

1 **Introducing CHiDO – a No Code Genomic Prediction Software implementation for the**
2 **Characterization & Integration of Driven Omics**

3
4 ¹Francisco González, ^{1,2}Julián García-Abadillo, & ^{1*}Diego Jarquín

5
6 ¹ Agronomy Department, University of Florida, Gainesville, FL 32611

7
8 ² Centro de Biotecnología y Genómica de Plantas, Universidad Politécnica de Madrid (UPM),
9 Campus de Montegancedo, Pozuelo de Alarcón, Madrid, 28223, Spain

10
11
12 *** Correspondence:**

13 Corresponding Author

14 jhernandezjarqui@ufl.edu

15
16
17 **Keywords:** Multi-Omics Integration, Genomic Selection - Genomic Prediction, R Shiny,
18 Climate Adaptation, Low-code-no-code (LCNC), Bayesian Statistics, High Dimensional
19 Interactions.

20 **Core ideas:**

- 21 1. The authors developed CHiDO, a platform for breeders to build predictive models
22 integrating multi-omics data.
23 2. CHiDO is a no-code tool that leverages the reaction norm model proposed by Jarquin et
24 al. (2014).
25 3. The platform aims to increase access to predictive analytics lowering relevant technical
26 and financial barriers.
27

28
29
30
31
32
33
34
35
36

37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53

54

55

56
57
58
59
60
61
62

63
64
65
66
67

ABSTRACT

Climate change represents a significant challenge to global food security by altering environmental conditions critical to crop growth. Plant breeders can play a key role in mitigating these challenges by developing more resilient crop varieties; however, these efforts require significant investments in resources and time. In response, it is imperative to use current technologies that assimilate large biological and environmental datasets into predictive models to accelerate the research, development, and release of new improved varieties. Leveraging large and diverse data sets can improve the characterization of phenotypic responses due to environmental stimuli and genomic pulses. A better characterization of these signals holds the potential to enhance our ability to predict trait performance under changes in weather and/or soil conditions with high precision. This paper introduces CHiDO, an easy-to-use, no-code platform designed to integrate diverse omics datasets and effectively model their interactions. With its flexibility to integrate and process data sets, CHiDO's intuitive interface allows users to explore historical data, formulate hypotheses, and optimize data collection strategies for future scenarios. The platform's mission emphasizes global accessibility, democratizing statistical solutions for situations where professional ability in data processing and data analysis is not available.

68

1 INTRODUCTION

69 As the global population continues to surge, projected to reach 10 billion by 2050 (Gu et al.,
70 2021), the imperative to increase agricultural yields becomes increasingly critical (Van Dijk et
71 al., 2021). This challenge is compounded by the escalating frequency and intensity of weather
72 variations due to climate change, posing a significant threat to food security worldwide (Lesk et
73 al., 2016). Such climatic extremes have already begun to impact the productivity of elite crop
74 varieties, with studies indicating a potential reduction of up to 6% for an increase of one degree
75 Celsius in average temperature (Zhao et al., 2017).

76 To meet these rising demands for available food products, agricultural production must
77 increase, and supply chain improvements must be achieved to reduce food waste at different
78 stages. Plant breeding can play a pivotal role in increasing total harvestable output through the
79 development of improved genotypes that can withstand changing environmental conditions.
80 More resilient crop varieties serve dual purposes: 1) fulfilling the nutritional demands of a
81 growing population, and 2) mitigating reliance on environmentally harmful inputs like fossil
82 fuels and synthetic agrochemicals (Foley et al., 2011). However, traditional plant breeding
83 methods, which predominantly rely on phenotypic and pedigree data, are resource-intensive and
84 time-consuming (Atlin et al., 2017). These conventional approaches require significant
85 investments in land and time, often taking up to eight years to develop a new variety for annual
86 crops (Jarquín et al., 2017). Moreover, genetic engineering, while a potential solution, is
87 surrounded by socio-economic and ecological concerns, as well as issues of accessibility,
88 corporate control and public acceptance (Clapp, 2018; Tsatsakis et al., 2017).

89 Recent advances in sequencing technologies have revolutionized our ability to
90 characterize genotypes with high precision through DNA-based marker profiles (Varshney et al.,
91 2014). This genomic information enables the characterization of genetic relationships between
92 individuals (Bernardo 1994), forming the foundation of genomic selection (GS). However, this
93 selection framework has its own set of challenges, such as handling large genomic datasets for
94 reduced number of phenotypic observations --the "large p , small n " problem--. Leveraging
95 genomic selection-by-genomic prediction (GS-GP) techniques allows the prediction of the
96 performance of unobserved genotypes based solely on their marker profiles (Meuwissen et al.,
97 2001). This approach, although a significant leap in breeding efficiency, overlooks the impact of
98 environmental factors. The next iteration of computational methods to accelerate and improve
99 breeding efforts was modeling Genotype-by-Environment ($G \times E$) interactions.

100 $G \times E$ analysis examines how genotypes respond to different environmental conditions
101 (change in the response patterns - rankings). However, similar problems ($p \gg n$) than for
102 conventional GS-GP models arise when considering the interaction between genes and
103 environmental factors increasing the computational demands of modeling a large number of
104 interactions via contrasts (Cossa et al., 2017). Utilizing the approach proposed by Jarquín et al.

105 (2014), we can overcome these challenges by analyzing and integrating all first degree
106 interactions between marker SNPs and environmental covariates via covariance
107 matrices/structures. This alternative significantly reduces the dimensionality of data by
108 leveraging correlations between genotype-by-environment combinations that are similar both
109 genetically and at the level of the environmental covariates (environmentally) rather than
110 computing individual contrast effects between markers and weather factors (Jarquín et al.,
111 2014). Several studies have shown the advantages of taking into consideration the G×E
112 interaction in prediction models in plant and animal breeding applications (Jarquin et al., 2020,
113 Tiezzi et al., 2017). The predictive power of the G×E interaction can be bolstered through the
114 inclusion of a broad spectrum of omics (or layers) data (e.g. genomics, proteomic, metabolomics,
115 enviromics, ionomics, high-throughput data, etc.), known as multi-omics analysis (Yang et al.,
116 2021).

117 Implementing models that effectively integrate and interpret this complex multi-omics
118 data can be challenging, often requiring specialized programming and statistical expertise that
119 may not be readily available in many breeding programs around the world, especially in
120 developing regions. To address this gap, we have developed CHiDO, a no-code platform
121 designed to facilitate the integration of multi-omics data to build, train and test complex G×E
122 prediction scenarios.

123 Across several Latin-American cultures, the word ‘chido’ (meaning ‘cool’ in English) is
124 a powerful and oversimplistic expression that succinctly describes all the positive aspects of an
125 action, event, thing, etc. In our case, CHiDO stands for **C**haracterization and **I**ntegration of
126 **D**riven **O**micS, and it enables breeders to use advanced analytical methods without having to
127 write code themselves; removing a technical barrier and democratizing access to the latest
128 predictive analytics used in breeding implementations. Our CHiDO development is not just a
129 “prediction software”, it also integrates a series of developments proposed by several Latin
130 American scientists (Drs. de los Campos, Crossa, Perez-Rodriguez, Gianola) that are recurrently
131 cited along this paper, and this is a way to acknowledge their contributions in the field. In this
132 paper, we discuss the development of this platform, its components and the statistical methods
133 leveraged for its functionality. Currently, the application can be accessed at
134 <https://jarquinlab.shinyapps.io/chido/>.

135

136 **2 MATERIALS AND METHODS**

137 **2.1 Platform Overview**

138 The implementation of elaborated prediction models, integrating data from multiple omics of
139 information (including interactions of several types), and their corresponding evaluation

140 considering different prediction scenarios (mimicking realistic scenarios) poses extra challenges
141 in many breeding programs.

