

1 **Rice Galaxy: an open resource for plant science**

2 Venice Juanillas¹, Alexis Dereeper², Nicolas Beaume¹, Gaetan Droc³, Joshua Dizon¹, John Robert

3 Mendoza⁸, Jon Peter Perdon⁸, Locedie Mansueto¹, Lindsay Triplett⁷, Jillian Lang⁷, Gabriel Zhou⁴,

4 Kunalan Ratharanjan⁴, Beth Plale⁴, Jason Haga⁵, Jan E. Leach⁷, Manuel Ruiz³, Michael Thomson^{1,6},

5 Nickolai Alexandrov^{1,10}, Pierre Larmande², Tobias Kretzschmar^{1,9}, Ramil P. Mauleon¹

6 **Author affiliations**

7 ¹ International Rice Research Institute, Manila, Philippines

8 ² Institut de recherche pour le développement (IRD), University of Montpellier, DIADE, IPME,

9 Montpellier, France

10 ³ CIRAD, UMR AGAP, F-34398 Montpellier, France

11 ⁴ Indiana University, 107 S Indiana Ave, Bloomington, IN 47405, USA

12 ⁵ National Institute of Advanced Industrial Science and Technology, AIST Tsukuba Central 1,1-1-1

13 Umezono, Tsukuba, Ibaraki 305-8560 JAPAN

14 ⁶ Department of Soil and Crop Sciences, Texas A&M University, Houston, USA

15 ⁷ Department of Bioagricultural Sciences and Pest Management, Colorado State University, Fort Collins,

16 CO 80523-1177

17 ⁸ Advanced Science and Technology Institute, Department of Science and Technology, Quezon City,

18 Philippines

19 ⁹ Southern Cross Plant Science, Southern Cross University, Lismore, Australia

20 ¹⁰Inari Agriculture Inc., 200 Sidney St, Cambridge, Massachusetts, USA

21 Correspondence should be addressed to R.P.M. (r.mauleon@irri.org, ORCID: 0000-0001-8512-144X)

22

23

24 **Abstract**

25 *Background*

26 Rice molecular genetics, breeding, genetic diversity, and allied research (such as rice-pathogen
27 interaction) have adopted sequencing technologies and high density genotyping platforms for genome
28 variation analysis and gene discovery. Germplasm collections representing rice diversity, improved
29 varieties and elite breeding materials are accessible through rice gene banks for use in research and
30 breeding, with many having genome sequences and high density genotype data available. Combining
31 phenotypic and genotypic information on these accessions enables genome-wide association analysis,
32 which is driving quantitative trait loci (QTL) discovery and molecular marker development. Comparative
33 sequence analyses across QTL regions facilitate the discovery of novel alleles. Analyses involving DNA
34 sequences and large genotyping matrices for thousands of samples, however, pose a challenge to non-
35 computer savvy rice researchers.

36 *Findings*

37 We adopted the Galaxy framework to build the federated Rice Galaxy resource, with shared datasets,
38 tools, and analysis workflows relevant to rice research. The shared datasets include high density
39 genotypes from the 3,000 Rice Genomes project and sequences with corresponding annotations from
40 nine published rice genomes. Rice Galaxy includes tools for designing single nucleotide polymorphism
41 (SNP) assays, analyzing genome-wide association studies, population diversity, rice-bacterial pathogen
42 diagnostics, and a suite of published genomic prediction methods. A prototype Rice Galaxy compliant to
43 Open Access, Open Data, and Findable, Accessible, Interoperable, and Reproducible principles is also
44 presented.

45 *Conclusions*

46 Rice Galaxy is a freely available resource that empowers the plant research community to perform state-
47 of-the-art analyses and utilize publicly available big datasets for both fundamental and applied science.

48 **Keywords**

49 Rice, Breeding, workflow, genomes, high-density genotypes, reproducibility, SNP, GWAS, Galaxy project

50 **Findings**

51 *Background*

52 With the decreasing cost of genome sequencing, rice molecular geneticists, breeders and diversity
53 researchers are increasingly adopting genotyping technologies as routine components in their
54 workflows, generating large datasets of genotyping and genome sequence information. Concurrently
55 international consortia have made re-sequencing or high density genotyping data from representative
56 diversity collections publically available. These include, but are not limited to the medium-depth (15-20x
57 coverage) resequencing data of the 3,010 accessions from the 3K Rice Genome (3K RG) Project (~1 – 2
58 million SNPs per accession) [1,2] and the 700,000 SNP Affymetrix array data for the 1,445 accessions of
59 the High Density Rice Array (HDRA) germplasm collections [3]. The corresponding accessions are
60 available at non-profit prices from the Genetic Resource Center (GRC) of the International Rice Research
61 Institute (IRRI) for phenotyping, allowing subsequent Genome-Wide Association Studies (GWAS).
62 Analysis of such datasets is a challenge to rice researchers due to (1) the fairly large data matrix and the
63 compute-intensive algorithms that requires specialized computing infrastructure (a fairly large RAM,
64 powerful CPU, and large disk space), and (2) the relative difficulty in using Open Source / free software
65 tools for analysis, which are commonly provided without graphical user interface and require proper
66 installation in a Linux operating system environment.

67 On the computational side, public web resources with specialized tools already exist, and are
68 maintained at different institutions. The Rice SNP-Seek database [4,5], largely developed and hosted by
69 IRRI, contains phenotypic, genotypic, and passport information for over 4,400 rice accessions from large
70 scale rice diversity projects such as the 3K RG and the HDRA collections. SNP-Seek ([http://snp-
seek.irri.org](http://snp-
71 seek.irri.org)) currently contains phenotype data for 70 different morphological and agronomic traits

72 and stores SNPs and small indels discovered by mapping the 3K RG accessions to four published rice
73 draft genome assemblies, collectively resulting in the discovery of ~11M new SNPs and ~0.5M new
74 indels. While SNP-Seek focused on delivery of prior analyzed content rather than providing an analysis
75 platform, Gigwa [6] (<http://gigwa.southgreen.fr/gigwa/>), hosted at the South Green portal [7]
76 (<http://www.southgreen.fr/>), is a scalable and user-friendly web-based tool which provides an easy and
77 intuitive way to explore large amounts of genotyping data from next-generation sequencing (NGS)
78 experiments. Gigwa allows for filtering of genomic and genotyping data from NGS analyses based not
79 only on variant features, including functional annotations, but also on genotype patterns to explore the
80 structure of genomes in an evolutionary context for a better understanding of the ecological adaptation
81 of organisms. Gramene [8] is a curated, open-source, integrated data resource for comparative
82 functional genomics in crops and model plant species that, among other species, includes rice. Data and
83 analysis tools are available as portals at the Gramene site (<http://gramene.org/>). In these resources
84 mentioned, the analyses methodologies are custom-built by the respective projects.

85 There are other freely available web-based bioinformatics and breeding informatics software tools,
86 optimized for plant species other than rice, including Araport (<https://www.araport.org/>) for
87 Arabidopsis, Cassavabase (<https://cassavabase.org/>) for cassava, and The Triticeae Toolbox (T3,
88 <https://triticeaetoolbox.org/>) for wheat and barley. While these tools are very useful, they are
89 species/crop specific and custom-built for the specialized requirements of their respective communities
90 (such as project datasets), making adoption in rice challenging for at least two reasons: (1) the need to
91 produce curated rice datasets that work seamlessly with the software system (e.g. genome-browser
92 ready data, curated genes, published QTLs from biparental crosses and GWAS and markers associated to
93 traits), and (2) the need for a dedicated software development team to customize the application for
94 rice-specific data and analyses.

95 The ability of software to automate repetitive analyses task is attractive for data analysts, and the public
96 sharing of the analytical methodology (as opposed to just the raw data and the results) enhances
97 reproducibility and is being supported by academic communities of practice such as FORCE11
98 (<https://www.force11.org>). Many research groups working with NGS data have a high demand for
99 computing infrastructure and their complex analyses often comprise several steps using different
100 software tools (pipeline). The deployment of these different software tools is a big challenge to small
101 institutions without dedicated scientific computing support staff. There is no single solution to address
102 these challenges. Our approach to help overcome them is the integration of a range of these different
103 bioinformatics tools into the Galaxy bioinformatics system. Galaxy [9] is a web-based analysis
104 workbench and workflow management system initiated at the Penn State University. It includes a
105 collection of software packages which can be operated via a web browser on a public server. Galaxy is a
106 mature community effort, supported by various high-powered institutions, is relatively easy to deploy
107 and maintain, and thus well-suited to serve low and moderately resourced institutions such as IRRI. The
108 graphical user interface of Galaxy means that no knowledge of code is needed, thus facilitating
109 bioinformatics analyses by researches without computational expertise.

110 We built a suite of federated Galaxy resources and tools, which we collectively named **Rice Galaxy**
111 (Figure 1). Rice Galaxy contains shared software tools and datasets tailored to the needs of rice
112 researchers and breeders, providing computing resources through an easy-to-use interface, and
113 allowing reproducibility and publication of analytical methodology and results.

