# Journal Pre-proof

Smart Breeding Platform: a web-based tool for high-throughput population genetics, phenomics, and genomic selection

Huihui Li, Xin Li, Peng Zhang, Yingwei Feng, Junri Mi, Shang Gao, Lele Sheng, Mohsin Ali, Zikun Yang, Liang Li, Wei Fang, Wensheng Wang, Qian Qian, Fei Gu, Wenbin Zhou

Please cite this article as: Li H., Li X., Zhang P., Feng Y., Mi J., Gao S., Sheng L., Ali M., Yang Z., Li L., Fang W., Wang W., Qian Q., Gu F., and Zhou W. (2024). Smart Breeding Platform: a web-based tool for high-throughput population genetics, phenomics, and genomic selection. Mol. Plant. doi: https://doi.org/10.1016/j.molp.2024.03.002.

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.

1    **Smart Breeding Platform: a web-based tool for high-throughput population genetics,**

2    **phenomics, and genomic selection**

3    Huihui Li[1,2,#,*], Xin Li[3,4,#], Peng Zhang[3,4], Yingwei Feng[1,2], Junri Mi[3,4], Shang Gao[1,2], Lele

4    Sheng[3,4], Mohsin Ali[1,2], Zikun Yang[3,4], Liang Li[1], Wei Fang[1], Wensheng Wang[1,2], Qian

5    Qian[1,2,5,6], Fei Gu[3,4,*], Wenbin Zhou[1]

6

7    [1] State Key Laboratory of Crop Gene Resources and Breeding, Institute of Crop Sciences,

8    Chinese Academy of Agricultural Sciences (CAAS), Beijing, China

9    [2] Nanfan Research Institute, CAAS, Sanya, Hainan, China

10   [3] DAMO Academy, Alibaba Group, Hangzhou, China

11   [4] Hupan Lab, Hangzhou, China

12   [5] Yazhouwan National Laboratory, No. 8 Huanjin Road, Yazhou District, Sanya City, Hainan

13   Province 572024, China

14   [6] State Key Laboratory of Rice Biology, China National Rice Research Institute, Hangzhou,

15   Zhejiang, China

16   * Corresponding author: Huihui Li, email: lihuihui@caas.cn

17                                   Fei Gu, email: gufei.gf@alibaba-inc.com

18   # These authors contributed equally

19

20   **Runing title**: Smart Breeding Platform

21   In the era of big data and artificial intelligence, "smart breeding" has become a broad

22   conceptual framework encompassing the paradigm shift of crop breeding to relying on

23   analysis of high-throughput population genetics and phenomics data to conduct genomic

24   selection, allowing identification and optimal use of the genetic potential in crop species

25   (Sharma et al., 2022; Xiao et al., 2022; Xu et al., 2022). Most existing tools for analyzing

26   high-throughput breeding data require extensive computational power, complex installation

27   processes, and command-line expertise, and are therefore challenging and inconvenient

28   for the majority of researchers and breeders (Brandies and Hogg, 2021). To overcome

29   these limitations, we developed Smart Breeding Platform (https://sbp.ibreed.cn), a user-

30    friendly, web-based tool for management and analysis of large-scale genetic, genomic, and

31    phenomic data. This platform is freely accessible through the internet and allows users to

32    import data, perform various statistical analyses, and conduct genome-wide association

33    studies and genomic selection using both classical machine learning and deep-learning

34    models. It will enable plant breeders to easily conduct the following steps: (1) efficiently

35    record, manage, and process raw phenotypic and genotypic data; (2) perform phenotypic

36    and population genetic analyses in highly customizable ways; and (3) easily conduct

37    GWAS and genomic selection using classical machine learning and deep-learning models.

38    Smart Breeding Platform contains four main sections (Figure 1): (1) Germplasm Data

39    Management, (2) Test Management, (3) Genomic Data Management, and (4) Data

40    Analysis. Each section is described in detail below.

41

42    **Germplasm Data Management**

43        This section contains tables for storing germplasm and intermediate material data and

44    pedigree data. The tables are directly editable and function similarly to a standard

45    spreadsheet.