142

143 CHiDO is a no-code platform that can fill this gap by allowing breeders to build, train, and
144 validate linear models that incorporate data derived from multiple omics of information in a
145 simple manner. It also addresses two challenges with leveraging G×E interaction models for
146 breeding efforts by 1) using a UI-based workflow to overcome the technical barrier associated
147 with multi-omics data handling and programming, and 2) reducing the dimensionality of G×E by
148 adopting the reaction norm model described in Jarquín et al., (2014) which is further described in
149 the *Statistical Background* subsection. The platform's user interface (UI) is divided into four
150 sections –data loading, model assembly, training/validation, and results view where each section
151 contains instructions and widgets to customize the metadata, parameters and model equations as
152 necessary. The drag-and-drop interface is a novel approach to building complex models where
153 users can add individual omics as main effects and form interactions between them (e.g., G×E)
154 by *collapsing* these effects without requiring any advanced programming knowledge.

155

156 Other key features of CHiDO include: 1) customizable data processing and parameter
157 tuning, 2) handling multiple input files within a single session, 3) viewing omics data and editing
158 associated metadata, 4) building and testing multiple models in a single session, 5) viewing
159 results in the UI with the option to download them as CSV and PNG files, and 6) exporting
160 models as R objects via an RDS file. These features and additional functionality are split into
161 four separate page views within the CHiDO platform (**Figure 1**).

162 **2.1.1 Design**

163

164 CHiDO was built using the *Shiny* framework (Winston Chang et al., 2023), a popular R package
165 for creating interactive web and desktop applications. Its architecture, however, diverges from
166 the typical UI-Server split typically found in most *Shiny* applications, emphasizing a modular
167 design methodology. This approach involves segmenting logical components into individual
168 functions to enhance the platform's long-term maintainability, support and expansion. The
169 software's architecture is based on modern development practices to prevent logical duplication,
170 reduce dependencies within the codebase, and minimize disruption as new versions are released.

171

172 Consequently, CHiDO's design integrates both R and JavaScript in its frontend and backend
173 logic. The packages *shinyjs* (Dean Attali, 2021) and *shinyjqui* (Yang Tang, 2022) are utilized to
174 introduce functionality that extends beyond the traditional capabilities of R/Shiny, including the
175 drag-and-drop interface for model assembly (Figure 2). Many features are made available by
176 leveraging a suite of R packages such as *ggplot2* (Hadley Wickham, 2016) for rich data
177 visualizations and *DT* (Yihui Xie et al., 2023) to display and handle tabular objects. The

178 selection of *DT* is deliberate, enabling table and data frame manipulation with either JavaScript
179 or R code. For an exhaustive list of libraries used by CHiDO, see Table 1.

180 2.1.2 Usage

181

182 The typical workflow for CHiDO (Figure 3) can be listed as the following steps:

183

184 **Step 1:** CHiDO accommodates the upload of tabular data in CSV and RDA format, with a
185 flexible approach to data requirements. The primary necessity is the phenotypic response file
186 (Y), which must contain columns for environment IDs, genotype IDs, and the target trait to
187 predict, at minimum. Dealing with data from multi-environment trials, the column corresponding
188 to the ID of the environments, and the genotypes should be specified in the interface. Also, if
189 omic data is collected at the plant or sample level (e.g. multispectral data collected with drones,
190 ionomic data, information collected on secondary traits, etc.) a column serving as a unique
191 identifier (*compound*) for that level should be included for alignment.

192

193 For omics data, each file must have an identifier column specified by the user to link back to the
194 matrix of phenotypes in the Y file; this could be a column with the genotype IDs, an environment
195 ID column, or a unique identifier (UID-*compound*) column (e.g. genotype-in-environment
196 combination, plant or sample ID). Once a data file is uploaded, users can modify its metadata
197 including its display label and linkage type (Environment ID, Genotype / Line ID, Compound
198 ID).

199

200 Users are responsible for ensuring their data is properly formatted prior to uploading it to
201 CHiDO. Extra care should be taken to ensure that all identifier levels of an omic are represented
202 in the Y file. For example, if molecular information is uploaded, all lines referenced in this
203 dataset should be present in the Y file as part of the Genotype / Line ID column, even if the
204 corresponding phenotypic values are missing. If these identifiers are not consistent across both
205 files, the covariances matrices cannot be constructed for the implementation to work.

206 **Step 2:** Upon upload, each file is recognized as a separate omic within CHiDO and is assigned a
207 unique label, if not specified by the user during upload. In the model assembly section, these
208 labels appear as draggable elements for the user to add as main effects into a linear model
209 formula box. Interaction effects can be added by dragging *-collapsing-* two or more of these
210 labels into the same box before adding them into the formula. Users can build and save as many
211 models as desired, facilitating comparative analysis of these to determine which set of effects can
212 best predict trait performance (e.g., including G×E interactions) for desired phenotypic
213 expression.

214

215 **Step 3:** In the training and validation section, users have the option to adjust convergence hyper-
216 parameters (e.g., number of iterations and burn-in) and data pre-processing steps on genomic

217 data (i.e., quality control – minor allele frequency and percentage of missing values) at their
218 discretion. These settings can be altered for each model or applied uniformly across the multiple
219 models created in the previous section. Once a model is selected, it can be tested with one or
220 more of the four distinct cross-validation (CV) schemes that mimic prediction scenarios of
221 interests to breeders; 1) CV2: predicting tested genotypes in observed environments; 2) CV1:
222 predicting untested genotypes in observed environments; 3) CV0: predicting tested genotypes in
223 new environments; and 4) CV00: predicting untested genotypes in new environments (Persa et
224 al. 2021). The implementation of the declared linear predictors (models) is done using the BGLR
225 (Bayesian Generalized Linear Regression) package developed by Perez and de los Campos
226 (2014).

227
228 **Step 4:** The results of the selected CV schemes are presented in the UI in both tabular and
229 graphical outputs, with the option to download these locally. The downloadable results are
230 delivered in a compressed zip folder where the contents are systematically sorted by model, with
231 each model's folder containing CSV files with the raw numeric data for each CV and PNG files
232 showing a graphical representation of the results. In addition to the CV data, the results also
233 include evaluation metrics to assist with model interpretation efforts and the corresponding
234 variance components derived from the full data analysis. The metrics available are *prediction*
235 *accuracy* (as the Pearson correlation between predicted and observed values), *root-mean-*
236 *squared-error* (RMSE), and *variance components* to evaluate the relative contribution of each
237 one of the omics to explain the phenotypic variability. The formulae for each metric are provided
238 in the *Statistical methods* section. The output and evaluation metrics are displayed in both tabular
239 and graphical formats. This data is available to view as overall model performance (Figure 4) or
240 split by environment (Figure 5).

241 **2.2 Statistical Background**

242
243 Modeling high-dimensional sets of covariates p (e.g., genomic, environmental, gene \times
244 environment interactions, etc.) using a reduced set of n phenotypic observations such that $p \gg$
245 n , poses extra challenges. Especially under the conventional prediction approaches based on
246 linear regressions of the ordinary least squares (OLS) framework. The phenotypic response y_i of
247 the i^{th} genotype ($i = 1, 2, \dots, n$) can be represented as the linear combination between p markers
248 x_{ij} ($j = 1, 2, \dots, p$) and their corresponding effects b_j such that

$$249 \quad 250 \quad 251 \quad y_i = \sum_{j=1}^p x_{ij} b_j + \varepsilon_i \quad [1]$$

252 Then, under the OLS framework the solution for the vector of marker effects is given by

$$253 \quad 254 \quad 255 \quad \hat{b} = (X'X)^{-1}X'y$$

256 A major challenge to obtain the solution of the vector of marker effects is the inversion of
257 singular matrices of the form $(XX')^{-1}$ which are not full rank due to the larger number of
258 coefficients to estimate (p) with respect to the reduced number of data points (n) available for
259 model fitting. Under the parametric context, several statistical approaches have been developed
260 to deal with the curse of the dimensionality ($p \gg n$). Two of the most popular statistical
261 frameworks are the penalized regressions and the Bayesian methods which in many cases are a
262 sort of Bayesian versions of the former ones. By design, the penalized methods delimit to n the
263 total number features or covariates to select in the model.