114 The Rice Galaxy federated resources are hosted at:

- 115 • IRRI Galaxy @ International Rice Research Institute: <http://galaxy.irri.org>
- 116 • Rice Galaxy (common) Toolshed: <http://52.76.88.51:8081/>

117 *DISCUSSION*

118 **1. Built-in / interoperable rice data**

119 The Rice Galaxy system is customized to provide rice-specific genomic and genotypic data. Of primary
120 importance is the gold-standard *japonica* variety reference genome (Nipponbare IRGSP release 1.0) [10],
121 to which the reference gene models and most of the SNPs published have been anchored. In addition,
122 eight medium to high quality published genomes from various sequencing projects and the respective
123 genome annotations for each are installed as alternative genome builds and are available as drop-down
124 menu choices in Rice Galaxy. These include four high-quality builds from *indica*-type varieties Minghui
125 63 and Zhengshan 97 [11], IR 8 (GenBank: MPPV00000000.1), Shuhui 498 [12], as well as an aus-type
126 variety N 22 (GenBank: LWDA00000000.1), as well as four medium to low quality genomes, two *indica*
127 (IR 64 , [13] and 93-11, [14]) and two aus-type rice genomes (DJ 123, [13] and Kasalath, [15]). While
128 these references were selected to represent diversity, they further represent variety groups that display
129 agronomically important characteristics, such as heat and drought tolerance, disease resistance,
130 submergence tolerance, adaptation to low-phosphorus soil, wide adaptability, good grain quality,
131 aerobic (upland) adaptation and deep roots [16-18]. Even though these genomes are highly similar to
132 each other, they each contain unique regions (from 12.3 Mbp to 79.6 Mbp) that may harbor genes
133 restricted to these variety-groups [5]. With the availability of several reference genomes, it becomes
134 relatively straightforward to custom design SNP assays that are either of broad utility across varietal
135 groups or specific to single groups.

136 Rice Galaxy includes genotyping data of the 3k RG (such as the 3K RG 3024 accessions x 4.8M filtered
137 SNPs, 440K core SNPs, 1M GWAS-ready SNPs, and 2.3M indels) useful for GWAS, region-specific
138 diversity analyses, and single locus allele mining in the shared data library.

139 **2. Toolkits Built (and detailed discussion of each toolkit)**

140 *SNP assay design: Lift-over of SNPs from one genome to another*

141 SNPs discovered relative to the gold-standard reference genome (Nipponbare IRGSP 1.0, [10]) are
142 commonly used in QTL mapping (either by GWAS or biparental cross). In order to develop robust

143 markers associated with the trait of interest, however, a SNP assay that works in the target varietal
144 groups is needed. Consequently there is a need to “lift-over” SNPs from one genome to another (for
145 example from Nipponbare *japonica* to an *indica* varietal group represented by IR 64). The workflow is as
146 follows: (1) Get flanking sequences surrounding the target SNP in source genome (the main reference
147 Nipponbare), (2) align these flanking sequences to target genome of variety of interest to verify if it hits
148 a unique region in the target genome of similar location from the source genome, allowing some
149 mismatches but not allowing multiple region hits, and (3) identify the flanking sequences surrounding
150 the lifted-over SNP in the context of the target genome, for SNP assay design. The shared workflow is
151 published in Rice Galaxy as [SNP lift-over].

152 *3k RG data access*

153 Rice Galaxy provides access to the raw variant call format (VCF) files of each accession in the 3K RG
154 project via connection (as data source in Rice Galaxy) to the 3,000 rice genomes at Amazon Web
155 Services (AWS) Public Data (<https://aws.amazon.com/public-datasets/3000-rice-genome/>), with tools
156 allowing region-specific download. In Rice Galaxy, tools in the [Get Data / FROM 3KRG] section allows
157 listing of the accessions in the 3K RG and retrieval of genotype data for a selected accession of interest
158 from the 3K RG collection. The subset genome region of interest (chromosome name – base start – base
159 end) can be specified and extracted from the VCF of the accession of interest stored in AWS Public
160 Datasets.

161 In addition, we developed an original Rice Galaxy component called Rapid Allelic Variant extractor
162 (RAVE), which allows simultaneous extraction of genotyping data from several accessions of the internal
163 3K RG resource. It relies on the PLINK software [20], which efficiently builds a user-adjusted genotyping
164 submatrix from a compressed PLINK binary biallelic genotype table (bed file + bim, fam files). Users can
165 customize the genotyping dataset vertically by choosing a subpopulation (*indica*, *japonica*, *aromatic*,
166 *aus*, *tropical*, *temperate*, etc.) or setting a list of varieties, and horizontally by restricting variations with a

167 list of genomic regions, or a list of gene names. Additionally, users can filter the SNP positions by
168 specifying thresholds for missing data or minor allele frequency (MAF). The extracted VCFs can be
169 directly generated by Rice Galaxy, stored as output into the history pane of the Galaxy interface, and can
170 be reformatted to Hapmap, a versatile file format for further analyses such as marker (SNP) design,
171 GWAS analyses, or visualization in within a JBrowse [21] genome browser (Vcf2jbrowse component).
172 External SNP datasets can also be imported into Rice Galaxy and merged with 3k accessions in order to
173 compare and look at the closest genotypes using SNIPlay [22] workflow.

174 *GWAS analysis using TASSEL*

175 Using this feature, it is relatively easy to construct a genotyping matrix for a subset of accessions from
176 the 3K RG and connect associated phenotypic information to perform GWAS analyses online, with
177 outputs being decorated with various graphical enhancements. For the 3K RG accessions, the subset 1M
178 GWAS and 440K Core SNPs that is usable for GWAS is already available as shared dataset in Rice Galaxy
179 (Figure 2). Researchers working on the 3K RG panel can generate new phenotyping data from their
180 respective experiments, upload the phenotype data into Rice Galaxy, and then perform GWAS using the
181 TASSEL bioinformatics tool [23]. The GWAS Rice Galaxy workflow implementing TASSEL and Multi-Locus
182 Mixed-Model package for association studies is shared from SNIPlay at Rice Galaxy (Figure 3).

183 Aside from GWAS with 3K RG datasets, researcher-generated marker and phenotype data (outside of 3K
184 RG) can also be uploaded to Rice Galaxy for GWAS analysis.

185 *Genomic selection using Oghma genome prediction tool*

186 Genomic selection (GS) is a promising breeding technique with potential to improve the efficiency and
187 speed of the breeding process in rice [24]. With the intent of enabling the GS analysis process used on
188 the 2 datasets in the Spindel et al. [24] study, (encoding data, filtering data to keep informative markers,
189 creating a model from training set, evaluating the model and finally, performing the prediction itself),
190 and to automate the analysis pipeline, the relevant packages (`methods`, `fpc`, `cluster`,

191 `vegan`, `pheatmap`, `pROC`, `randomForest`, `miscTools`, `pRF`, `e1076`,
192 `rrBLUP`, and `glmnet`) for the R Statistical language (<https://www.r-project.org/>) were installed in
193 Rice Galaxy and the tool suite was collectively named Oghma (Operators for Genome decipHering by
194 Machine learning). Quality control tool (based on PLINK) and imputation tool using Beagle [25]
195 (<https://faculty.washington.edu/browning/beagle/beagle.html>) were also installed. Four phenotype
196 prediction/classifier methods (rrBLUP, random forest, SVM and lasso) were identified as relevant and
197 deployed as tools in Rice Galaxy (Figure 4).

198 Figure 5 shows the overall GS analysis workflow using Oghma. Genotypes are encoded through [encode]
199 tool. For the training set, an encoded genotype and the corresponding phenotype files are used by a
200 classifier tool to train a model, which can be used with another encoded genotype file to predict trait
201 values (the genomic prediction). It is important to note that (1) both genotype for training and
202 genotype to predict must have the same markers (and thus, genotype files must have the same number
203 of columns) to make a prediction, and (2) the "evaluation" option of the classifier tool can have any
204 value except 1 (it is recommended to keep the default value = 0).

205 A big challenge when using machine learning approaches for genomic prediction is the optimization of
206 the model based on training data, specifically setting the best parameters of the methods mentioned
207 prior. Oghma was designed to automate the optimization of the parameter(s) of the classifiers on the fly
208 (as opposed to manual tweaking), thus allowing users without experience of machine learning to easily
209 optimize a model for their own data. Oghma includes some tools to evaluate prediction accuracy to
210 allow the user to choose the most accurate method on their data by performing a cross-validation with
211 a user-uploaded training set. Two metrics, the coefficient of determination (R^2) and the correlation
212 between predicted and observed phenotype, and a visualization (scatterplot of predicted vs observed)
213 have been implemented to evaluate the methods. The [computeR2] and [plotPrediction] tools are used
214 to compute R^2 and visualize the accuracy of prediction. These tools both take the true phenotypes and

215 the predicted outputs as inputs (take note that both predictions and phenotypes data must be in the
216 same order), and return the computed R^2 or the scatterplot display of true phenotype vs prediction.
217 Oghma can be used to evaluate a classifier (Figure 6). Like the general GS workflow, genotype and
218 phenotype are used as input for any classifier, but the "evaluation" option must be set to 1. Fold for
219 cross-validations are designed through the [fold] tool, which take as input the encoded file. These folds
220 are used as extra argument by the classifier tools. The chosen classifier tool produces a file, which is not
221 a model but the prediction of the test set for each cross-validation. This output is used as input, along
222 with the phenotypes and folds, by [evaluation], which output some performances index (R^2 and
223 correlation). Although it does give a real indication of performances, trying to predict the training set
224 (i.e. using the same genotype file in the pipeline described above), or at the least, showing if the
225 classifier is not under-fitting the data.

226 We installed several classifiers in Oghma to allow users to test the best one(s) suited for their dataset, as
227 our literature survey shows that no method seems to outperform the others on all genomic prediction
228 tasks. It was noticed that Random Forest was the most accurate and the most stable classifier on Spindel
229 dataset [24], thus we set this as default in Oghma. An original aggregation method is also implemented
230 in Oghma, aggregating outputs of multiple classifiers to improve prediction. This tool takes as input the
231 prediction of n classifiers and tries to aggregate them through weighted mean of the prediction (weight
232 optimized by genetic algorithm) or regression (multiple type of regression have been implemented,
233 based on decision tree, SVM and Random Forest). Limited testing shows that this approach is promising,
234 matching Random Forest in some cases, especially with a meta-SVM, with polynomial or linear model, as
235 aggregation method, but still needs some improvement as the accuracy remains unstable when
236 evaluated through cross-validation (data not shown). The aggregation method can also be evaluated
237 using the aforementioned evaluation tools.