46    *Germplasm and intermediate material table*

47        In the germplasm and intermediate material table, germplasm metadata can be

48    entered in a single row each, with the following parameters (columns) provided by default:

49    germplasm name, year of seed production, breeding station, storage location, quantity

50    harvested, quantity currently available, and serial number. Additional parameters, including

51    images, can be added and customized by the user. The data in each row are directly

52    exportable. Advanced lines can be promoted, allowing users to track advancement choices

53    over time. A summary bar graph shows the number of plant lines across years, generations,

54    and storage locations. All tables support efficient row-level filtering and sorting for quick

55    data retrieval.

56    *Pedigree module*

57        The pedigree module contains pedigree records for existing germplasm. Information

58    about parents and offspring can be viewed in either table or graph form. Graphs can be

59    used to visualize the lineage of one or more lines, including inbred lines, hybrids, or both.

60  Graphs can also be customized to show information for only parental lines or progeny. In

61  instances involving more than one generation, the depth of the visualized pedigree can

62  also be adjusted. The entire pedigree for a breeding program can be displayed as a

63  network graph that highlights the most popular lines.

64  *Location module*

65  By default, this table includes the following fields: year, season, breeding station,

66  location name, longitude, latitude, size (e.g., the number of rows in a field trial), and

67  environmental factors such as maturity zone, soil characteristics, and agronomic practices.

68  These data can be used to compare variables between sites and track relevant factors that

69  may contribute to phenotypic outcomes.

70  *Warehouse in–out module*

71  This module is used to track all seeds that enter and exit a specific research station.

72  Users can log events such as seed allocation for yield trials, seed transfers to other

73  breeders, or receipt of new germplasm. Relevant data including seed quantity and the time

74  and date of transfers can be included. A bar chart allows users to examine changes in seed

75  stocks over time.

76  In summary, the Germplasm Data Management section allows for efficient and intuitive

77  storage and analysis of metadata for all germplasm used in a breeding program.

78  **Test Management**

79  Broadly, this section stores data about all field testing and nursery locations. Details

80  of the two specific modules, the field testing module and the crossing nursery module, are

81  included below.

82  *Field testing module*

83  This section includes tools for management of field testing. It contains functionality for

84  specifying relevant lines, experimental designs, replicate numbers, and field layouts.

85  Seven commonly-used experimental designs are included: completely randomized design,

86  randomized complete block design (RCBD), augmented design, spatial design, sparse

87  design, alpha-lattice design, and row-column design. After a design is selected by the user,

88  the module can be used to automatically generate a suitable field layout based on the

89  number of entries in the experiment and the number of available field plots. By including

90  information about the physical location of a trial, multiple experiments from different

91  breeders can be automatically placed within a single field. The field layout is output as a

92  table containing the coordinates of each plot within the field. This module also includes

93  heatmaps, which show the distribution of values for each trait across the field; a stability

94  analysis, which shows the performance of specific lines across locations; and a testing

95  history, demonstrating the trials and locations in which a specific line has been tested.

96  *Crossing nursery*

97  The crossing nursery module can be used to plan new lines and pedigrees. The user

98  selects sets of female and male parental lines. The module then generates a crossing

99  matrix, with options for user input regarding specific cross combinations (e.g., crossing

100  patterns) and harvest instructions (for each row or plant). The module auto-generates

101  inventory entries to be added to the germplasm table and adds the pedigree of each cross

102  combination to the pedigree record table.

103  The Test Management section has tools to track experimental locations and plant

104  research materials (i.e., seeds) with ease. It allows researchers to easily visualize available

105  stock and to plan field experiments and crosses. Intuitive organization of these resources

106  in a single location enables researchers to focus on planning and conducting high-level

107  experiments.