264

265 On the other hand, the Bayesian methods consider distributional assumptions of the
266 marker effects, allowing (in principle) all features to be included in the final prediction model. In
267 both cases, the inversion of matrices with large dimensions ($p \times p$) is accomplished by adding a
268 value to the diagonal elements of the (XX') matrix to “break” the singularity. Another option is
269 to consider prior distributions for the marker effects such that $b_j \sim N(0, \sigma_b^2)$. This will help to
270 reduce the uncertainty of their estimation (*prediction*) by adding a bias. In both cases, the general
271 solution takes the following form

272

$$273 \hat{b} = (X'X + \lambda \times I_p)^{-1} X'y$$

274

275 The value to add in the diagonal matrix of $X'X$ is conveniently selected in a trade-off
276 between model goodness of fit and model complexity, where $\lambda \sim \frac{\sigma_\varepsilon^2}{\sigma_b^2}$, σ_b^2 and σ_ε^2 are the
277 corresponding variance components of the effect of the genomic covariates (genes/SNPs) and of
278 the error term.

279

280 Although the previous implementations allow us to get a solution when considering main
281 effects only, these still deal with large matrices and do not solve the problem of including
282 interactions between high dimensional sets of covariates (e.g., p genomic or phenomic and Q
283 environmental features for a total of $p \times Q$ first order contrasts). To tackle this problem, first we
284 examine an alternative parameterization proposed in animal breeding (VanRaden, 2008) to
285 include main effects in a computationally-convenient manner, then we provide a few details of
286 the implemented method for including interactions between groups of covariables.

287

288 The Genomic Best Linear Predictor (G-BLUP) attempts to directly compute the genomic
289 effect of the i^{th} individual g_i resulting from the linear combination between p marker and their
290 corresponding effects such that $g_i = \sum_{j=1}^p x_i b_{ij}$. Hence, instead of focusing in obtaining first the
291 marker effects b_{ij} to be used later in the linear combination, the genomic effect is obtained in
292 one step. The solution to this model requires the inversion of matrices of the type $(XX')^{-1}$ and
293 order $n \times n$ instead of $p \times p$, facilitating the handling of information derived from large matrices,
294 with X centered and scaled by columns (rows-genotypes; columns-marker SNPs). Under this

295 parameterization, the vector of genetic effects is modeled as $\mathbf{g} = \{g_i\} \sim N(\mathbf{0}, \mathbf{G}\sigma_g^2)$, where $\mathbf{G} =$
296 $\frac{XX'}{p}$ and $\sigma_g^2 = p \times \sigma_b^2$. Here, \mathbf{G} corresponds to the kinship matrix whose entries describe the
297 genomic similarities between pairs of individuals (VanRaden 2008). The resulting model in a
298 matrix parameterization is as follows

299

$$300 \quad \mathbf{y} = \mathbf{g} + \boldsymbol{\varepsilon} \quad [2]$$

301

302 A similar idea based on covariance structures can be considered to include high-
303 dimensional interactions between factors (more details are provided below in 2.3 section).
304 Jarquin et al. (2014) proposed the reaction norm model that allows the inclusion of all first order
305 interactions between genomic and weather factors. First, it was shown that the main effects of
306 weather covariables can be introduced into models in a similar fashion than the main effect of
307 genes or marker SNPs. Here, the environmental similarities between pairs of environments can
308 be characterized using weather information. This is analogous to considering marker SNPs to
309 conducting the genomic characterization between pairs of genotypes. Then the interactions
310 between markers and environmental factors are introduced via covariance structures computed as
311 the element-to-element product between the previous covariance matrices for genotypes and
312 environments.

313

314 Indirectly, this model, below described, allows to include the interaction between each
315 marker SNP and each weather covariate by modeling the interaction between their corresponding
316 linear combinations via covariances structures following the G-BLUP model fashion. The
317 resulting covariance structure of this interaction component that considers genomic and weather
318 factors is computed as the Hadamard product, which is the cell-by-cell product between two or
319 more covariance structures of the same dimension. In this case, the corresponding covariance
320 structures are redistributed/extended according to the vector of phenotypes and levels of the
321 corresponding factors (genotypes and environments) to ensure these are conformable.

322

323 In summary, modeling the G×E interactions can be both computationally and statistically
324 expensive due to the high dimensionality of the number of contrasts that can be formed between
325 genetic markers and environmental covariates (ECs). There are equivalent methods that reduce
326 such dimensionality by introducing markers and ECs via covariance structures as described in
327 Crossa et al. (2017). Interactions can be introduced through covariance structures computed via
328 the Hadamard product between these. Although these methods were already developed, there is
329 no simple method for capturing and integrating interactions among different omics.

330

2.3 Statistical Methods – Model Building

331

332 As mentioned previously, CHiDO was developed as a way to easily build models that can
333 capture main effects of diverse omics and incorporate the interactions between these, such as
334 those derived between genomic markers and ECs. CHiDO'S drag-and-drop interface simplifies
335 the process of creating complex models and adds a layer of abstraction for the methodology
336 established by Jarquin et al. (2014).

337

338 2.3.1 Main effects

339

340 Upon uploading the phenotypic response file, CHiDO automatically recognizes the
341 environment (E) and genotypic line (L) data to make them available as random effects, E_j and
342 L_i , respectively. These random effects can be added as terms in the model assembly section to
343 capture the inherent variability in phenotypic responses due to environmental and genetic
344 differences. Therefore, a base model with no additional omics data can be represented as

345

$$346 y_{ij} = \mu + E_j + L_i + \varepsilon_{ij} \quad [3]$$

347

348 where y_{ij} is the phenotypic observation (target trait) of the i^{th} genotype ($i = 1, \dots, L$) in the j^{th}
349 environment ($j = 1, \dots, J$), μ is the overall mean and ε_{ij} is the error term capturing the non-
350 explained variability by the other model terms.

351

352 When an omic data set \mathbf{O} is uploaded, CHiDO attempts to compute its specific vector of
353 effects $\mathbf{o} = \{o_k\}$, transforming the data into a covariance matrix $\mathbf{\Omega}$ that captures the similarities
354 among the pairs of entries for the different factor values (e.g., genotypes, environments,
355 genotypes-in-environments, etc.) For instance, if a file containing p ($m = 1, \dots, p$) genetic
356 markers $\mathbf{X} = \{x_{im}\}$ is uploaded, CHiDO attempts to modeling the vector of genomic effects $\mathbf{g} =$
357 $\{g_j\}$ as described for the G-BLUP model by constructing a genomic relationship matrix \mathbf{G} whose
358 entries describe genomic relationships between pairs of genotypes. For a given factor f (e.g.,
359 genotype, environment, genotype-environment combination) with T levels ($t = 1, \dots, T$), the
360 generalized form of the vector of effects associated to an omic-type \mathbf{o} can be calculated as a
361 linear combination between M covariates O_{tm} and their corresponding effects τ_m (e.g., SNP
362 markers, weather covariates, soil features, multispectral, Near InfraRed NIR, etc.)

