zhao, Jiawei Li, Mingli Yan, Qian Yang, Hao Chen, Zhongsong Liu, Zhengjie Wan, Zhiquan Yang, Qing-Yong Yang

PII: S2590-3462(24)00195-0

DOI: https://doi.org/10.1016/j.xplc.2024.100925

Reference: XPLC 100925

To appear in: PLANT COMMUNICATIONS

Please cite this article as: Zhang, L., Xiao, J., Liang, C., Chen, Y., Yu, C., Zhao, X., Li, J., Yan, M., Yang, Q., Chen, H., Liu, Z., Wan, Z., Yang, Z., Yang, Q.-Y., BjulR: A Multi-omics Database with Various Tools for Accelerating Functional Genomics Research in *Brassica juncea*, *PLANT COMMUNICATIONS* (2024), doi: https://doi.org/10.1016/j.xplc.2024.100925.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2024



	1100		D		5	\sim		
	սու	lai			D	U	υ	

1	BjuIR: A Multi-omics Database with Various Tools for							
2	Accelerating Functional Genomics Research in Brassica juncea							
3								
4	Linna Zhang ^{1,2,3,†} , Jinyuan Xiao ^{1,2,†} , Congyuan Liang ^{1,2,3,†} , Yifan Chen ^{1,2,4} , Changchun							
5	Yu ⁵ , Xinle Zhao ^{1,2} , Jiawei Li ^{1,2} , Mingli Yan ⁶ , Qian Yang ⁶ , Hao Chen ⁷ , Zhongsong Liu ⁷ ,							
6	Zhengjie Wan ^{5,*} , Zhiquan Yang ^{1,2,4,*} and Qing-Yong Yang ^{1,2,3,4,*}							
7								
8	¹ National Key Laboratory for Germplasm Innovation & Utilization of Horticultural Crops							
9	College of Informatics, Huazhong Agricultural University, Wuhan, 430070, China							
10	² Hubei Key Laboratory of Agricultural Bioinformatics and Hubei Engineering Technology							
11	Research Center of Agricultural Big Data, College of Informatics, Huazhong Agricultural							
12	University, Wuhan 430070, China							
13	³ National Key Laboratory of Crop Genetic Improvement, Hubei Hongshan Laboratory,							
14	Huazhong Agricultural University, Wuhan 430070, China							
15	⁴ Yazhouwan National Laboratory, Sanya 572025, China							
16	⁵ National Key Laboratory for Germplasm Innovation & Utilization of Horticultural Crops,							
17	College of Horticulture & Forestry Sciences, Huazhong Agricultural University, Wuhan							
18	430070, China							
19	⁶ Crop Research Institute, Hunan Academy of Agricultural Sciences, Changsha 410125, China							
20	⁷ College of Agronomy, Hunan Agricultural University, Changsha 410128, China							
21	[†] These authors contributed equally to this article.							
22	*Correspondence: Zhengjie Wan (<u>wanzj@mail.hzau.edu.cn</u>), Zhiquan Yang							
23	(<u>yang_zq@foxmail.com</u>), Qing-Yong Yang (<u>yqy@mail.hzau.edu.cn</u>)							

24 Dear Editor,

Advancements in high-throughput omics technologies, along with methodologies for 25 integrating multi-omics datasets, have substantially enhanced the efficiency of 26 identifying candidate genes in breeding (Gupta et al., 2019; Gusev et al., 2018). 27 However, this process is often complex and laborious. To address this challenge, 28 databases that integrate extensive data and enable convenient and efficient function 29 genomics studies are being developed (Ma et al., 2021; Yang et al., 2023). Brassica 30 31 juncea (B. juncea), commonly known as mustard, is an economically significant agricultural species for its diverse uses, including vegetables, resilient oilseeds, and 32 distinctively flavored condiments (Yang et al., 2018). This diversity of applications has 33 34 spurred the accumulation of substantial multi-omics data in fundamental research of mustard, yet there lacks a specialized platform to harness these data fully for mustard's 35 genetic improvement. Addressing this gap, we have developed BjuIR (Brassica juncea 36 Information Resource, available at https://yanglab.hzau.edu.cn/BjuIR), integrating the 37 38 most comprehensive mustard omics datasets to date from over 2,000 accessions, including data from genomics, variomics, transcriptomics, phenomics, and 39 metabolomics. BjuIR provides sophisticated analyses for these multi-omics datasets of 40 mustard with user-friendly interfaces, enabling rapid querying of "variant/gene 41 expression-phenotype" associations for the quick identification of candidate genes and 42 greatly benefiting functional genomics research. 43