238 *Diversity and population structure analysis of end-user datasets*

239 SNP datasets - such as those extracted from the 3K RG resource after a filtering by the RAVE module or
240 custom sets directly uploaded in Rice Galaxy environment (Figure 7) can be processed for a complete
241 exploration and large scale analysis thanks to the SNIPlay Rice Galaxy workflow (Figure 8). The workflow
242 is available through the instance, requiring a VCF file as input. This workflow allows various analyses: (i)
243 SNP annotation by snpEff (<http://snpeff.sourceforge.net/>) wrapper preconfigured for RGAP release 7.0
244 [26] (<http://rice.plantbiology.msu.edu/>) gene models (ii) variant filtration using PLINK wrapper, (iii)
245 general statistics such as Transition-Transversion ratio, levels of heterozygosity and missing data for
246 each variety using VCFtools, (iv) SNP density analysis, (v) diversity indices calculation in sliding windows
247 along the genome using VCFtools (Pi, Tajima's D, FST if subpopulations provided), (vi) linkage
248 disequilibrium, (vii) population structure by sNMF ([http://membres-](http://membres-timc.imag.fr/Olivier.Francois/snmf/index.htm)
249 [timc.imag.fr/Olivier.Francois/snmf/index.htm](http://membres-timc.imag.fr/Olivier.Francois/snmf/index.htm)), (viii) Principal Component Analysis and Identity By State
250 (IBS) clustering of varieties by PLINK, and (ix) SNP-based distance phylogenetic tree by FastME
251 (<http://www.atgc-montpellier.fr/fastme/>). Most key steps are decorated with sophisticated
252 visualizations using a dedicated plugin. Visualization can be displayed by clicking on the [visualization]
253 icon.

254 In practice, this workflow can be processed for many applications such as the identification of possible
255 introgression events, the identification of putative genomic regions involved in the control of qualitative
256 trait through a FST approach, the investigation for potential duplicates in the 3K RG accessions dataset
257 and custom datasets, or the estimate of closest varieties of new sequenced accessions, by ranking a list
258 of varieties from the database most closely matching the given sample. It can be used also for the close
259 inspection of genomic region of interest after a GWAS analysis, through a linkage disequilibrium focus or
260 the haplotyping of candidate genes.

261 *Uniqprimer module*

262 Uniqprimer is a workflow for comparative genomics-based diagnostic primer design, developed from a
263 pipeline used in-house at Colorado State University to develop novel species and subspecies-level
264 diagnostic tools for bacterial plant pathogens including pathovars of *Xanthomonas translucens* [27],
265 geographical variants of rice-associated *Xanthomonas* spp. [28-30], and the genetically diverse rice
266 pathogen *Pseudomonas fuscovaginae* [31]. Uniqprimer is now deployed in Rice Galaxy for user-friendly
267 diagnostic primer design from draft or complete pathogen genomes. The user inputs multiple bacterial
268 genomes from diagnostic target species as well as non-target species (i.e. “include” and “exclude”
269 genome files), and the tool performs comparative alignment, primer design, and primer validation to
270 output a list of primers that are specific to the target genomes (Figure 9). The uniqprimer standalone
271 program is written in Python and is available at the Southgreen github repository
272 (<https://github.com/SouthGreenPlatform/Uniqprimer>), along with the detailed documentation for
273 developers and end-users.

274 **3. Rice Galaxy OA: a Prototype for Open Access**

275 IRRI, as a member center of the Consultative Group for International Agricultural Research (CGIAR,
276 <https://www.cgiar.org/>), complies with the CGIAR policy on Open Access and Open Data
277 (<https://www.cgiar.org/how-we-work/accountability/open-access/>). In collaboration with Indiana
278 University in the United States and National Institute of Advanced Industrial Science and Technology in
279 Japan, and carried out through grants from the National Science Foundation (NSF) in the US and the
280 MacArthur Foundation through the Research Data Alliance (RDA - <https://www.rd-alliance.org/>), the
281 team undertook a prototyping effort to bring the Rice Galaxy system to maximum compliance with the
282 CGIAR policy.

283 The basis for the design to add open access to Rice Galaxy is a foundational technical idea emerging
284 from activities occurring in the international RDA. This idea acknowledges that for open data access to
285 be broadly realized, all meaningful data objects must have a globally unique and persistent identifier

286 (PID). Globally unique means the name is not shared with other objects on a global scale. An identifier is
287 persistent when the PID itself cannot be destroyed, and when the relationship between the identifier
288 and the data object it points to is permanent. Through an international working group in RDA, a team of
289 researchers is advancing the notion of PID Kernel Information, which injects a tiny amount of carefully
290 selected metadata into a PID record. This technique has the potential to stimulate an entirely new
291 ecosystem of third party services that can process the billions of expected PIDs. The key challenge of this
292 working group is to determine which from amongst thousands of relevant metadata are suitable to
293 embed in the PID record.

294 Our design draws on earlier work by us in data provenance capture and representation [32-34] and
295 employs a hands-off technique (*data provenance capture*) to gather information about a researcher's
296 rice genomics analysis as the analysis is running. Through this technique, information acquired while the
297 analysis is running is compiled and combined with pre-analysis information that is available at the
298 beginning of the analysis workflow. Such information includes who performed the analysis, when it was
299 performed, and under what conditions.

300 There have been earlier approaches to capture provenance of Galaxy workflows. Geocks *et al* [35]
301 developed a history panel for users to facilitate reproducibility. Gaignard *et al*. [36] proposes the SHARP
302 toolset, a semantic web (i.e. linked data) approach of harmonizing provenance collected from both the
303 Galaxy and Taverna workflow systems. Kanwal *et al*. [37] captured the activity of a workflow (called a
304 *provenance trace*) including the version of analysis tools run, the software parameters used, and the
305 data objects produced at each workflow step. This work also targets increased reproducibility of past
306 workflow instances. Missier *et al*. [38] proposes the "Golden Trail" architecture to describe and store
307 workflow runtime provenance retrieved from Galaxy. The golden trail of provenance that is collected
308 can be used to construct a virtual experiment view of past workflow runs. The four research
309 contributions described further underline the need for the capture of provenance from workflow

310 systems. They propose different but equally important uses of data provenance, that is, to facilitate the
311 improvement of science through reproducibility and construction of virtual views of an experiment once
312 it has completed.

313 Our design for Rice Galaxy Open Access (OA) shares similarities with these other techniques, however its
314 end goal is different, which is to advance open access, hence making Rice Galaxy consistent with CGIAR's
315 open access policy. To do this, we focus on each piece of data and information deemed valuable that
316 emerges from workflow runs deemed to be of importance. This particular data and information must be
317 retained and shared with others, while being subject to reasonable restrictions. This is a highly selective
318 approach to provenance capture, and one that makes our work unique. We briefly outline the solution
319 here and identify resources for those interested in pursuing the topic in more detail.

320 The architecture of Rice Galaxy OA (Figure 10A) utilizes the Handle system [39] and two standards
321 emerging from the Research Data Alliance, RDA PID Type [40] and the Data Type Registry [41]. It
322 additionally uses storage and compute resources provisioned through the NSF funded project, Pacific
323 Rim Applications and Grid Middleware Assembly (PRAGMA).

324 A researcher interacts with the open access enhanced Rice Galaxy system as follows:

- 325 (1) Researcher performs an analysis in Rice Galaxy
- 326 (2) Data objects (input data, output data, information such as configuration parameters) are
327 extracted from Rice Galaxy OA into a PRAGMA Data Repository Database (MongoDB) (Figure
328 10A),
- 329 (3) The data objects are assigned Persistent Identifiers, the PID Kernel Information is assigned into
330 the PID record at this time, and a landing page created for each (Figure 10B).
- 331 (4) Data objects can be downloaded from the Data Identity server and re-loaded to the Rice Galaxy
332 server for full faithful reproduction of the analysis

333 The resulting system appears to be promising and addresses a number of the recommendations from
334 CGIAR. The Rice Galaxy OA system is a user transparent means of harvesting digital objects from
335 applications and assigning PIDs to scientific outcomes. The architecture is modular and built with default
336 PID information types and metadata using RDA products (Figure 10A). Although this proof-of-concept
337 prototype successfully demonstrates the feasibility of this approach, there remains some future work.
338 The community needs to provide feedback on which data and information products are most important
339 to retain and make available. Additionally, not all workflow runs are important to a researcher as they
340 could be system tests or new workflow tests. Thus, how a researcher identifies the items he/she wishes
341 to make available to others and when, remains an important consideration for this system. For more
342 information, points of contact to the team, the underlying software for Rice Galaxy OA, and the link to
343 the prototype server can be found at <https://github.com/Data-to-Insight-Center/RDA-PRAGMA-Data-Service/wiki/Welcome-to-PRAGMA-Data-Service-Prototype> .