108  **Genomic Data Management**

109  In this section, users can easily upload and manage all genomic sequencing data,

110  reference genome files, and genomic variant files. The data stored in this section can then

111  be used in the Data Analysis section.

112  **Data Analysis**

113  *Phenotypic statistical analysis module*

114  This module is used for analyzing high-throughput phenotypic data collected in the

115  field. Multi-year, multi-location data can be extracted directly from the field testing module

116  or can be uploaded separately. Based on the experimental design, the module can be used

117  to fit a mixed linear model (MLM) to calculate best linear unbiased estimation (BLUE) or

118  best linear unbiased prediction (BLUP) (Bates et al., 2015). The model fits two-dimensional

119  spatial patterns for spatial designs to account for soil heterogeneity (Covarrubias-Pazaran,

120  2016; Rodriguez-Alvarez et al., 2018). Data can be analyzed separately for each location

121  or as an integrated dataset including points from all locations. For each genotype, the

122  module outputs the BLUP and BLUE values of the included traits. Variance components,

123  heritability, and trait correlations can also be calculated. Entry-mean heritability and plot-

124  mean heritability of each trait are derived from the variance components of random models.

125  These two metrics can help breeders to assess the precision of trait values both at single-

126  plot level and across locations.

127  This module also automatically calculates correlations for all pairs of phenotypic traits.

128  The phenotypic correlation between each pair of traits is calculated as the Pearson

129  correlation coefficient of the raw phenotypic data, whereas the genetic correlation between

130  each pair of traits is calculated as the correlation of genetic effects in a model fitting both

131  traits and residual correlation effects (Muñoz and Sanchez, 2020). Examination of trait

132  correlations enables breeders to identify traits that can be bred independently (i.e., traits

133  that have low correlations with other traits) and traits that must be separated or bred jointly

134  (i.e., traits that have strong positive or negative correlations with other traits). For each

135  linear model, the goodness of fit and the validity of the residual normality assumption can

136  be assessed using diagnostic plots, including raw data distribution histograms, residual

137  histograms, plots showing residual compared to fitted values, and residual Q–Q plots.

138  Breeders can then select the best lines (those with ideal values across traits) using

139  scatterplots that display the distribution and correlation of BLUP or BLUE values for pairs

140  of traits. Overall, this module includes advanced single-trait and multi-trait analyses that

141  can be conducted in an automated, user-friendly manner.

142  *Genetic variation analysis module*

143  This module can be used to efficiently identify genetic variants based on high-

144  throughput genome sequencing data. In comparison to the standard pipeline for sequence

145  alignment and germline variant-calling analysis (BWA+GATK) (Yin et al., 2021), the

146  analysis method used here is significantly faster on our platform, due to the boosted tools

147  with novel acceleration algorithm on the NVIDIA CUDA platform. Results of the new

148  method are highly consistent with the standard BWA+GATK pipeline (99.9% accuracy) and

149  are completed ~100× faster when two NVIDIA Turing T4 graphics cards are used. Inclusion

150  of additional graphics cards would further improve the processing speed. The sequence

151  alignment and sequencing depth can be visualized with Integrative Genomics Viewer (IGV)

152  (Robinson et al., 2011), which has been optimized to load large genome dataset.

153  *Genomic statistical analysis module*

154  This module facilitates analyses of genetic diversity for a specific population. It takes SNP

155  data as input, either as VCF files produced by the variant-calling module or as user-

156  uploaded HapMap or VCF files. The module outputs some or all of the following 10

157  population genetics measures as specified by the user: allele frequency values, genotype

158  frequency values, population divergence (*F*st) values, nucleotide diversity values,

159  population structure results, a kinship matrix, a neighbor-joining tree, unweighted pair

160  group method with arithmetic mean (UPGMA) clustering results, linkage disequilibrium (LD)

161  values ($r^2$, D, and D'), and an LD graph. These analyses enable breeders to evaluate

162  germplasm diversity and select the best lines for future crosses to maintain long term

163  genetic gain. For example, the neighbor-joining tree (Paradis and Schliep, 2019) and

164  UPGMA clustering show the genetic similarities among individuals in the population and

165  enable breeders to assess the genetic diversity in the population.