363

$$364 o_t = \sum_{m=1}^M O_{tm} \tau_m$$

365

366 Using this form, we can describe general main effects for a given factor. For example,
367 modeling the genomic effect of the i^{th} ($i = 1, 2, \dots, L$) genotype using marker information $X =$
368 $\{x_{il}\}$ on p molecular markers and their corresponding effects b_l ($l = 1, 2, \dots, p$) we have

369

370
$$o_i = \sum_{l=1}^p x_{il} b_l$$

371
372 Similarly, for modeling the effect of the j^{th} environment based on Q weather covariates
373 $W = \{W_{jl}\}$ and their corresponding effects γ_l ($l = 1, 2, \dots, Q$) we have

374
$$o_j = \sum_{l=1}^Q w_{jl} \gamma_l$$

375
376 For an omic data observed at the particular/specific level *-compound-* (e.g. genotype-in-
377 environment combinations; the i^{th} genotype in the j^{th} environment), such as those derived from
378 high-throughput phenotyping platforms, the information $Z = \{z_{ijl}\}$ on s features (e.g., images)
379 can be modeled also as a linear combination considering their corresponding effects δ_l ($l = 1,$
380 $2, \dots, s$) as follows

381
$$o_{ij} = \sum_{l=1}^s z_{ijl} \delta_l$$

382
383 In addition, the information of covariance structures relevant to the factors of study (e.g.,
384 genotype, environment, etc.), can be also included in the models. For example, genetic effects
385 based on the pedigree matrix \mathbf{A} , or the environmental effects based on an environmental kinship
386 matrix \mathbf{C} . In these cases, it is necessary to specify the factor ID in the phenotypic matrix to
387 connect with the associated covariance structure. Hence, the alignment of the data will be
388 conducted as previously described, and also similar distributional assumptions (normality) as
389 before will be considered such that

390
391
$$\mathbf{o} = \{o_t\} \sim N(\mathbf{0}, \mathbf{\Omega} \sigma_{\mathbf{\Omega}}^2)$$

392
393 where $\mathbf{\Omega}$ is the corresponding covariance structure whose entries describe similarities between
394 pairs of levels (genotypes, environments, genotype-in-environments, etc.), and $\sigma_{\mathbf{\Omega}}^2$ is the
395 associated variance component. In this case, $\mathbf{\Omega}$ might represent the pedigree matrix (\mathbf{A}) whose
396 entries describe genetic similarities between pairs of individuals. Also, $\mathbf{\Omega}$ can represent an
397 environmental kinship matrix (\mathbf{C}) whose entries describe environmental similarities between
398 pairs or environments. If a covariance structure derived from soil information (\mathbf{S}) is available, it
399 can be also introduced into the models in a similar manner.

400
401 Models including only main effects can be easily constructed by adding the information
402 of the different omics into the linear predictor. For example, a linear model created using two
403 omics, one generic of type \mathbf{o} with T -levels ($t = 1, 2, \dots, T$) and M covariates, and another based
404 on p genetic markers for L individuals ($i = 1, \dots, L$) can be represented by

405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444

$$y_{it} = \mu + g_i + o_t + \varepsilon_{it} \quad [4]$$

CHiDO's capacity to handle different omic types of variable dimension extends to the creation and alignment of distinct covariance matrices for each associated dataset. This is done by calculating the matrix cross-product that reflects/expands the specific relationships across the levels of that omic type according to the matrix phenotypic responses. For this, on each of the different F omics \mathbf{O}_f ($f = 1, 2, \dots, F$) it is necessary to compute the incidence matrix Z_f that connects phenotypes with the T different levels of the omics (e.g., genotype, family, environment, mega-environment, farm, herd, genotype-in-environment combination, etc.) Then, the resulting aforementioned covariance matrices are aligned and/or expanded across all phenotypic records by computing $Z_f \mathbf{O}_f Z_f'$ with the `tcrossprod(tcrossprod(Z_f, O), Z_f)` instruction.

2.3.2 Multiplicative interactions

Interactions between different omics are modeled by calculating the Hadamard product of their corresponding covariance structures. For example, the interaction between the covariance matrices \mathbf{G} and \mathbf{O} denoted by $(\mathbf{G} \# \mathbf{O})$, represents the interaction between genotypes (using molecular marker information \mathbf{X}) and any other related omic - \mathbf{O} -. The corresponding covariance matrix of this interaction term is represented with $\mathbf{B}_{\mathbf{G} \# \mathbf{O}} = Z_g \mathbf{G} Z_g' \circ Z_o \mathbf{O} Z_o'$, and modeled as

$$\mathbf{g} \times \mathbf{o} \sim N(\mathbf{0}, \mathbf{B}_{\mathbf{G} \# \mathbf{O}} \sigma_{\mathbf{G} \# \mathbf{O}}^2)$$

where Z_g and Z_o are the corresponding incidence matrices that connect the phenotypic observations with the different levels of the omics (e.g., genotypes, environments, genotype-in-environment, families, etc.)

For any given covariance matrix of the main and interaction effects, CHiDO performs the spectral decomposition using the `eigen()` function to retrieve its *eigenvalues* and *eigenvectors*. The *eigenvalues* reveal the magnitude of variance in the omic data along the directions defined by their corresponding eigenvectors. For G×E predictions, this information could provide insights into the major factors that contribute towards trait variation. This factorization is conveniently implemented to save computing time when fitting different linear predictors and prediction scenarios in BGLR R-package. Each time the BGLR function is used, and the covariance matrices are provided, it internally computes the eigen-value decomposition before starting the model fitting. Using datasets with a large number of phenotypic observations (n) this procedure might be time consuming, especially in those cases where the cross-validations involve exhaustive scenarios and/or folds. Thus, by providing the resulting factorization of these matrices a considerable amount of time and resources are saved avoiding extra computational burden.

2.3.3 Cross-validation schemes

445

446 Prior to implementing prediction models in real-world applications such as GS, it is necessary to
447 evaluate their usefulness integrating different omics to deliver accurate and reliable results.

448 Cross-validation studies are a common, time-tested method to perform such evaluation. Hence,
449 after the models are created and saved in CHiDO, users can select from a range of cross-
450 validation (CV) schemes (based on their specific research objectives) how to train and evaluate
451 the performance of their model(s).

452

453 These CV schemes mimic real life prediction problems that breeders face at different stages
454 along the breeding pipeline for the development of improved genotypes. As discussed, CV1
455 considers the prediction of ‘newly’ or untested genotypes in environments where other genotypes
456 were already observed. CV2 (or incomplete field trials) mimics the prediction of already tested
457 genotypes observed in other environments but not in the target environment (where other
458 genotypes were also already tested). CV0 (or forward prediction) emulates the prediction of
459 already tested (in other environments) genotypes in novel environments where no phenotypic
460 records on any of the lines have been collected. CV00 is similar to the previous scheme with the
461 main difference that the genotypes to predict have not been observed at any of the environments
462 in the training sets. This last prediction scenario is the most challenging and probably the most
463 interesting for breeders.

464

465 The manner to create the different partitions representing training and testing sets depends on the
466 prediction problem (cross-validation scheme). Here, the folds are defined by the user according
467 to the selected CV scheme to partition the phenotypic data (training/testing). For instance, in a k -
468 fold cross-validation setting such as in CV1 and CV2, the dataset D is divided into k mutually
469 exclusive subsets (D_1, D_2, \dots, D_k) , with each subset serving as a testing set -one at a time- while
470 the remaining subsets are aggregated to form the training set. Under CV2 scheme, the
471 phenotypes are randomly assigned to the folds, while under the CV1 scheme extra care is taken
472 to assign genotypes to folds ensuring that all the phenotypic records from the same individual
473 appear in the same fold. On the other hand, under CV0 and CV00 each environment naturally
474 becomes a fold and care is taken to ensure similar training sample sizes to those in the previous
475 schemes (CV2 and CV1) according to Persa et al., (2021). When performing the different CV
476 schemes, CHiDO loops the folds until all folds are considered as testing or prediction sets using
477 the BGLR function.