44 DATA AND FUNCTIONAL MODULES IN BjuIR

BjuIR boasts a rich repository of large-scale multi-omics datasets (Figure 1A), 45 encompassing 40 genome assemblies (Supplemental Table 1), 8,869,856 single 46 nucleotide polymorphisms (SNPs) and short insertions/deletions (InDels) across 1,614 47 accessions (Supplemental Figure 1 and Supplemental Table 2), 941 RNA-seq libraries 48 (Supplemental Table 3), 412 metabolites (Supplemental Table 4), phenotypic data of 49 628 accessions spanning 16 traits (Supplemental Table 5), and 1,841 mustard-centric 50 51 literature entries. Various analysis methods were applied to fully explore the value of these datasets and the analysis results are organized and accessible in eight modules 52 within BjuIR (Figure 1B). 53

54 The "Genomics" module provides queries for syntenic relationships between genomes and gene annotation; the "Population" module details accession information 55 and provides selective signals of populations; the "Variations" module allows for the 56 exploration of variants, and their associations with phenotypes and gene expression 57 58 levels; The "Transcriptomics" module features gene expression profiles, co-expression networks and differential expression analysis; the "Phenomics" module presents 59 phenotype data; the "Metabolomics" module provides metabolite information; The 60 "Multi-omics" module facilitates quick queries for "variation-gene expression-61 phenotype" associations generated from genome-wide association study (GWAS), 62 transcriptome-wide association analysis (TWAS), expression quantitative trait loci 63 (eQTL) mapping analysis, colocalization analysis, and summary-based mendelian 64 randomization (SMR); and the "Literature" module supports literature studies on 65 66 mustard. Each module is equipped with user-friendly interfaces for result visualization, data downloads, and seamless navigation among modules and other databases. 67

3

68 APPLICATIONS AND ANALYSIS TOOLS IN BjuIR

69 Comprehensive datasets and thoughtfully designed modules in BjuIR are 70 of great utility for those involved in functional genomics and candidate 71 genes/variations identification efforts.

In the "Genomics" module, users can visualize global genome alignments in a dotplot (Figure 1C) and local genome alignments in Gbrowse in the "Genome synteny" interface (Supplemental Figure 2A) and query homologous gene clusters by entering gene ID or gene name in the "Gene cluster" interface (Supplemental Figure 2B). Additionally, they can explore annotations of gene families and biological pathways for mustard genes in the "Gene family" and "Pathway" interfaces (Supplemental Figure 2C and 2D).

79 The "Transcriptomics" module offers gene expression profiles across populations and tissues, co-expression analysis, and comparative transcription analysis, as 80 demonstrated in Supplemental Figure 3. These capabilities are exemplified through the 81 82 gene BjuA09.WRI1. By entering the gene BjuVA09G4986 (WRI1) in the "Tissue expression profile" and "eFP" interface, users can access the tissue expression profiles 83 of BjuA09.WRI1 and its homologous genes, presented in a heatmap (Supplemental 84 Figure 3A) and eFP viewer (Figure 1D). Users can also explore gene-gene and lncRNA-85 mRNA co-expression networks related to BjuA09. WRI1 in the "Co-expression network" 86 interface (Supplemental Figure 3B and 3C). Moreover, they can access differentially 87 expressed genes/lncRNAs in the "Differential expression" interface (Supplemental 88 Figure 3D and 3E). 89

The "Population" module provides queries for detailed information of 1,614 accessions, including their subpopulations, usages, and origins (Supplemental Figure 4 and Supplemental Table 2). Furthermore, this module enables the querying of selective signals, such as π , Tajima's D, F_{ST} , and cross-population extended haplotype homozygosity (XP-EHH), which are calculated using variations to identify candidate regions and genes that may be under selection (Supplemental Figure 5).