166  *GWAS analysis module*

167  The GWAS module implements the 'GAPIT' R package (Wang and Zhang, 2021) to

168  identify SNPs underlying phenotypic variations. Phenotype and marker data can be

169  transferred directly from other modules in the platform or can be uploaded individually by

170  the user. Seven models from the 'GAPIT' package are included: Generalized Linear Model

171  (GLM), MLM, Compressed Mixed Linear Model (CMLM), Multi-Locus Mixed Model

172  (MLMM), Settlement of MLM Under Progressively Exclusive Relationship (SUPER), Fixed

173  and Random Model Circulating Probability Unification (FarmCPU), and Bayesian-

174  information and Linkage-disequilibrium Iteratively Nested Keyway (BLINK). The user can

175  define the minor allele frequency (MAF) for filtering SNPs and the number of principal

176  components (PCs) to include in the model as fixed effects. The output includes genome-

177  wide and single-chromosome Manhattan plots. Manhattan plots can also be generated for

178  multiple traits to determine whether nearby SNPs control differing traits. Q–Q plots show

179  the observed compared to the expected (i.e., uniformly-distributed) *p*-values. Additional

180  summary graphs show the distributions of traits and markers in the genome and the LD

181  values between nearby markers. Statistically significant markers are considered strong

182  candidates for marker-assisted selection (MAS) for desired traits or for fine-mapping to

183  identify causal genes for a specific phenotype.

184  *Genomic selection analysis module*

185      The genomic selection module is used to predict trait values among inbred lines,

186  hybrids, or progenies based on molecular markers. It consists of three steps: dataset

187  formation, training, and prediction. To generate a training dataset, the phenotype and

188  marker data are uploaded or retrieved from the phenotype analysis module and the variant-

189  calling module. In the dataset formation step, the program calculates the number of

190  samples with both phenotype and marker data. In the model training step, the user selects

191  a dataset, trait(s) of interest, and a model. The latter is either a statistical model such as

192  genomic BLUP (GBLUP) (VanRaden, 2008; Endelman, 2011), a classical machine learning

193  model, or a deep-learning model. After the model is trained, cross-validation is performed

194  and the predictive accuracy is displayed. The user can then choose a trained model (e.g.,

195  the model with the highest predictive accuracy) to predict the performance of offspring from

196  a cross or of the corresponding parental lines. This module yields predicted trait values for

197  each specified line, allowing a breeder to select lines that are predicted to have optimal

198  performance and to discard lines with undesirable traits. The breeder can thus select the

199  most promising potential crosses from many possible combinations, saving time and

200  resources. The predicted high performance inbreds or hybrids can be directly exported to

201  a germplasm table for crossing or field evaluation.

202  **Case study**

203      To showcase the Smart Breeding Platform's capabilities, we utilized a rice dataset

204  (Wang et al., 2018) with 100 varieties, each featuring multi-year, multi-location phenotypic

205  data (Supplemental Table S1). Germplasm data were uploaded to the "rice100" table in

206  the Germplasm Management module. The first five varieties advanced to

207  "Advancement2023." A pedigree table simulated 35 records from ERS470485 and

208  ERS470543. Three testing locations were added to the Location Management table, each

209  with 20 rows and 20 ranges (totaling 400 plots). Field testing experiments, rice_2022 and

210  rice_2023, followed RCBD designs across different locations. Phenotypic data underwent

211  analysis in the Phenotypic Statistical Analysis module, generating BLUP and BLUE values,

212  assessing genetic and phenotypic variance, heritability, and trait correlations.