478

479 Since the models are fitted under the Bayesian framework, the users can define additional
480 training hyper-parameters for BGLR such as the number of iterations and the burn-in rate. These
481 parameters influence the convergence and stability of the Bayesian models. As mentioned above,
482 the cross-validations are executed using the `BGLR()` function, which applies the user-defined
483 settings. The *eigenvalues* and *eigenvectors* for each omic-matrix, carrying the information of the

484 different model terms, are incorporated into the ETA object to compose the readable linear
485 predictor for BGLR.

486

487

488 **2.3.4 Metrics**

489

490 Upon completion of the BGLR analysis, CHiDO employs the model outputs to calculate several
491 metrics essential for evaluating the performance of the different linear models. Custom functions
492 have been developed within the CHiDO framework to facilitate these calculations, ensuring
493 accuracy and efficiency in metric derivation.

494

495 *Prediction accuracy (PA) measured on a trial basis:* It is obtained by computing the Pearson's
496 moment correlation ρ between predicted and observed (phenotypic) values within each
497 trial/environment/year/location/etc. This metric helps to determine how well a given model can
498 predict phenotypic traits based on the multi-omics data associated to the provided model terms.
499 The formula for PA in the j^{th} environment (or grouping factor) is given by

500

$$501 \quad \rho_j = \frac{\sum_{i=1}^{n_j} (\hat{y}_{ij} - \hat{\bar{y}}_j) (y_{ij} - \bar{y}_j)}{\sqrt{\sum_{i=1}^{n_j} (\hat{y}_{ij} - \hat{\bar{y}}_j)^2} \sqrt{\sum_{i=1}^{n_j} (y_{ij} - \bar{y}_j)^2}}$$

502

503 where \hat{y}_{ij} and y_{ij} are the predicted and the observed values of the i^{th} genotype at the j^{th}
504 environment, $\hat{\bar{y}}$ and \bar{y}_j are their corresponding means, and n_j represents the number of
505 observations at the j^{th} environment.

506

507 For an easier assessment of the model's performance across environments, the weighted mean
508 correlation is computed accounting for the uncertainty and the sample size of the environments
509 according to Tiezzi et al. (2017) as follows:

$$510 \quad \rho_\varphi = \frac{\sum_{j=1}^J \frac{\rho_j}{V(\sigma_j)}}{\sum_{j=1}^J \frac{1}{V(\sigma_j)}}$$

511 where $V(\sigma_j) = \frac{1-\rho_j^2}{n_j}$ corresponds to the sampling variance.

512

513 *Root Mean Squared Error (RMSE):* Quantifies the average magnitude of prediction error,
514 measures a model's precision, and penalizes large errors to a greater extent by squaring the
515 difference between predicted and observed values. The formula for RMSE for the j^{th}
516 environment is given by:

517

518

$$RMSE_j = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_{ij} - y_{ij})^2}$$

519

520

521

522 *Variance Components*: This metric measures the portion of variance explained by each model
523 term associated to an omic with respect to the overall phenotypic variability. This estimation is
524 critical for understanding which main or interaction effects influence the most the phenotypic
525 expression/variability of target traits. It is computed considering a full data analysis (i.e., no
526 missing values are generated on the phenotypic information).

527

528 The variance component of each term is computed as the percentage of the total variance
529 explained, which for the f^{th} ($f=1, 2, \dots, F$) omic \mathbf{O}_f it corresponds to the ratio between the current
530 variance component and the sum of all the F variance components plus the unexplained residual
531 variance σ_ε^2

532

533

$$\% \tilde{\sigma}_{\mathbf{O}_f}^2 = \frac{\text{Specific Variance}_i}{\text{Total Variance}} \times 100 = \frac{\tilde{\sigma}_{\mathbf{O}_f}^2}{\sum_{f=1}^F \tilde{\sigma}_{\mathbf{O}_f}^2 + \tilde{\sigma}_\varepsilon^2} \times 100$$

534

535 Here, under a given model, the specific variance refers to the variance attributable to the
536 particular model term f^{th} and total variance is the sum of variances of all terms, including the
537 residual variance. The total variance corresponds to the 100% of the phenotypic variability.

538

539

3 RESULTS AND DISCUSSION

540

541 Dealing with prediction analyses for breeding applications, usually an important amount of time
542 (~85%) is dedicated to the data preparation (quality control, alignment, cross-validation
543 scenarios, etc.) and the remaining time (~15%) is for the development and implementation of
544 these models. Therefore, the availability of low-code, no-code (LCNC) applications such as
545 CHiDO can help breeders save time and obtain expedited results by automating and assisting
546 with many of these tasks, allowing them to focus on specific research questions derived from
547 initial quick analyses.

548

549 In this paper we discussed the reason for CHiDO's development, the technical and statistical
550 methods applied, and its potential benefits to breeders. The CHiDO platform is a significant
551 contribution to empower breeders and democratize access to modern solutions by enabling the

552 modeling of different interaction types such as the G×E interaction without the need for in-depth
553 programming knowledge. Increasing access to advanced analytics and prediction tools can not
554 only accelerate research for new improved varieties/individuals, but also enable broader
555 participation in agricultural research. CHiDO reflects a growing trend towards more accessible
556 and flexible computational tools in genomics, as evidenced by recent literature advocating for the
557 democratization of data science (Shang et al., 2019).

558
559 The practical implications of the CHiDO platform extend significantly beyond the immediate
560 sphere of plant breeding. By enabling more accurate and efficient selection processes, CHiDO
561 contributes to the development of crops with improved yields and environmental resilience. This
562 capacity is particularly crucial in the context of climate change and the increasing demands for
563 sustainable agricultural practices. The forthcoming introduction of interactive graphics for model
564 evaluation further underscores CHiDO's potential to enhance understanding and application of
565 complex genomic data in breeding strategies.

566
567 LCNC platforms such as CHiDO are becoming increasingly popular and offer various benefits
568 for researchers (Sufi, 2023). Some benefits include 1) ease of adoption through a reduced
569 learning curve, 2) accelerated development speed, and 3) circumventing resource scarcity,
570 among many others listed in (Sufi, 2023; Yan, 2021). Despite these benefits, LCNC solutions are
571 not without their challenges. Some notable drawbacks to LCNC are recurring costs and vendor
572 lock-in. Similarly, developers can learn how to use the platform effectively but are bound to the
573 limitations of said platform without the potential to extend its functionalities as opposed to
574 custom developed alternatives.

575
576 We are addressing these drawbacks in CHiDO by ensuring the platform remains a free-to-use
577 service and providing users the ability to submit issues or product feature requests on GitHub
578 (<https://github.com/jarquinalab/CHiDO>). In addition to this, we are evaluating the potential
579 release of CHiDO's backend logic as an R package or API for more advanced users to extend
580 CHiDO's functionalities or integrate the tools with other packages when scripting.

581
582 In addition to the aforementioned features, future updates to CHiDO aim to enhance its
583 functionality to cover a broader array of plant and animal breeding prediction scenarios, with the
584 potential to extend these to public health applications such as personalized medicine. However,
585 working on these proposed developments below detailed while maintaining an optimum
586 functionality of the software will require of the investment of resources. We will seek for
587 funding opportunities and partnerships to secure the needed resources to continue these and other
588 future developments.

589
590 A few of the key additions we would like to integrate are modules for sparse testing designs,
591 estimation of G×E markers using weather data -enabling a focused analysis on the relevance of

592 genetic markers and ECs influencing target traits-, hybrid prediction via general and specific
593 combining ability (GCA, SCA) terms and their corresponding interactions. Separately, CHiDO
594 will incorporate options for selecting from multiple artificial intelligence (AI/ML) algorithms to
595 facilitate the modeling of complex, non-linear relationships within multi-omics datasets. The use
596 of Deep Learning and ML algorithms (e.g., RandomForest) is already being evaluated for their
597 robustness in capturing intricate G×E interactions (Crossa et al., 2019), potentially leading to
598 more accurate genomic selections. The launch of CHiDO online, alongside comprehensive
599 documentation, is poised to democratize access to these advanced tools, stimulating worldwide
600 collaboration and further research.