345 ***4. Rice Galaxy architecture discussion***

346 We deployed IRRRI Galaxy in an AWS EC2 instance (t2.large instance 2 vCPU, 4 GiB RAM) for the
347 production server deployment in the cloud with Linux Ubuntu release 12.04.2 LTS (GNU/Linux 3.2.0-40-
348 virtual x86_64) using Galaxy release 14 as described in the Galaxy documentation.
349 External data from the 3K RG Project files stored in the 3K RG AWS Simple Storage Service (S3) Public
350 Data resource hosted at <http://s3.amazonaws.com/3kricegenome/> (or s3:// 3kricegenome/) is accessed
351 using AWS S3 Command Line Interface, a command line tool utility in AWS that provides an interface to
352 access AWS S3 objects (CLI, <https://docs.aws.amazon.com/cli/latest/reference/s3/>). First, Rice Galaxy
353 connects to the 3K RG AWS bucket using s3API and allows the objects inside the bucket to be
354 transparent to Galaxy. VCF files (and the accompanying index files) are downloaded to Rice Galaxy using
355 the S3 CLI with the `aws s3 cp` command, executed as:

```
356     aws -profile user s3 cp  
357     s3://3kricegenome/REFERENCE/VCF_FILE.snp.vcf.gz* .
```

358 The subset region of the VCF file (chromosome:start-end) is then extracted using BCftools
359 (<http://samtools.github.io/bcftools/>) wrapped in Rice Galaxy and exported to the history pane as
360 bgzipped, indexed BCF file, which can then be converted back to VCF using [VCFTOOLS] in Rice Galaxy.
361 Standard methods for tool wrapper development and deployment were followed. All tool wrapper XMLs
362 developed specifically for Rice Galaxy are deposited and shared in a project-specific Rice Galaxy toolshed
363 repository at <http://52.76.88.51:8081/> (Figure 11) and will also be deposited in the central Galaxy
364 toolshed (<https://toolshed.g2.bx.psu.edu/>). All developments and testing of Rice Galaxy and Rice Galaxy
365 Toolshed were done in Docker containers hosted in virtual machines at the Advanced Science and
366 Technology Institute, Department of Science and Technology of the Philippine Government (ASTI –
367 DOST) prior to final deployment to the AWS instance.

368 This resource will empower the rice research community to benefit from publicly available datasets (e.g.
369 3K RG) and materials (seed/accessions), to enhance or even drive their own respective institutional
370 genetic/genomic/breeding efforts. The Rice Galaxy instance (data, tools, computing resources) is free for
371 use by all.

372 In addition to the integration of these tools, new Galaxy wrappers and visualization plugins are being
373 developed for visualizing chromosomes and their information (SNP density, structural variants,
374 translocations) either in linear or circular mode, using recent web technologies (Ideogram.js [42] ,
375 BioCircos.js [43], respectively).

376 Finally, a Docker container of Rice Galaxy is under development so that it can be easily shared and
377 deployed through the Galaxy Docker flavor initiative ([https://github.com/bgruening/docker-galaxy-](https://github.com/bgruening/docker-galaxy-stable)
378 [stable](https://github.com/bgruening/docker-galaxy-stable)).

379 **Conclusion**

380 Rice Galaxy is a federated Galaxy resource specialized for rice genetics, genomics, and breeding. The
381 resource empowers the rice research community to utilize publicly available datasets (3K RG), materials
382 (seed/accessions), and their own data, allowing complex data analyses to be performed even without
383 investment in their own computational infrastructure and software development team. Rice research –
384 related tools are also hosted in Rice Galaxy (i.e. Uniqprimer rice pathogen diagnostic design).
385 Rice Galaxy is freely accessible to all and we invite the rice research community to participate in
386 enriching the tools hosted by the resource. It can serve as a repository for data, analyses results, and
387 new bioinformatics tools coming from institutions that have used the publicly available rice diversity
388 panels from 3K RG, or have developed rice genomic/genetic analyses tools that they wish to share to the
389 community, and a computing infrastructure for small institutes without in-house computing capability.

390 **Availability and requirements**

391 Project name: RICE GALAXY

392 Project home page: <https://github.com/InternationalRiceResearchInstitute/RiceGalaxy>

393 Operating system(s): Linux Ubuntu release 12.04.2 LTS

394 Programming language: Python

395 Other requirements: R release 3.2.3 and following packages: methods, fpc, cluster, vegan, pheatmap,
396 pROC, randomForest, miscTools, pRF, e1076, rrBLUP, glmnet ;TASSEL release 5.2.40; plink v1.90b3k;
397 JBrowse 1.14.1; snpEff 4.3T; sNMF 1.2 (and as R package LEA); FastME 2.0

398 License: Rice Galaxy tools are released under GNU GPL. All software from external sources is bound by
399 their respective licenses.

400 Any restrictions to use by non-academics: Rice Galaxy tools are without restriction to non-academics. All
401 software from external sources is bound by their respective non-academic restrictions

402 Code availability: Tool wrappers at Rice Galaxy Toolshed (<http://52.76.88.51:8081/>). Rice Galaxy is
403 available at IRRI Github (<https://github.com/InternationalRiceResearchInstitute/RiceGalaxy>).

404 **Availability of supporting data**

405 3,000 Rice Genomes Project at Gigascience database (<http://gigadb.org/dataset/200001>)

406 3K RG BAM and VCF files available from Amazon Public data and ASTI-DOST IRODs site, instructions at

407 <http://iric.irri.org/resources/3000-genomes-project> .

408 SNP sets and morpho-agronomic characterization from 3K RG at SNP-Seek download site (<http://snp->

409 [seek.irri.org/ download.zul](http://snp-irri.org/download.zul))

410 **Availability of supporting source code and requirements**

411 Project name: Uniqprimer

412 Project home page: <https://github.com/SouthGreenPlatform/Uniqprimer>

413 Operating system(s): Linux OS

414 Programming Language: Python

415 Other requirements: MUMmer 3

416 License: GNU GPL

417 Project name: PRAGMA Data Service

418 Project home page: repository <https://github.com/Data-to-Insight-Center/RDA-PRAGMA-Data->

419 [Service/wiki/Welcome-to-PRAGMA-Data-Service-Prototype](https://github.com/Data-to-Insight-Center/RDA-PRAGMA-Data-Service/wiki/Welcome-to-PRAGMA-Data-Service-Prototype)

420 Operating system(s): Platform independent

421 License: Apache License 2.0

422 **Declarations**

423 **Abbreviations**

424 3K RG:3,000 Rice Genomes;HDRA: High Density Rice Array; SNP: single nucleotide polymorphism; GWAS:

425 Genome-Wide Association Studies; RAM:random access memory; CPU: central processing unit;

426 IRRI:International Rice Research Institute;NGS: next-generation sequencing;QTL:quantitative trait loci;

427 IRGSP:International Rice Genome Sequencing Project; RGAP:Rice Genome Annotation Project;KASP:

428 Kompetitive Allele Specific PCR; VCF : variant call format ;AWS: Amazon Web Services; RAVE: Rapid
429 Allelic Variant extractor; MAF: minor allele frequency;TASSEL: Trait Analysis by aSSociation, Evolution
430 and Linkage; GS:genomic selection;Oghma: Operators for Genome decipHering by MACHine learning;
431 rrBLUP:ridge regression best linear unbiased predictor; SVM:support vector machine; FST:fixation index;
432 NSF:National Science Foundation;CGIAR: Consultative Group for International Agricultural Research;
433 RDA: Research Data Alliance;PID:persistent identifier;OA:open access; PRAGMA:Pacific Rim Applications
434 and Grid Middleware Assembly; EC2:elastic computing cloud; CLI: command line interface; S3: Simple
435 Storage Service; API: Python Application Programming Interfaces; XML: eXtensible Markup Language;

436 **Competing interests**

437 The author(s) declare that they have no competing interests.

438 **Funding**

439 Components of the project are supported by the following grants: Taiwan Council of Agriculture Grant to
440 IRRI, International Rice Informatics Consortium, and CGIAR Excellence in Breeding Platform for financial
441 support to the Rice Galaxy main server, the USA National Science Foundation PRAGMA grant number:
442 NSF OCI 1234983, the RDA/US-sponsored adoption program funded by the MacArthur Foundation, and
443 the AIST ICT International Collaboration Grant.

444 **Authors' contributions**

445 VJ and AD equally contributed to create Rice Galaxy. NB contributed the genomic prediction tools. AD,
446 GD, PL, and MR contributed the RAVE and SNIPLAY tools. JD, JRM, JPP created the development Rice
447 Galaxy cloud instances hosted at DOST-ASTI. LM created the SNP-Seek interfaces. LT, JL, JEL contributed
448 the Uniqprimer tool. GZ, KR, BP, and JH contributed the Rice Galaxy Open Access, MT, NA, and TK
449 contributed to funding acquisition and writing, RM coordinated the conceptualization of the project and
450 the writing process.

451 **Acknowledgments**

452 The authors are grateful to the following people and institutions/agencies for their support: DOST-ASTI
453 for hosting the Rice Galaxy toolshed server, Jay Santos and Denis Diaz for assistance with AWS
454 architecture.

455 **References**

- 456 1. 3,000 rice genomes project. The 3,000 rice genomes project. *Gigascience*. 2014;3:7.
- 457 2. Wang, W-S, Mauleon R, Chebotarov, D, et al. Genomic variation in 3,010 diverse accessions of
458 Asian cultivated rice. *Nature*. 2018;557: 43–49 .doi:10.1038/s41586-018-0063-9.
- 459 3. McCouch S, Wright M, Tung C-W, Maron L, McNally K, Fitzgerald M, et al. Open Access
460 Resources for Genome Wide Association Mapping in Rice. *Nature Comm*. 2016;7: 10532, doi
461 10.1038/ncomms10532.
- 462 4. Alexandrov N, Tai S, Wang W, Mansueto L, Palis K, Fuentes RR, et al. SNP-Seek database of SNPs
463 derived from 3000 rice genomes. *Nucleic Acids Res*. 2015;63:2–6.
- 464 5. Mansueto L, Fuentes RR, Chebotarov D, Borja FN, Detras J, Abriol-Santos JM, et al. SNP-Seek II: A
465 resource for allele mining and analysis of big genomic data in *Oryza sativa*. *Curr. Plant Biol*.
466 2016;6628:16–25.
- 467 6. Sempéré G, Philippe F, Dereeper A, Ruiz M, Sarah G, Larmande P. Gigwa-Genotype investigator
468 for genome-wide analyses. *Gigascience*. 2016;5:25.
- 469 7. The South Green Collaborators. The South Green portal: a comprehensive resource for tropical
470 and Mediterranean crop genomics. *Curr Plant Biol. Elsevier*. 2016;7–8: 6–9.
471 doi:10.1016/J.CPB.2016.12.002.
- 472 8. Tello-Ruiz MK, Naithani S, Stein JC, Gupta P, Campbell M, et al. Gramene 2018: unifying
473 comparative genomics and pathway resources for plant research. *Nucleic Acids Res*. 2017;
474 PMID: 29165610. doi: 10.1093/nar/gkx1111.