Journal Pre-proot

The "Variations" module provides detailed information on genetic variants and 96 assessment of how specific variants or haplotypes affect phenotypes and gene 97 expression (Supplemental Figure 6). For instance, by inputting the gene ID 98 "BjuVB08G59610" into the "Variations/Single-locus model" interface, the result page 99 displays annotations for all SNPs and InDels within the gene region (Supplemental 100 Figure 6A-6D). Choosing a particular variant, such as "BB Chr08:62443776" 101 (Supplemental Figure 6D), users can examine its allele frequency across diverse 102 subpopulations and geographical locations (Supplemental Figure 6E and 6F), and 103 explore its correlation with phenotypic traits, such as thousand seed weight, as well as 104 with gene expression levels (Supplemental Figure 6G and 6H). 105

Integrated analyses of multi-omics data in the "Multi-omics" module vastly improve 106 the efficiency of candidate gene discovery in B. juncea. Here, users can submit a gene 107 name, ID, or trait name to unveil associations between variations, gene expression, and 108 phenotypes. This module includes "variation-trait" associations identified by GWAS 109 (Supplemental Figure 7A and Supplemental Table 6), "variation-gene expression" 110 111 associations from eQTL (Supplemental Figure 7B and Supplemental Table 7), "gene expression-trait" associations identified by TWAS (Supplemental Figure 7C 112 Supplemental Table 8), as well as colocalization analyses (Supplemental Figure 7D-7F 113 and Supplemental Table 9). The reliability of these integrated results is demonstrated 114 by a reproducibility rate of 60.42% in comparison to prior findings (Harper et al., 2020), 115 as listed in Supplemental Table 10. Additionally, the "variation-gene expression-trait" 116 117 associations are also viewable in a network format, exemplified by entering the "BIM1" gene within the "Multi-omics/Association networks" interface, which showcases all 118 related associations in a visual network (Figure 1E). 119

The "Literature" module offers advanced search capabilities based on keywords, journal names, and publication years, allowing users to efficiently access research advancements related to mustard in a specific field. For instance, by entering the keyword "flowering" in this module, users can retrieve relevant literature on "flowering", with detailed information displayed in a table. In addition, the statistics of

5

Journal Pre-proof

the literature, organized by publication year or journal, are visually presented in line
graphs and bar charts. An additional feature is the provision of a keyword co-occurrence
network for visualizing trends in studies related to the flowering of mustard
(Supplemental Figure 8).

The "Tools" module incorporates 15 essential bioinformatics analysis tools applications, supporting user-initiated analyses, including Gene Ontology (GO) enrichment, linkage disequilibrium (LD) calculations, SNP matching for germplasm identification, sequence extraction, primer design, among other functions (Supplemental Figure 9).