601
602 Ultimately, CHiDO stands at the forefront of integrating multi-omics data for plant breeding,
603 representing a critical advancement in computational tools within agriculture. Its development is
604 timely, addressing the urgent need for innovative solutions in plant breeding to meet the global
605 challenges of food security and sustainability.

606
607

608 ACKNOWLEDGMENTS

609
610 JGA would like thank to the ‘Programa Propio’ of the ‘Universidad Politécnica de Madrid’
611 (UPM) for the financial support during his PhD studies and internship at the University of
612 Florida. All authors would like to thank the UF Strawberry Breeding Program for their
613 contribution in the technical discussion during the design phase.

614

615 AUTHOR CONTRIBUTIONS

616 **FG** Methodology, Software, Validation, Formal analysis, Writing – Original draft, Writing –
617 Reviewing and Editing; **JGA** Software, Visualization, Writing – Reviewing and Editing; **DJ**
618 Conceptualization, Resources, Supervision, Writing – Reviewing and Editing

619

620 DATA AVAILABILITY

621 There are no original data associated with this article. CHiDO is a web-based application
622 accessible at <https://jarquinlab.shinyapps.io/chido/> where users can upload their own data to
623 develop predictive models. Data uploaded to CHiDO is not stored anywhere and is only used
624 during the active session while users interact with the platform.

625 For demos and testing purposes, users can use sample data available at
626 <https://github.com/jarquinlab/CHiDO>. These data sets were extracted from Trachsel et al. (2019)
627 and correspond to a maize experiment comprising 97 genotypes (double haploid) tested in four
628 environments (two reps, and only rep was used for the demo) and scored for grain yield (GY),

629 plant height (PH), anthesis silk interval (ASI), and day to anthesis (DA). Also, genomic (551
630 marker SNPs) and hyperspectral (five flights - 62 bands per-fly) data were available for analysis.
631 In addition, a kinship matrix was computed using a random sample of SNPs to emulate a
632 pedigree matrix.

633

634

CONFLICT OF INTEREST

635 The authors express no conflict of interest with any of the components involved in this
636 publication.

637

638

REFERENCES

- 639 Atlin, G.N., Cairns, J.E., Das, B., 2017. Rapid breeding and varietal replacement are critical to
640 adaptation of cropping systems in the developing world to climate change. *Glob. Food*
641 *Secur.* 12, 31–37. <https://doi.org/10.1016/j.gfs.2017.01.008>
- 642 Bernardo, R. 1994. Prediction of maize single-cross performance using RFLPs and information
643 from related hybrids. *Crop Science.* 34: 20-25.
- 644 Clapp, J., 2018. Mega-Mergers on the Menu: Corporate Concentration and the Politics of
645 Sustainability in the Global Food System. *Glob. Environ. Polit.* 18, 12–33.
646 https://doi.org/10.1162/glep_a_00454
- 647 Crossa, J., Martini, J.W.R., Gianola, D., Pérez-Rodríguez, P., Jarquin, D., Juliana, P.,
648 Montesinos-López, O., Cuevas, J., 2019. Deep Kernel and Deep Learning for Genome-
649 Based Prediction of Single Traits in Multienvironment Breeding Trials. *Front. Genet.* 10,
650 1168. <https://doi.org/10.3389/fgene.2019.01168>
- 651 Crossa, J., Pérez-Rodríguez, P., Cuevas, J., Montesinos-López, O., Jarquín, D., De Los Campos,
652 G., Burgueño, J., González-Camacho, J.M., Pérez-Elizalde, S., Beyene, Y., Dreisigacker,
653 S., Singh, R., Zhang, X., Gowda, M., Roorkiwal, M., Rutkoski, J., Varshney, R.K., 2017.
654 Genomic Selection in Plant Breeding: Methods, Models, and Perspectives. *Trends Plant*
655 *Sci.* 22, 961–975. <https://doi.org/10.1016/j.tplants.2017.08.011>
- 656 Dean Attali, 2021. shinyjs: Easily Improve the User Experience of Your Shiny Apps in Seconds.
- 657 Foley, J.A., Ramankutty, N., Brauman, K.A., Cassidy, E.S., Gerber, J.S., Johnston, M., Mueller,
658 N.D., O’Connell, C., Ray, D.K., West, P.C., Balzer, C., Bennett, E.M., Carpenter, S.R.,
659 Hill, J., Monfreda, C., Polasky, S., Rockström, J., Sheehan, J., Siebert, S., Tilman, D.,
660 Zaks, D.P.M., 2011. Solutions for a cultivated planet. *Nature* 478, 337–342.
661 <https://doi.org/10.1038/nature10452>
- 662 Gu, D., Andreev, K., Dupre, M.E., 2021. Major Trends in Population Growth Around the World.
663 *China CDC Wkly.* 3, 604–613. <https://doi.org/10.46234/ccdcw2021.160>
- 664 Hadley Wickham, 2016. ggplot2: Elegant Graphics for Data Analysis.
- 665 Jarquín, D., Crossa, J., Lacaze, X., Du Cheyron, P., Daucourt, J., Lorgeou, J., Piraux, F.,
666 Guerreiro, L., Pérez, P., Calus, M., Burgueño, J., De Los Campos, G., 2014. A reaction
667 norm model for genomic selection using high-dimensional genomic and environmental
668 data. *Theor. Appl. Genet.* 127, 595–607. <https://doi.org/10.1007/s00122-013-2243-1>
- 669 Jarquín, D., Lemes Da Silva, C., Gaynor, R.C., Poland, J., Fritz, A., Howard, R., Battenfield, S.,