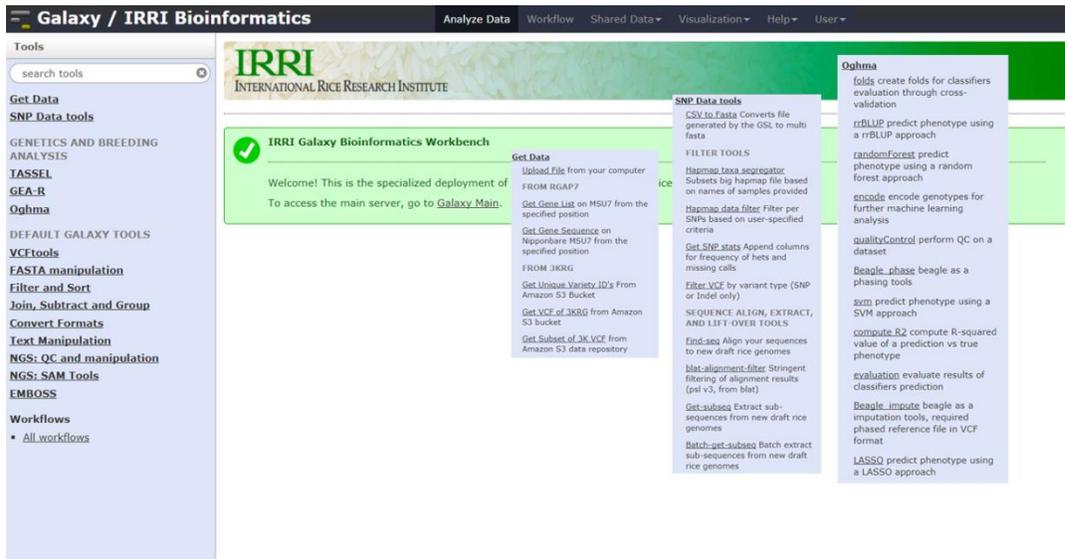
- 475 9. Afgan, E, Baker D. van den Beek M, Blankenberg D, Bouvier D, et al. The Galaxy platform for
476 accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids*
477 *Research*. 2016;44(W1): W3-W10 doi:10.1093/nar/gkw343.
- 478 10. Kawahara, T., et al. Improvement of the *Oryza sativa* Nipponbare reference genome using next
479 generation sequence and optical map data. *Rice*. 2013;6:4 .
- 480 11. Zhang J, Chen L-L, Xing F, Kudrna DA, Yao W, Copetti D, et al. Extensive sequence divergence
481 between the reference genomes of two elite *indica* rice varieties Zhenshan 97 and Minghui 63.
482 *Proc Natl Acad Sci U S A*. 2016;113: E5163–71.
- 483 12. Du, H, Yu Y, Ma Y, Gao Q, et al. Sequencing and *de novo* assembly of a near complete *indica* rice
484 genome. *Nature Communications*. 2017;8 (15324). doi:10.1038/ncomms15324.
- 485 13. Schatz, M., et al. Whole genome *de novo* assemblies of three divergent strains of rice, *Oryza*
486 *sativa*, document novel gene space of *aus* and *indica*. *Genome Biology*. 2014;15, 506.
- 487 14. Gao, Z.Y., et al. Dissecting yield-associated loci in super hybrid rice by resequencing recombinant
488 inbred lines and improving parental genome sequences. *PNAS*. 2013; 110 (35), 14492-14497.
- 489 15. Sakai, H., et al. Construction of pseudomolecule sequences of the *aus* rice cultivar Kasalath for
490 comparative genomics of Asian cultivated rice. *DNA Research*. 2014; do:10.1093/dnares/dsu006.
- 491 16. Xu K, Xu X, Fukao T, Canlas P, Maghirang-Rodriguez R, Heuer S, et al. *Sub1A* is an ethylene-
492 response-factor-like gene that confers submergence tolerance to rice. *Nature*. 2006;442: 705–
493 708.
- 494 17. Gamuyao R, Chin JH, Pariasca-Tanaka J, Pesaresi P, Catausan S, Dalid C, et al. The protein kinase
495 *Pstol1* from traditional rice confers tolerance of phosphorus deficiency. *Nature*. 2012;488: 535–
496 539.

- 497 18. Uga Y, Sugimoto K, Ogawa S, Rane J, Ishitani M, Hara N, et al. Control of root system
498 architecture by *DEEPER ROOTING 1* increases rice yield under drought conditions. Nat Genet.
499 2013;45: 1097–1102.
- 500 19. Thomson MJ, Singh N, Dwiyantri MS, Wang DR, Wright MH, et al. Large-scale deployment of a
501 rice 6 K SNP array for genetics and breeding applications. Rice. 2017;10:40 doi:10.1186/s12284-
502 017-0181-2.
- 503 20. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, et al. PLINK: a toolset for whole-
504 genome association and population-based linkage analysis". American Journal of Human
505 Genetics. 2007. 81: 559–75. doi:10.1086/519795.
- 506 21. Skinner ME, Uzilov AV, Stein LD, Mungall CJ, Holmes IH. JBrowse: a next-generation genome
507 browser. Genome Res. 2009;19:1630–8.
- 508 22. Dereeper A, Homa F, Andres G, Sempere G, Sarah G, Hueber Y, et al. SNIPlay3: a web-based
509 application for exploration and large scale analyses of genomic variations. Nucleic Acids Res.
510 2015;43:W295-300.
- 511 23. Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. TASSEL: Software for
512 association mapping of complex traits in diverse samples. Bioinformatics. 2007; 23:2633-2635.
- 513 24. Spindel J, Begum H, Akdemir D, Virk P, Collard B, Redona E, Atlin G, Jannink JL, McCouch SR.
514 Genomic selection and association mapping in rice (*Oryza sativa*): effect of trait genetic
515 architecture, training population composition, marker number and statistical model on accuracy
516 of rice genomic selection in elite, tropical rice breeding lines. PLOS Genetics. 2015; 11(2):
517 e1004982. doi:10.1371/journal.pgen.1004982.
- 518 25. Browning BL, Browning SR. Genotype imputation with millions of reference samples. Am J Hum
519 Genet. 2016. 98:116-126. doi:10.1016/j.ajhg.2015.11.020.

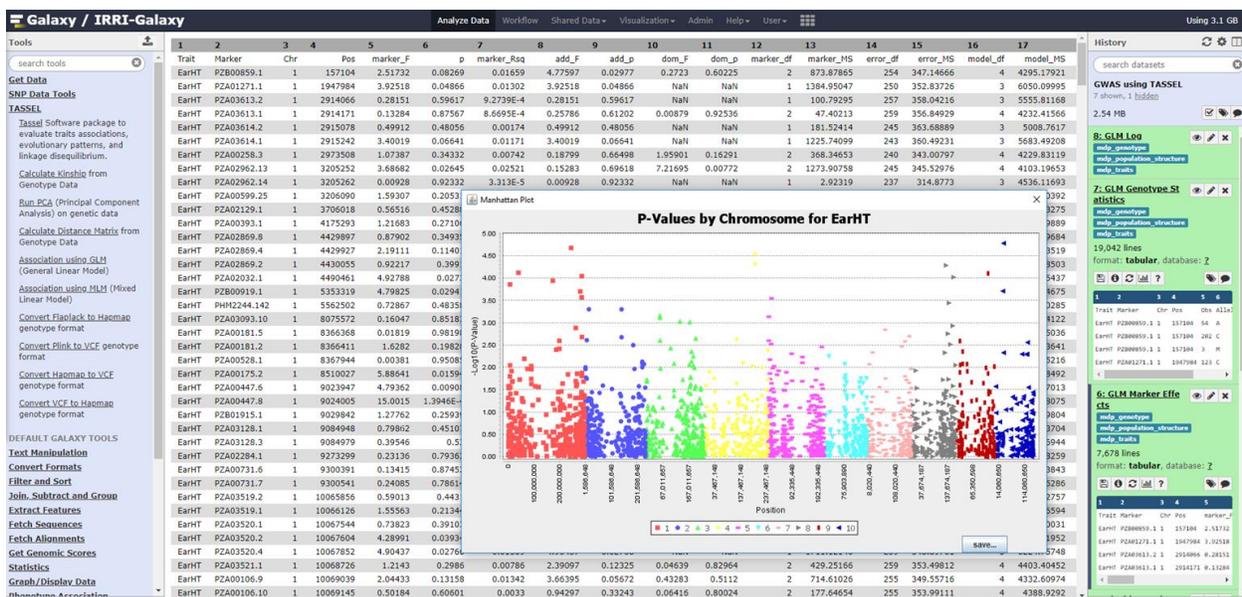
- 520 26. Rice Genome Annotation Project (RGAP) release 7. 2013; <http://rice.plantbiology.msu.edu/>.
521 Accessed 3 May 2018.
- 522 27. Langlois, PA, Snelling J, Hamilton JP, Bragard C, Koebnik R, Verdier V, et al. Characterization of
523 the *Xanthomonas translucens* complex using draft genomes, comparative genomics,
524 phylogenetic analysis, and diagnostic LAMP assays. *Phytopathology*. 2017; 107: 519-527.
- 525 28. Triplett, L, Hamilton JP, Buell CR, Tisserat NA, Verdier V, Zink F, Leach JE. Genomic Analysis of
526 *Xanthomonas oryzae* from US rice reveals substantial divergence from known *X. oryzae*
527 pathovars. *Appl. Environ. Microbiol.* 2011.;77(12):3930-3937. doi:10.1128/AEM.00028-11.
- 528 29. Lang, JM, Langlois P, Nguyen MHR, Triplett LR, Purdie L, et al. Sensitive detection of
529 *Xanthomonas oryzae pv. oryzae* and *X. oryzae pv. oryzicola* by Loop-Mediated Isothermal
530 Amplification. *Applied and Environmental Microbiology*. 2014; 80:4519-4530.
- 531 30. Triplett L, Verdier V, Campillo T, Van Malderghem C, et al. Characterization of a novel clade of
532 *Xanthomonas* isolated from rice leaves in Mali and proposal of *Xanthomonas maliensis* sp. nov.
533 2015; *Antonie van Leeuwenhoek* 107:869-81. doi: 10.1007/s10482-015-0379-5.
534 <http://link.springer.com/journal/10482/onlineFirst/page/1>.
- 535 31. Ash GJ, Lang JM, Triplett LR, Stodart BJ, Verdier V, et al. Development of a genomics-based
536 LAMP (Loop-1 mediated isothermal amplification) assay for detection of *Pseudomonas*
537 *fuscovaginae* from rice. 2014.; *Plant Dis* 98:909-915 doi.org/10.1094/PDIS-09-13-0957-RE .
- 538 32. Yogesh LS, Plale B, Gannon D. A survey of data provenance in e-science. *ACM Sigmod Record* .
539 2005;34.3, p. 31-36.
- 540 33. Zhou Q, Ghoshal D ,Plale B. Study in Usefulness of Middleware-Only Provenance. 2014 IEEE 10th
541 International Conference on e-Science, Sao Paulo. 2014; pp. 215-222.
542 doi:10.1109/eScience.2014.49.