134 CASE STUDY: MINING NOVEL CANDIDATE VARIANTS AND GENES 135 ASSOCIATED WITH TOCOPHEROL CONTENT USING BjuIR

We illustrate the utility of BjuIR in mining candidate genes and variations using the 136 example of tocopherol, a crucial vitamin E component vital for seed quality and human 137 nutrition. Tocopherol exists in various forms, such as a-tocopherol, which is recognized 138 139 for having the highest vitamin E activity in mammals and can be derived from γ tocopherol (Tucker and Townsend, 2005). Initiating with the query " γ -/ α -tocopherol 140 content in seed" in the "Multi-omics/GWAS" interface, the analysis rendered a 141 Manhattan plot revealing two genomic loci associated with γ -/ α -tocopherol content 142 located on chromosomes AA Chr02 and AA Chr06. While the AA Chr02 locus had 143 been previously reported (Harper et al., 2020), the locus on AA Chr06 emerged as a 144 novel finding, comprising an LD block between 5.88 to 5.95 Mb (Figure 1F and 1G), 145 within which seven genes reside (Figure 1H and Supplemental Table 11). Specifically, 146 BjuVA06G10820, notable for containing the largest number of GWAS-SNPs in its 147 coding and 3 kb upstream regions, was identified. Its homolog in Arabidopsis thaliana 148 (AT1G15125) is known to code for S-adenosyl-L-methionine-dependent 149 methyltransferases that participate in converting γ -tocopherol to α -tocopherol (Tavva 150 et al., 2007), suggesting that BjuVA06G10820 could be a prime candidate gene within 151 the AA Chr06 locus. Further, haplotype analysis of BjuVA06G10820 through the 152 "Variations" module's "Single-locus model" interface uncovered two prevalent 153

Journal Pre-proot

haplotypes (Figure 1I). Notably, accessions carrying Haplotype1 were significantly associated with reduced γ -/ α -tocopherol content compared to those with Haplotype2, delineated in Figure 1J. Such findings will provide a valuable reference for future breeding strategies aimed at boosting α -tocopherol levels in mustard seeds and underscore BjuIR's aptitude for identifying candidate genes and variants associated with specific traits.

In conclusion, BjuIR stands as the most extensive and comprehensive multi-omics 160 161 database to date for functional genomics research in mustard. Its key features include (1) expedited access to each omics dataset and complete analysis results; (2) quick 162 mining of candidate genes and variants via robust "variant-gene expression-phenotype" 163 164 associations; (3) multiple user-friendly, online bioinformatic tools; and (4) navigationfriendly interfaces for efficient data mining. With its rich database and thoughtful 165 design, BjuIR proves to be a highly efficient and convenient platform for functional 166 genomics research and candidate gene identification. Looking ahead, BjuIR will persist 167 168 in incorporating novel omics data, reinforcing its status as an indispensable platform for furthering functional genomics and genetic improvement in mustard. 169

170 **D**A

DATA AVAILABILITY

Sources of all datasets are described in supplemental information. All datasets are
available at <u>https://yanglab.hzau.edu.cn/BjuIR/download</u>.

173 FUNDING

This research was supported by the National Natural Science Foundation of China (32322061 and 32070559); the National Key Research and Development Plan of China (2021YFF1000100); the Fundamental Research Funds for the Central University HZAU (2662023XXPY001); the Hubei Hongshan Laboratory (2021HSZD004); and the Developing Bioinformatics Platform in Hainan Yazhou Bay Seed Lab (no. JBGS-B21HJ0001).

180 AUTHOR CONTRIBUTIONS

- 181 Q.-Y.Y., Z.Y., and Z.W. designed the project. L.Z. and Y.C. collected the datasets. L.Z.,
- 182 C.L., Y.C., J.X., and J.L. performed the bioinformatics analysis. J.X., Y.C., and X.Z.
- developed the BjuIR database. C.Y. provided the pictures of *B. juncea* germplasms.
- 184 L.Z., C.L., Z.Y., and Q.-Y.Y. wrote the manuscript. Q.-Y.Y., Z.Y., Z.W., Z.L., H.C., Q.Y.,
- and M.Y. directed the project. All authors read and approved the manuscript.

Johnaldre

ACKNOWLEDGMENTS

- 186 We thank the bioinformatics computing platform of the National Key Laboratory of
- 187 Crop Genetic Improvement, Huazhong Agricultural University, managed by Hao Liu.
- 188 No conflict of interest is declared.