- 670 Crossa, J., 2017. Increasing Genomic-Enabled Prediction Accuracy by Modeling
671 Genotype \times Environment Interactions in Kansas Wheat. *Plant Genome* 10,
672 *plantgenome2016.12.0130*. <https://doi.org/10.3835/plantgenome2016.12.0130>
- 673 Jarquin, D., Howard, R., Crossa, J., Beyene, Y., Gowda, M., Martini, J.W.R., Covarrubias
674 Pazaran, G., Burgueño, J., Pacheco, A., Grondona, M., Wimmer, V., Prasanna, B.M.,
675 2020. Genomic Prediction Enhanced Sparse Testing for Multi-environment Trials. *G3*
676 *GenesGenomesGenetics* 10, 2725–2739. <https://doi.org/10.1534/g3.120.401349>
- 677 Lesk, C., Rowhani, P., Ramankutty, N., 2016. Influence of extreme weather disasters on global
678 crop production. *Nature* 529, 84–87. <https://doi.org/10.1038/nature16467>
- 679 Meuwissen, T.H.E., Hayes, B.J., Goddard, M.E., 2001. Prediction of Total Genetic Value Using
680 Genome-Wide Dense Marker Maps. *Genetics* 157, 1819–1829.
681 <https://doi.org/10.1093/genetics/157.4.1819>
- 682 Pérez, P., De Los Campos, G., 2014. Genome-Wide Regression and Prediction with the BGLR
683 Statistical Package. *Genetics* 198, 483–495. <https://doi.org/10.1534/genetics.114.164442>
- 684 Persa, R., Iwata, H., and Jarquin, D., 2020. Use of family structure information in interaction
685 with environments for leveraging genomic prediction models. *Crop J.* 8, 843–854.
686 [doi:10.1016/j.cj.2020.06.004](https://doi.org/10.1016/j.cj.2020.06.004)
- 687 Shang, Z., Zraggen, E., Buratti, B., Kossmann, F., Eichmann, P., Chung, Y., Binnig, C., Upfal,
688 E., Kraska, T., 2019. Democratizing Data Science through Interactive Curation of ML
689 Pipelines, in: *Proceedings of the 2019 International Conference on Management of Data*.
690 Presented at the SIGMOD/PODS '19: International Conference on Management of Data,
691 ACM, Amsterdam Netherlands, pp. 1171–1188.
692 <https://doi.org/10.1145/3299869.3319863>
- 693 Sufi, F., 2023. Algorithms in Low-Code-No-Code for Research Applications: A Practical
694 Review. *Algorithms* 16, 108. <https://doi.org/10.3390/a16020108>
- 695 Tiezzi, F., de Los Campos, G., Gaddis, K.P., Maltecca, C., 2017. Genotype by environment
696 (climate) interaction improves genomic prediction for production traits in us holstein
697 cattle. *J. Dairy Sci.* 2017, 100, 2042–2056
- 698 Trachsel, S., Dhliwayo, T., Perez, L.G., Lugo, J.A.M., and Trachsel, M., 2019. Estimation of
699 physiological genomic estimated breeding values (PGEbV) combining full hyperspectral
700 and marker data across environments for grain yield under combined heat and drought
701 stress
702 in tropical maize (*Zea mays* L.). *PLoS One* 14:e0212200
- 703 Tsatsakis, A.M., Nawaz, M.A., Tutelyan, V.A., Golokhvast, K.S., Kalantzi, O.-I., Chung, D.H.,
704 Kang, S.J., Coleman, M.D., Tyshko, N., Yang, S.H., Chung, G., 2017. Impact on
705 environment, ecosystem, diversity and health from culturing and using GMOs as feed and
706 food. *Food Chem. Toxicol.* 107, 108–121. <https://doi.org/10.1016/j.fct.2017.06.033>
- 707 Van Dijk, M., Morley, T., Rau, M.L., Saghai, Y., 2021. A meta-analysis of projected global food
708 demand and population at risk of hunger for the period 2010–2050. *Nat. Food* 2, 494–
709 501. <https://doi.org/10.1038/s43016-021-00322-9>
- 710 VanRaden P.M., 2008. Efficient methods to compute genomic predictions. *J Dairy Sci* 91:4414–
711 4423. [doi:10.3168/jds.2007-0980](https://doi.org/10.3168/jds.2007-0980)
- 712 Varshney, R.K., Terauchi, R., McCouch, S.R., 2014. Harvesting the Promising Fruits of
713 Genomics: Applying Genome Sequencing Technologies to Crop Breeding. *PLoS Biol.*
714 12, e1001883. <https://doi.org/10.1371/journal.pbio.1001883>
- 715 Winston Chang, Joe Cheng, JJ Allaire, Carson Sievert, Barret Schloerke, Yihui Xie, Jeff Allen,

716 Jonathan McPherson, Alan Dipert, Barbara Borges, 2023. shiny: Web Application
717 Framework for R.
718 Yan, Z., 2021. The Impacts of Low/No-Code Development on Digital Transformation and
719 Software Development.
720 Yang, Y., Saand, M.A., Huang, L., Abdelaal, W.B., Zhang, J., Wu, Y., Li, J., Sirohi, M.H.,
721 Wang, F., 2021. Applications of Multi-Omics Technologies for Crop Improvement.
722 *Front. Plant Sci.* 12, 563953. <https://doi.org/10.3389/fpls.2021.563953>
723 Yang Tang, 2022. shinyjq: “jQuery UI” Interactions and Effects for Shiny.
724 Yihui Xie, Joe Cheng, Xianying Tan, 2023. DT: A Wrapper of the JavaScript Library
725 “DataTables.”
726 Zhao, C., Liu, B., Piao, S., Wang, X., Lobell, D.B., Huang, Y., Huang, M., Yao, Y., Bassu, S.,
727 Ciais, P., Durand, J.-L., Elliott, J., Ewert, F., Janssens, I.A., Li, T., Lin, E., Liu, Q.,
728 Martre, P., Müller, C., Peng, S., Peñuelas, J., Ruane, A.C., Wallach, D., Wang, T., Wu,
729 D., Liu, Z., Zhu, Y., Zhu, Z., Asseng, S., 2017. Temperature increase reduces global
730 yields of major crops in four independent estimates. *Proc. Natl. Acad. Sci.* 114, 9326–
731 9331. <https://doi.org/10.1073/pnas.1701762114>

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

FIGURES AND TABLES

758

759 Table 1. List of all R packages used in the creation of CHiDO.

760

| Package | Function in CHiDO | URL |
|----------------|--------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------|
| shiny | Main framework for displaying and organizing the web application | https://CRAN.R-project.org/package=shiny |
| shinydashboard | Simplified the creation of dashboards within the Shiny framework | https://CRAN.R-project.org/package=shinydashboard |
| gridExtra | Align widgets, plots, and data in grid-like format | https://CRAN.R-project.org/package=gridExtra |
| dplyr | Perform data processing and transformations in a consistent manner | https://CRAN.R-project.org/package=dplyr |
| DT | Handle and render tabular objects using R and/or JavaScript syntax | https://CRAN.R-project.org/package=DT |
| ggplot2 | Generate graphics of cross-validation results and evaluation metrics | https://ggplot2.tidyverse.org |
| shinyjs | Integrating JavaScript into Shiny application to extend functionalities of UI | https://CRAN.R-project.org/package=shinyjs |
| shinyjqui | Enable animation effects needed for the drag-and-drop interface in the model assembly page | https://CRAN.R-project.org/package=shinyjqui |