- 543 34. Suriarachchi I, Zhou Q, Plale B, Komadu. A Capture and Visualization System for Scientific Data
544 Provenance . Journal of Open Research Software. 2015;3 p . e4,. doi:10.5334/jors.bq.
- 545 35. Goecks J, Nekrutenko A, Taylor J. Galaxy: a comprehensive approach for supporting accessible,
546 reproducible, and transparent computational research in the life sciences. Genome biology.
547 2010; 11(8):R86.
- 548 36. Gaignard A, Belhajjame K, Skaf-Molli H. Sharp: Harmonizing and bridging cross-workflow
549 provenance. In: Blomqvist E, Hose K, Paulheim H, Lawrynowicz A, Ciravegna F, Hartig O, editors.
550 The Semantic Web: ESWC 2017 Satellite Events. 2017.; p. 219-234, Cham. Springer International
551 Publishing.
- 552 37. Kanwal S, Zaib Khan F, Lonie A, Sinnott RO. Investigating reproducibility and tracking provenance
553 - a genomic workflow case study. BMC bioinformatics. 2017; 18(1):337.p. 25.
- 554 38. Missier P, Ludascher B, Dey S, Wang M, McPhillips T, et al. Golden trail: Retrieving the data
555 history that matters from a comprehensive provenance repository. International Journal of
556 Digital Curation. 2012; 7(1):139-150.
- 557 39. Kahn R, Wilensky R. A Framework for Distributed Digital Object Services. Int. J. Digit. Libr. 2006;
558 6, 2: 115–123. doi:10. 1007/s00799- 005- 0128- x.
- 559 40. Research Data Alliance PID Kernel Information Working Group. PID Kernel Information guiding
560 principles. 2018; [https://www.rd-alliance.org/group/pid-kernel-information-wg/wiki/pid-](https://www.rd-alliance.org/group/pid-kernel-information-wg/wiki/pid-kernel-information-guiding-principles)
561 [kernel-information-guiding-principles.](https://www.rd-alliance.org/group/pid-kernel-information-wg/wiki/pid-kernel-information-guiding-principles)). Accessed 15 May-2018.
- 562 41. Research Data Alliance Data Type Registry Working Group. RDA Data Type Registries Working
563 Group Output. 2016; doi:10.15497/ A5BCD108-ECC4-41BE-91A7-20112FF77458. Accessed 15
564 May 2018.
- 565 42. Dereeper A, Bocs S, Rouard M, Guignon V, Ravel S, Tranchant-Dubreuil C, et al. The coffee
566 genome hub: a resource for coffee genomes. Nucleic Acids Res. 2015; 43:D1028-35.

567 43. Cui Y, Chen X, Luo H, Fan Z, Luo J, He S, et al. BioCircos.js: an interactive Circos JavaScript library
 568 for biological data visualization on web applications. *Bioinformatics*. 2016; 32(11):1740-2.
 569 doi:10.1093/bioinformatics/btw041.
 570



571
 572 Figure 1. Rice Galaxy @ IRR with customized analyses tools for genetics, breeding, and custom data
 573 sources (i.e. 3,000 Rice Genomes project).
 574

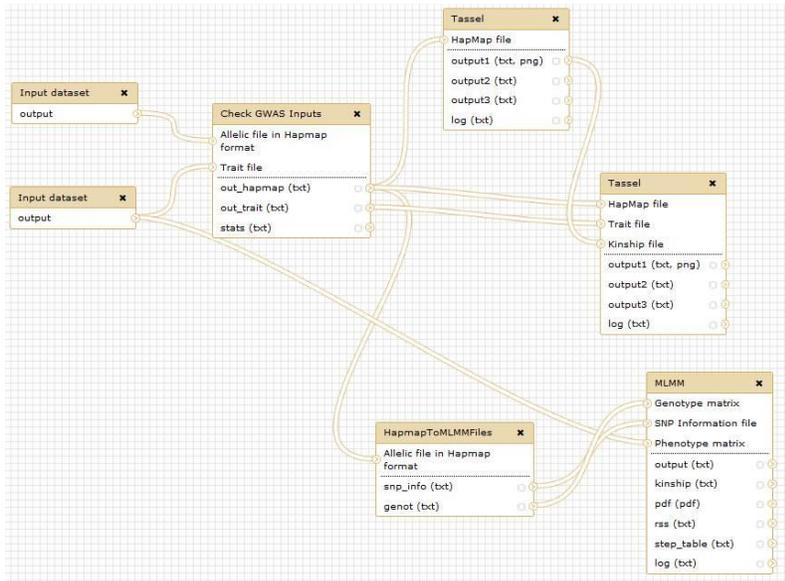


575
 576
 577
 578
 579
 580
 581
 582
 583
 584
 585
 586
 587
 588
 589
 590
 591
 592
 593
 594
 595
 596
 597
 598
 599
 600
 601
 602
 603
 604
 605
 606
 607
 608
 609
 610
 611
 612
 613
 614
 615
 616
 617
 618
 619
 620
 621
 622
 623
 624
 625
 626
 627
 628
 629
 630
 631
 632
 633
 634
 635
 636
 637
 638
 639
 640
 641
 642
 643
 644
 645
 646
 647
 648
 649
 650
 651
 652
 653
 654
 655
 656
 657
 658
 659
 660
 661
 662
 663
 664
 665
 666
 667
 668
 669
 670
 671
 672
 673
 674
 675
 676
 677
 678
 679
 680
 681
 682
 683
 684
 685
 686
 687
 688
 689
 690
 691
 692
 693
 694
 695
 696
 697
 698
 699
 700
 701
 702
 703
 704
 705
 706
 707
 708
 709
 710
 711
 712
 713
 714
 715
 716
 717
 718
 719
 720
 721
 722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811
 812
 813
 814
 815
 816
 817
 818
 819
 820
 821
 822
 823
 824
 825
 826
 827
 828
 829
 830
 831
 832
 833
 834
 835
 836
 837
 838
 839
 840
 841
 842
 843
 844
 845
 846
 847
 848
 849
 850
 851
 852
 853
 854
 855
 856
 857
 858
 859
 860
 861
 862
 863
 864
 865
 866
 867
 868
 869
 870
 871
 872
 873
 874
 875
 876
 877
 878
 879
 880
 881
 882
 883
 884
 885
 886
 887
 888
 889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917
 918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971
 972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000

576

577 Figure 2. Genome-Wide Association Studies analysis (implemented by TASSEL software) in Rice Galaxy.

578



579

580 Figure 3. Genome-Wide Association Studies analysis workflow in SNIPlay as implemented in Rice Galaxy.

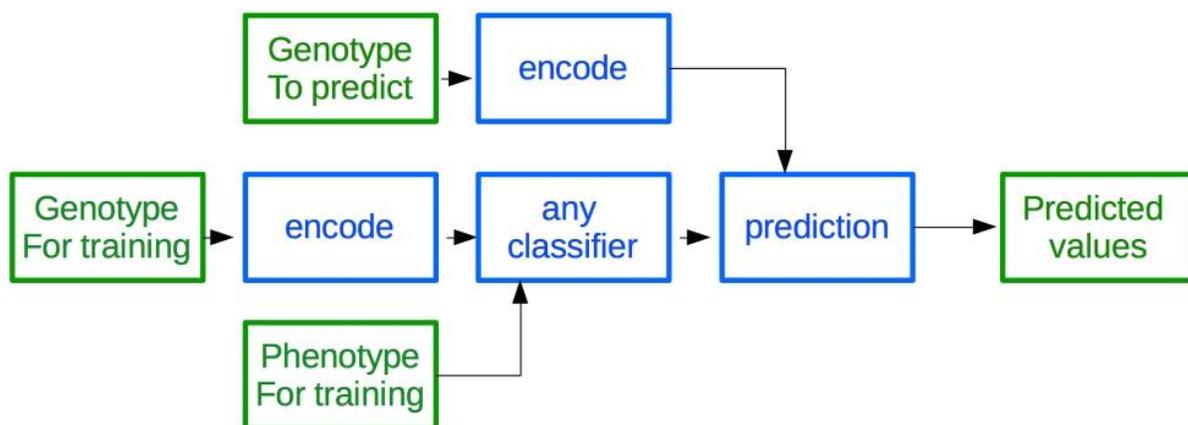
581

The screenshot shows the Galaxy / IRRi Bioinform interface. The top navigation bar includes 'Galaxy / IRRi Bioinform', 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Help', and 'User'. On the left, a 'Tools' sidebar lists various tools under the 'Oghma' category, including 'folds', 'rrBLUP', 'randomForest', 'encode', 'qualityControl', 'Beagle_phase', 'svm', 'compute R2', 'evaluation', 'Beagle_impute', and 'LASSO'. On the right, a green banner reads 'IRRI Galaxy Bioinformatics Workbench' with a checkmark icon, followed by a welcome message: 'Welcome! This is the specialized deployment of Galaxy at the International Rice Research Institute (IRRI). To access the main server, go to [Galaxy Main](#).'