8

189 REFERENCE

- Gupta, P.K., Kulwal, P.L., and Jaiswal, V. (2019). Association mapping in plants in
 the post-GWAS genomics era. Adv. Genet. 104: 75-154.
- Gusev, A., Mancuso, N., Won, H., Kousi, M., Finucane, H.K., Reshef, Y., Song, L.,
 Safi, A., Schizophrenia Working Group of the Psychiatric Genomics
 Consortium, and McCarroll, S., et al. (2018). Transcriptome-wide association
 study of schizophrenia and chromatin activity yields mechanistic disease
 insights. Nat. Genet. 50: 538-548.
- Harper, A.L., He, Z., Langer, S., Havlickova, L., Wang, L., Fellgett, A., Gupta, V.,
 Kumar Pradhan, A., and Bancroft, I. (2020). Validation of an associative
 transcriptomics platform in the polyploid crop species *Brassica juncea* by
 dissection of the genetic architecture of agronomic and quality traits. Plant J.
 103: 1885-1893.
- Ma, S., Wang, M., Wu, J., Guo, W., Chen, Y., Li, G., Wang, Y., Shi, W., Xia, G., Fu,
 D., et al. (2021). WheatOmics: A platform combining multiple omics data to
 accelerate functional genomics studies in wheat. Mol. Plant 14: 1965-1968.
- Tavva, V.S., Kim, Y.-H., Kagan, I.A., Dinkins, R.D., Kim, K.-H., and Collins, G.B.
 (2007). Increased α-tocopherol content in soybean seed overexpressing the
 Perilla frutescens γ-tocopherol methyltransferase gene. Plant Cell Rep. 26: 61 70.
- Tucker, J., and Townsend, D. (2005). Alpha-tocopherol: roles in prevention and
 therapy of human disease. Biomed. Pharmacother. 59: 380-387.
- Yang, J., Zhang, C., Zhao, N., Zhang, L., Hu, Z., Chen, S., and Zhang, M. (2018).
 Chinese root-type mustard provides phylogenomic insights into the evolution
 of the multi-use diversified allopolyploid *Brassica juncea*. Mol. Plant 11: 512514.
- Yang, Z.Q., Wang, S.B., Wei, L.L., Huang, Y.M., Liu, D.X., Jia, Y.P., Luo, C.F., Lin,
 Y.C., Liang, C.Y., Hu, Y., et al. (2023). BnIR: A multi-omics database with
 various tools for research and breeding. Mol. Plant 16: 775-789.

Journal Pre-proof

Figure 1. Overview of BjuIR. (A) Large-scale datasets collected in BjuIR. (B) Eight 218 modules and their functions in BjuIR. (C) Comparative genomics analysis between 219 T84-66.V2.0 and AU213.V1.0 genome in the "Genomics" module. (D) eFP viewer 220 displaying tissue-specific expression profiles of the gene BjuVA09G4986 in the 221 "Transcriptomics" module. (E) "Variation-gene expression-phenotype" associations 222 related to gene BIM1 in the "Multi-omics" module. (F-J) Identification of novel 223 candidate genes/variants associated with γ -/ α -tocopherol content in BjuIR. (F) GWAS 224 225 for γ -/ α -tocopherol content in seed. The *P*-value threshold was set at 6.54e-6 based on 1/n, where n represents the number of independent SNPs (n =152,884). (G) Local 226 Manhattan plot of GWAS for γ -/ α -tocopherol content and heatmap of linkage 227 disequilibrium (LD) blocks. The color of the dots represents the degree of LD with the 228 lead SNP. (H) Genes in the LD block are significantly associated with γ -/ α -tocopherol 229 content. (I) Haplotypes formed by combinations of GWAS-SNPs on the coding region 230 of BjuVA06G10820 and its 3 kb upstream flanking region. (J) Comparison of γ -/ α -231 tocopherol content between accessions with different haplotypes. * indicates P < 0.05232 233 (Wilcoxon rank sum test).

Iournal Pre-proof



r² color kev

1.0

0.0

BjuVA06G10790

5

1 29E-06

Haplotype1 Haplotype2 (129) (9)

AT1G67800