213  In the Crossing Nursery module, 10 female and 10 male lines produced 28 two-way

214  crosses. Germplasm data populated the 'Nursery23' table, and pedigree records filled the

215  respective table. The Genomic Data Management module received the Os-Nipponbare-

216  Reference-IRGSP-1.0 file and paired-end sequencing data. The Genetic Variation Analysis

217  module conducted variant calling; 100 VCF files merged with criteria (depth > 50, quality

218  value > 500), resulting in 55,589 SNPs. The Genomic Statistical Analysis module used

219  SNPs for diverse population results. BLUP and BLUE values, along with merged marker

220  data, identified significant SNPs in the GWAS module for plant height and heading date.

221  The same datasets predicted values for new inbred and hybrid lines in the Genomic

222  Selection module. Details of the analyses and results can be found in the platform manual.

223  All data and results can be viewed by clicking "Enter as Guest" button on the login page.

224

225  In conclusion，this novel intelligent breeding platform integrates numerous data types

226  (seed inventory, field testing, phenotypes, SNP markers, and plant crosses) with key

227  analyses (GWAS, population genetic parameters, and genomic selection) in a single

228  seamless system. All analytical tools have user-friendly interfaces and are simple to

229  configure and run. The computing speeds for the genomic data analyses are substantially

230  faster in this platform than in conventional tools. Smart Breeding Platform provides a

231  comprehensive tool for the storage and management of germplasm data, experiments,

232  and statistical analyses, allowing breeders to more easily identify and generate optimal

233  germplasm, ultimately increasing the speed of genetic gain.
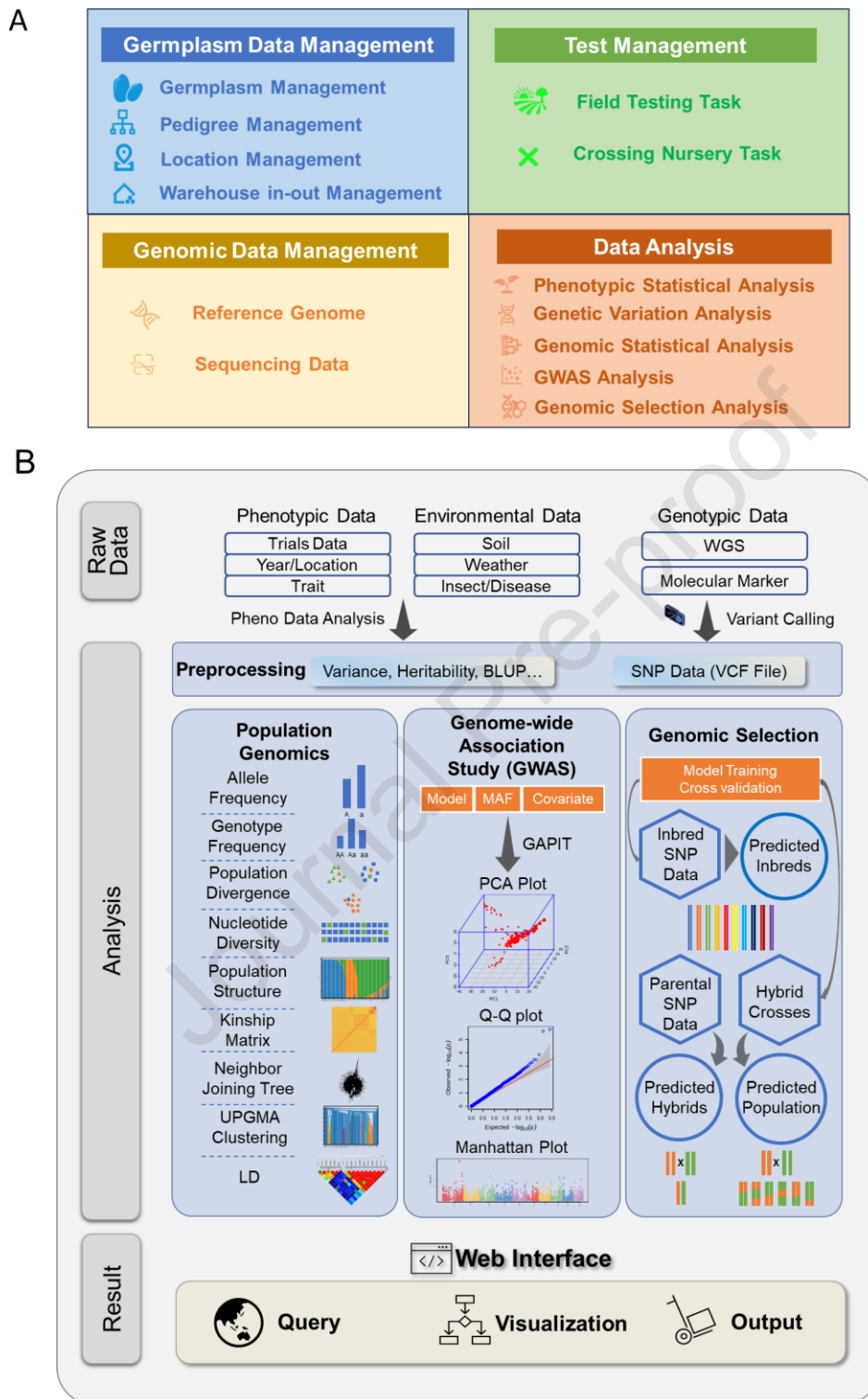
234