582

583 Figure 4. Oghma genomic prediction and selection tools in Rice Galaxy with various classifier tools
584 installed.

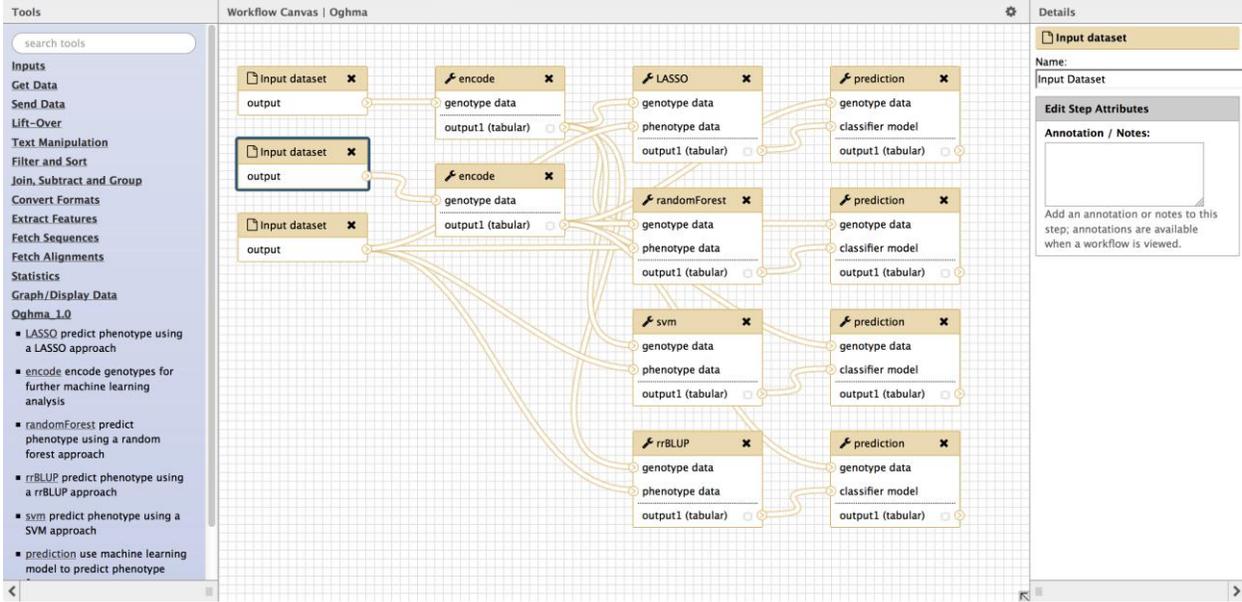
585



586

587 A. Overview of the Genomic Selection analyses workflow as implemented in Oghma tool suite.

588

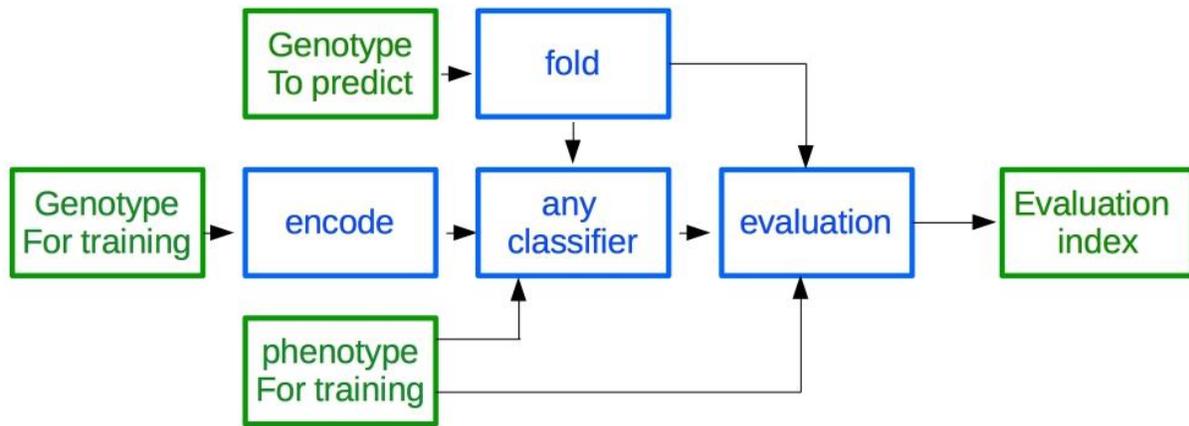


589

590 B. Rice Galaxy workflow for genome prediction using Oghma tool suite.

591 Figure 5: Genomic Selection analyses workflow as implemented by Oghma tool suite.

592

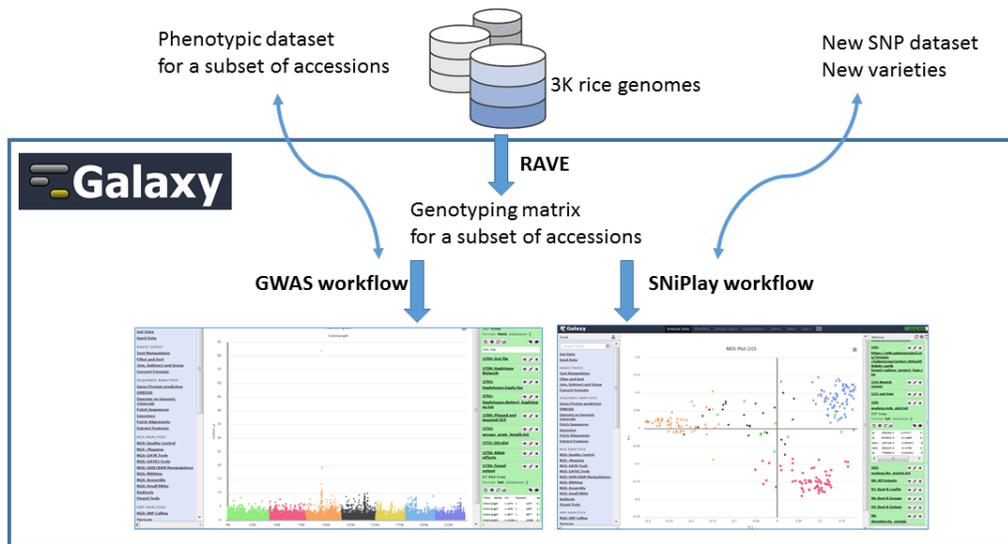


593

594

595 Figure 6. Workflow for classifier evaluation in the genome prediction tool suite implemented by Oghma.

596



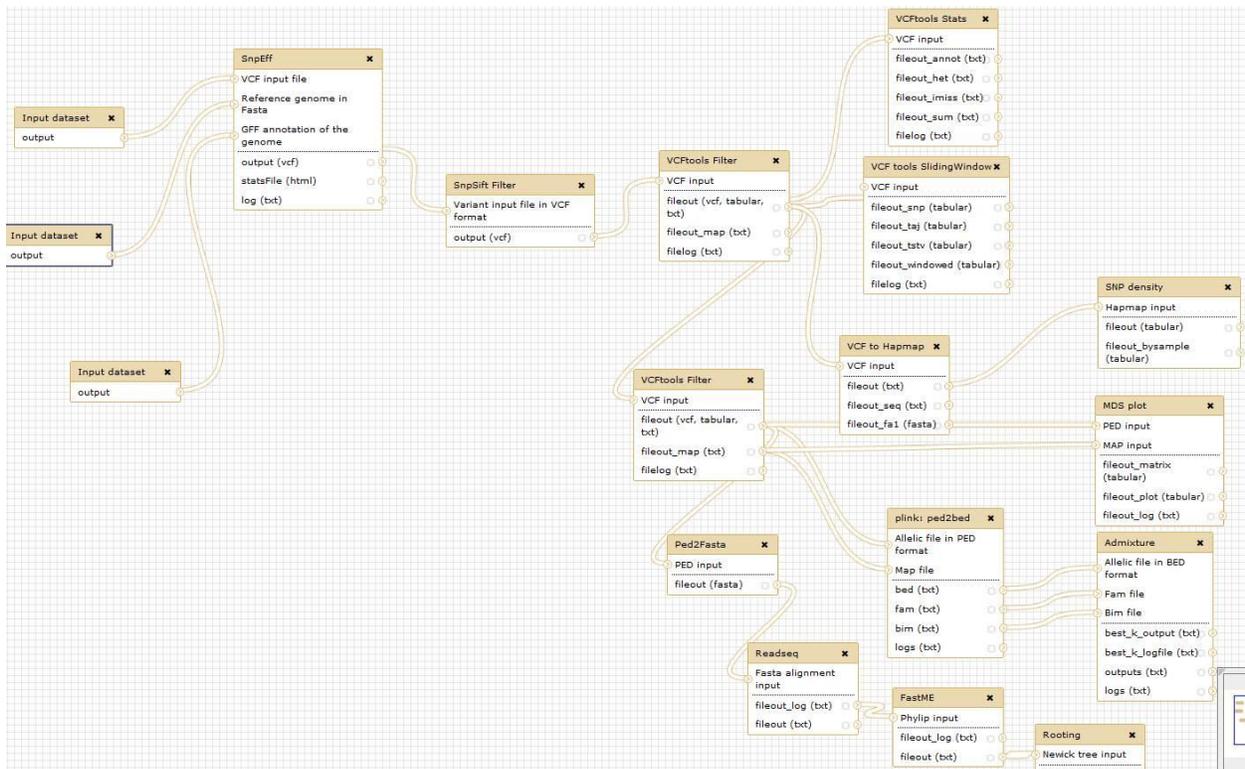
597

598 Figure 7. Overview schematic showing the integration of the 3K Rice Genomes project genotyping

599 database and rapid extraction of subset SNPs by RAVE module for use by analyses workflows installed in

600 Rice Galaxy.

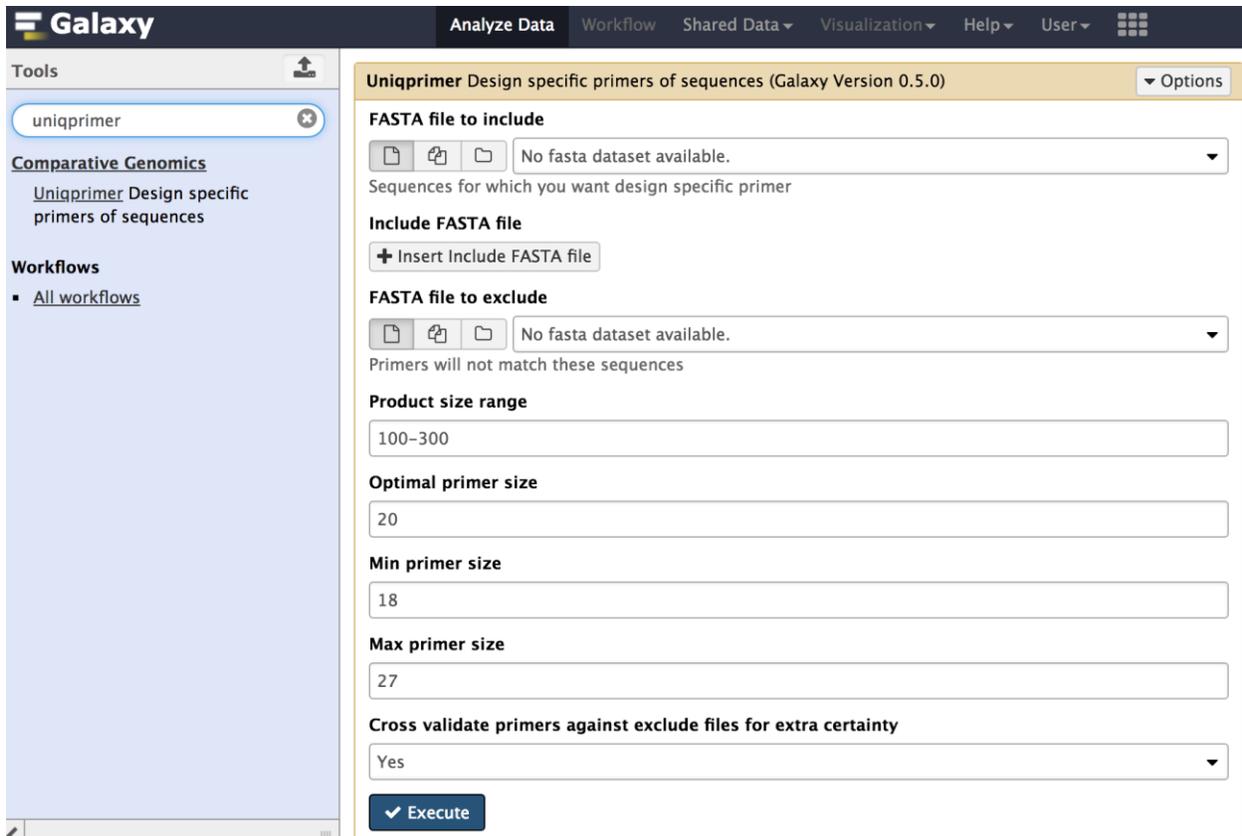
601



602

603 Figure 8. Rice Galaxy SNiPlay workflow for diversity and population structure analyses using various
604 software tools.

605

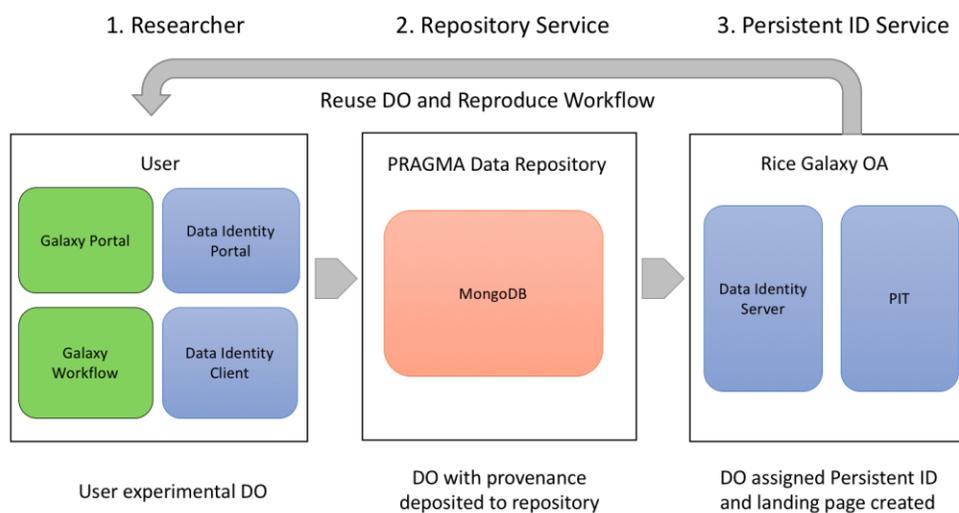


606

607 Figure 9. Uniqprimer comparative genomics-based diagnostic primer design tool for microbial pathogen

608 detection installed in Rice Galaxy.

609



610

611 A. The underlying software infrastructure for the components of Rice Galaxy Open Access.

Figure 10A: Galaxy Workflow

Inputs: HapMap file, Trait file, Kinship file

Tools: Tassel

Outputs: output1 (txt, png), output2 (txt), output3 (txt), log (txt)

mml.json:

```

{
  "a_galaxy_workflow": "true",
  "annotation": "",
  "format-version": "0.1",
  "name": "MLM",
  "steps": [
    {
      "annotation": "",
      "id": 0,
      "input_connections": {},
      "inputs": [
        {
          "description": "",
          "name": "HapMap file"
        }
      ],
      "name": "Input dataset",
      "outputs": [],
      "position": {
        "xact": 295.5,
        "yop": 138.5
      },
      "tool_errors": null,
      "tool_id": null,
      "tool_state": "{}",
      "tool_version": null,
      "type": "data_input",
      "user_outputs": []
    },
    {
      "id": 1
    }
  ]
}

```

Figure 10B: IRRI Data Repository beta

Data Repository beta

Search Data Object Persistent Identifier (PID)

Filter by

Creator: [Enter creator name]

Time Range: [Choose time range]

Workflow: [Enter workflow]

Data Type - IRRI Rice Genomes tassel workflow

12 Data Object found in 532 milliseconds

DO Name	Creator	Timestamp	Download Data Object	Download Metadata Object
kunalan_glm_2018-04-27-20:19:36	kunalan	2018-04-27-20:19:36	Download Data Object	Download Metadata Object
luoyu_glm_2018-04-27-20:47:29	luoyu	2018-04-27-20:47:29	Download Data Object	Download Metadata Object
luoyu_glm_2018-04-30-18:42:26	luoyu	2018-04-30-18:42:26	Download Data Object	Download Metadata Object
luoyu_glm_2018-04-30-19:53:18	luoyu	2018-04-30-19:53:18	Download Data Object	Download Metadata Object

MLM DO Metadata

luoyu_mlm_2018-04-30-18:42:26

11723/362b15b1-9334-4f0d-8e6d-8fb3d8d55e45

DO Name: luoyu_mlm_2018-04-30-18:42:26

Creator: luoyu

Timestamp: 2018-04-30-18:42:26

612

613 B. Digital Object flow in Rice Galaxy Open Access. A Galaxy analysis workflow (exported as JSON file) is
 614 deposited to the DO repository, and the data identity server publishes the deposited DO + meta-data for
 615 discoverability.

616 Figure 10. The components (A) and the flow of Digital Objects from upload to discoverability (B) in the
 617 prototype Rice Galaxy Open Access.

618

Galaxy Tool Shed

26 valid tools on Feb 08, 2018

Repositories Help User

Categories

search repository name, description

Name	Description	Repositories
EIB Hackathon 2018	tools from 2018 Hackathon	7
File Conversion	tools for converting file formats	1
Genomic Selection	tools for the genomic selection project	6
GWAS	gwas tools	1
Sequence	repositories for design, validation and analyses of sequences	2
SNP Calling	tools for SNP Discovery	2

Search

- Search for valid tools
- Search for workflows

Valid Galaxy Utilities

- Tools
- Custom datatypes
- Repository dependency definitions
- Tool dependency definitions

All Repositories

- Browse by category

Available Actions

- Login to create a repository

619

620

621 Figure 11. Rice Galaxy Toolshed with the various available tools.

622