235

236

237

238

239

240 **Figure**



241
242 Figure 1. Main sections (A) and analysis workflow (B) of the Smart Breeding Platform.
243 WGS, whole genome sequencing; BLUP, best linear unbiased prediction; SNP, single
244 nucleotide polymorphism; UPGMA, unweighted pair group method with arithmetic mean;
245 LD, linkage disequilibrium; and PCA, principal component analysis.
246

**Funding**

This work was supported by the Alibaba Foundation, the National Natural Science Foundation of China (32188102, 32361143514), Innovation Program of Chinese Academy of Agricultural Sciences, and the Project of Hainan Yazhou Bay Seed Lab (B21HJ0223).

**Author contributions**

W.Z., Q.Q., H.L., and F.G. conceived the project. X.L., P.Z., Y. F., J.M., S.G., L.S., M.A., H.L., and F.G. conducted data analyses and platform development. Z.Y. organized data used in case study. W.Z., Q.Q., L.L., W.W., and W.F. provided data and technical guidance. H.L., X.L., P.Z., M.A., and F.G. wrote the manuscript. All authors read and approved the manuscript.

**Acknowledgements**

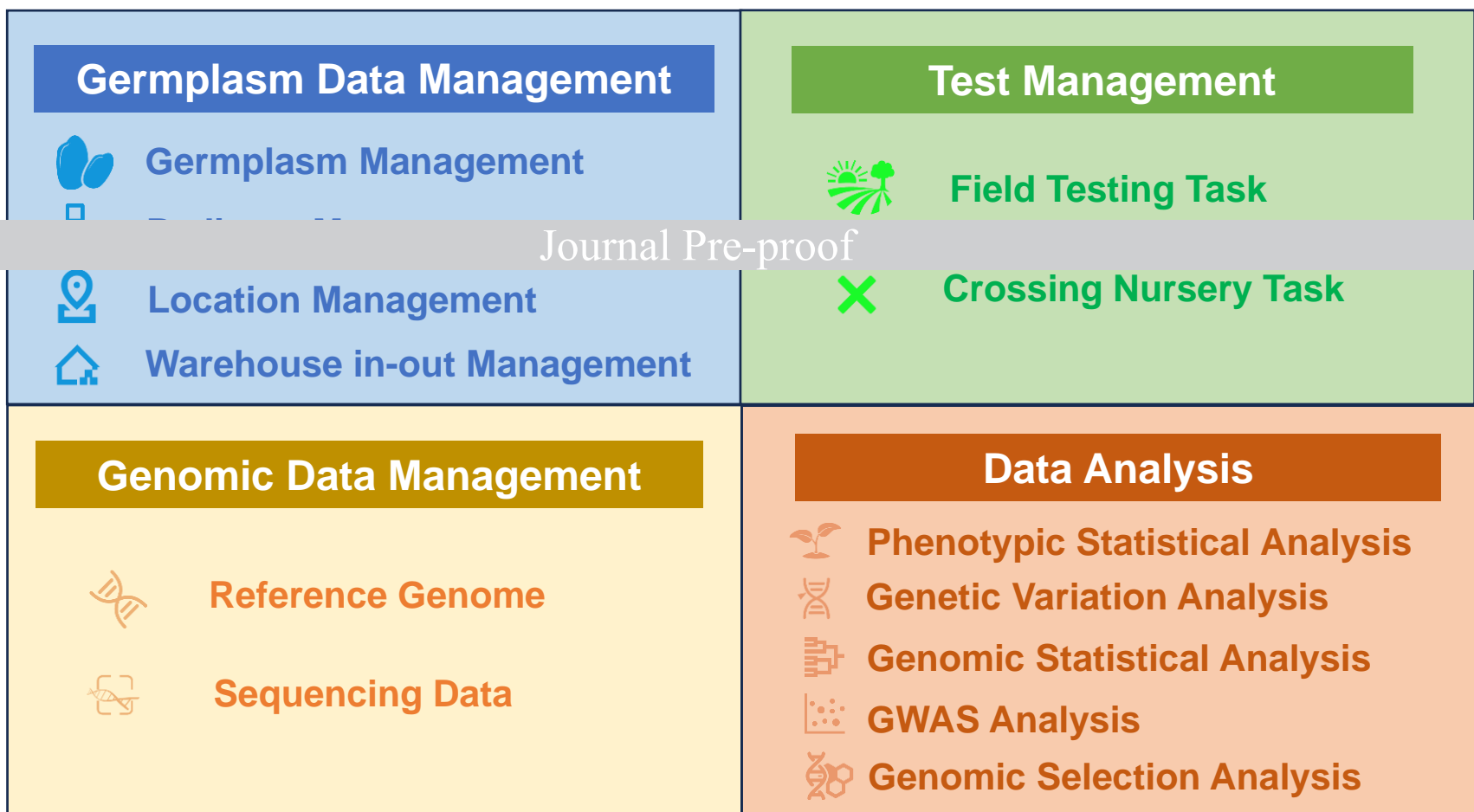The authors declare that they have no conflict of interest.

**References**

**Bates, D., Mächler, M., Bolker, B., and Walker, S.** (2015). Fitting Linear Mixed-Effects Models Using lme4. J. Stat. Softw. **67**:1-48. https://doi.org/10.18637/jss.v067.i01.

**Brandies, P.A., and Hogg, C.J.** (2021). Ten simple rules for getting started with command-line bioinformatics.   Public Library of Science San Francisco, CA USA. https://doi.org/10.1371/journal.pcbi.1008645

**Covarrubias-Pazaran, G.** (2016). Genome-Assisted Prediction of Quantitative Traits Using the R Package sommer. PLoS One **11**:e0156744. 10.1371/journal.pone.0156744.

**Endelman, J.B.** (2011). Ridge Regression and Other Kernels for Genomic Selection with R Package rrBLUP. The Plant Genome **4**:250-255 https://doi.org/10.3835/plantgenome2011.08.0024.

**Muñoz, F., and Sanchez, L.** (2020). breedR: Statistical Methods for Forest Genetic Resources Analysts. R package version 0.12-5.

**Paradis, E., and Schliep, K.** (2019). ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R. Bioinformatics **35**:526-528. https://doi.org/10.1093/bioinformatics/bty633

276 **Robinson, J.T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E.S., Getz, G.,**
277      **and Mesirov, J.P.** (2011). Integrative genomics viewer. Nat. Biotechnol. **29**:24-26.
278      https://doi.org/10.1038/nbt.1754

279 **Rodriguez-Alvarez, M.X., Boer, M.P., van Eeuwijk, F.A., and Eilers, P.H.** (2018). Correcting
280      for spatial heterogeneity in plant breeding experiments with P-splines. Spatial. Stat.
281      **23**:52-71. https://doi.org/10.1016/j.spasta.2017.10.003

282 **Sharma, S., Kumar, K., Tribhuvan, K., Rita, Kumar, S., Jain, P., Saxena, S., Vijayan, J.,**
283      **Srivastava, H., and Gaikwad, K.** (2022). High-throughput Genotyping Platforms.
284      Genotyping by Sequencing for Crop Improvement. New York: John Wiley & Sons 22-37.
285      https://doi.org/10.1002/9781119745686.ch2

286 **VanRaden, P.M.** (2008). Efficient methods to compute genomic predictions. J. Dairy Sci.
287      **91**:4414-4423. https://doi.org/10.3168/jds.2007-0980

288 **Wang, J., and Zhang, Z.** (2021). GAPIT version 3: boosting power and accuracy for genomic
289      association and prediction. Genomics Proteomics Bioinformatics **19**:629-640.
290      https://doi.org/10.1016/j.gpb.2021.08.005

291 **Wang, W., Mauleon, R.** (2018). Genomic variation in 3,010 diverse accessions of Asian
292      cultivated rice. Nature **557**:43-49. https://doi.org/10.1038/s41586-018-0063-9

293 **Xiao, Q., Bai, X., Zhang, C., and He, Y.** (2022). Advanced high-throughput plant phenotyping
294      techniques for genome-wide association studies: A review. J. Adv. Res. **35**:215-230.
295      https://doi.org/10.1016/j.jare.2021.05.002

296 **Xu, Y., Zhang, X., Li, H., Zheng, H., Zhang, J., Olsen, M.S., Varshney, R.K., Prasanna,**
297      **B.M., and Qian, Q.** (2022). Smart breeding driven by big data, artificial intelligence and
298      integrated genomic-enviromic prediction. Mol. Plant **15**:1664-1695.
299      https://doi.org/10.1016/j.molp.2022.09.001

300 **Yin, L., Zhang, H., Tang, Z., Xu, J., Yin, D., Zhang, Z., Yuan, X., Zhu, M., Zhao, S., and Li,**
301      **X.** (2021). rMVP: a memory-efficient, visualization-enhanced, and parallel-accelerated
302      tool for genome-wide association study. Genomics Proteomics Bioinformatics **19**:619-
303      628. https://doi.org/10.1016/j.gpb.2020.10.007

304
305

**A**

**Germplasm Data Management**
- Germplasm Management
- Location Management
- Warehouse in-out Management

**Test Management**
- Field Testing Task
- Crossing Nursery Task

**Genomic Data Management**
- Reference Genome
- Sequencing Data

**Data Analysis**
- Phenotypic Statistical Analysis
- Genetic Variation Analysis
- Genomic Statistical Analysis
- GWAS Analysis
- Genomic Selection Analysis

**B**

**Raw Data**

Phenotypic Data
- Trials Data
- Year/Location
- Trait

Environmental Data
- Soil
- Weather
- Insect/Disease

Genotypic Data
- WGS
- Molecular Marker

Pheno Data Analysis

Variant Calling

**Analysis**

**Preprocessing** Variance, Heritability, BLUP... | SNP Data (VCF File)

**Population Genomics**
- Allele Frequency
- Genotype Frequency
- Population Divergence
- Nucleotide Diversity
- Population Structure
- Kinship Matrix
- Neighbor Joining Tree
- UPGMA Clustering
- LD

**Genome-wide Association Study (GWAS)**

Model | MAF | Covariate

GAPIT

PCA Plot

Q-Q plot

Manhattan Plot

**Genomic Selection**

Model Training Cross validation

Inbred SNP Data → Predicted Inbreds

Parental SNP Data | Hybrid Crosses

Predicted Hybrids | Predicted Population

**Result**

**Web Interface**

Query | Visualization | Output