761

762

763

764

765

766

767

768

769

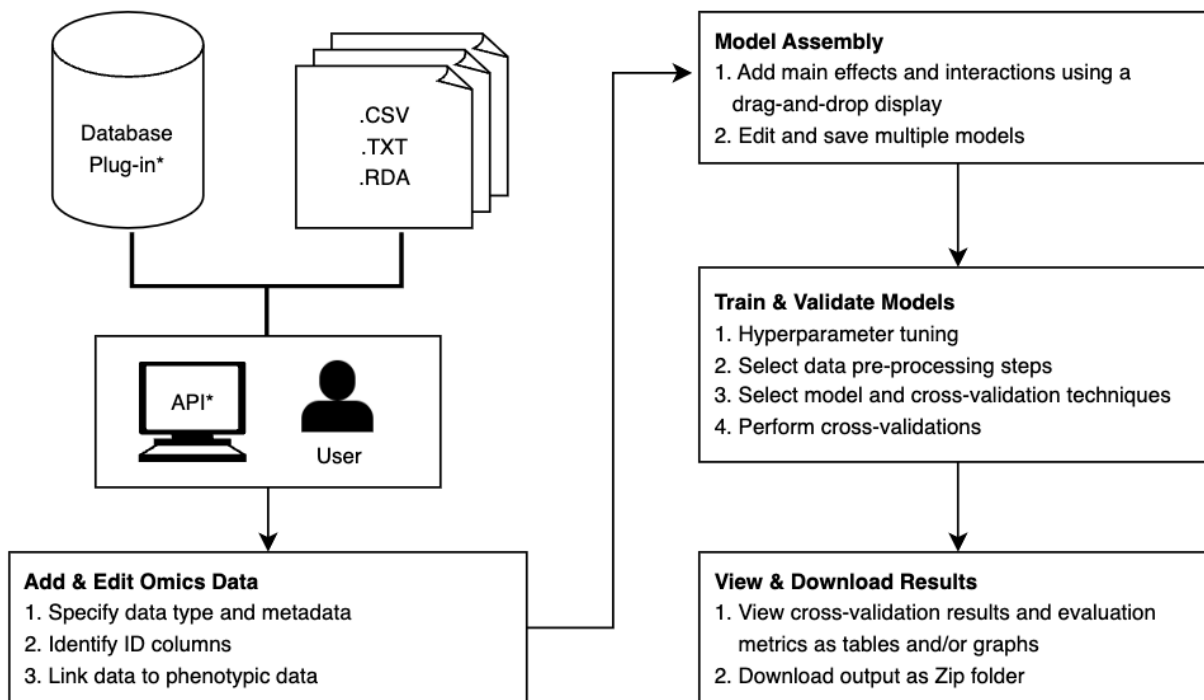
770

771

772

773

774

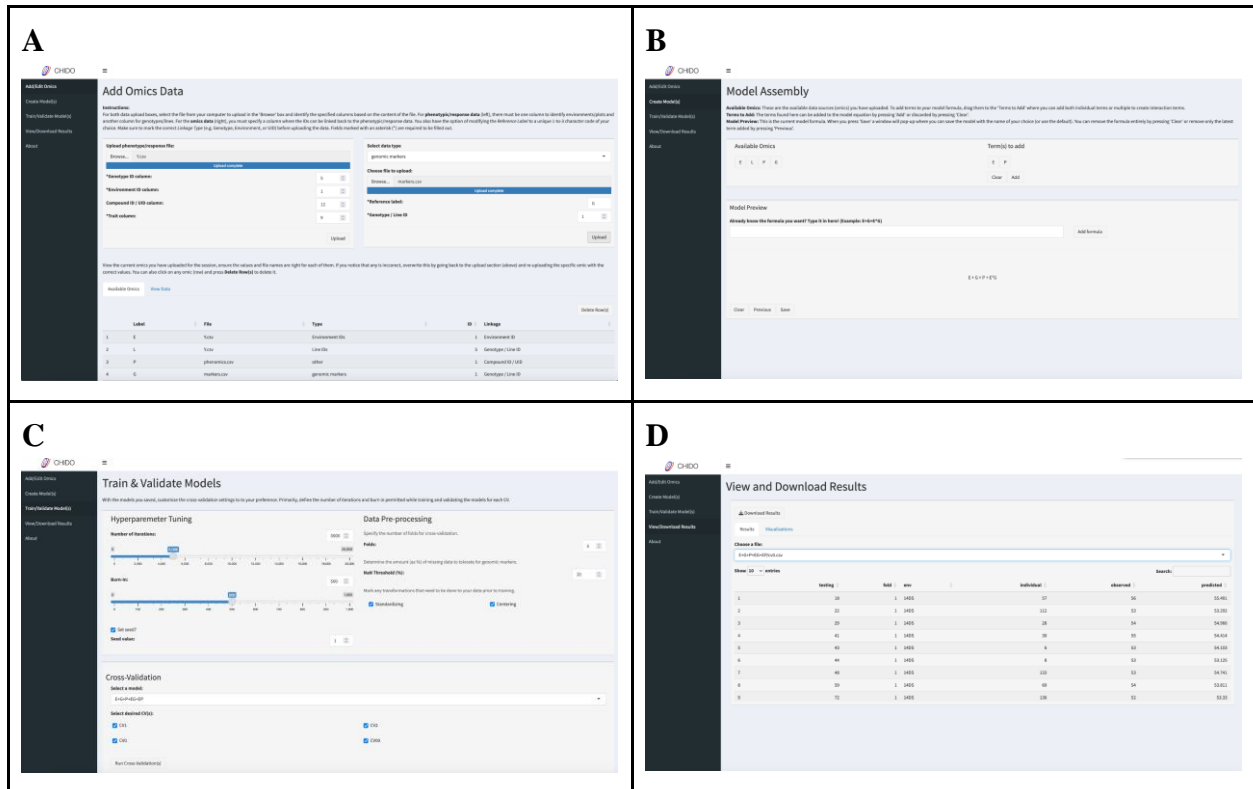


(*) denotes this development is still pending

Figure 1. Overview of the different components and functionalities within the CHiDo platform

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 Figure 2. User interface for CHiDO; (A) The *Add Omics Data* page is where users upload their
800 files and define metadata for each of them such that the platform can treat them as separate
801 omics; (B) The *Model Assembly* page lets users create multiple models using the uploaded data
802 as main effects or combining them with interaction terms; (C) Users can tune training and
803 validation parameters, apply quality control on the genomic data, as well as selecting the
804 different cross-validation schemes to employ; and (D) the *View and Download Results* page
805 allows users to view prediction outputs and evaluation metrics in tabular and graphical formats
806 before downloading them to the user's local environment.

807

808

809

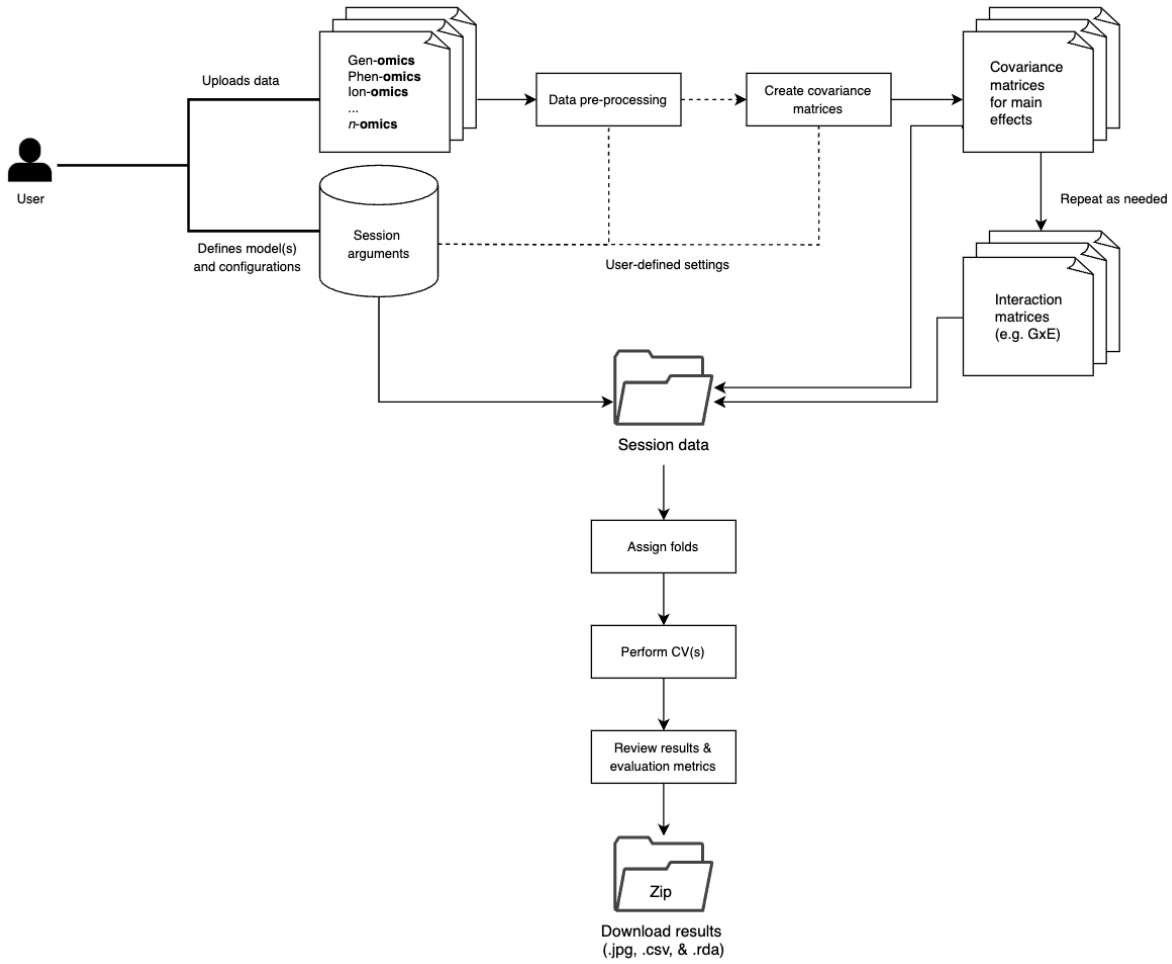
810

811

812

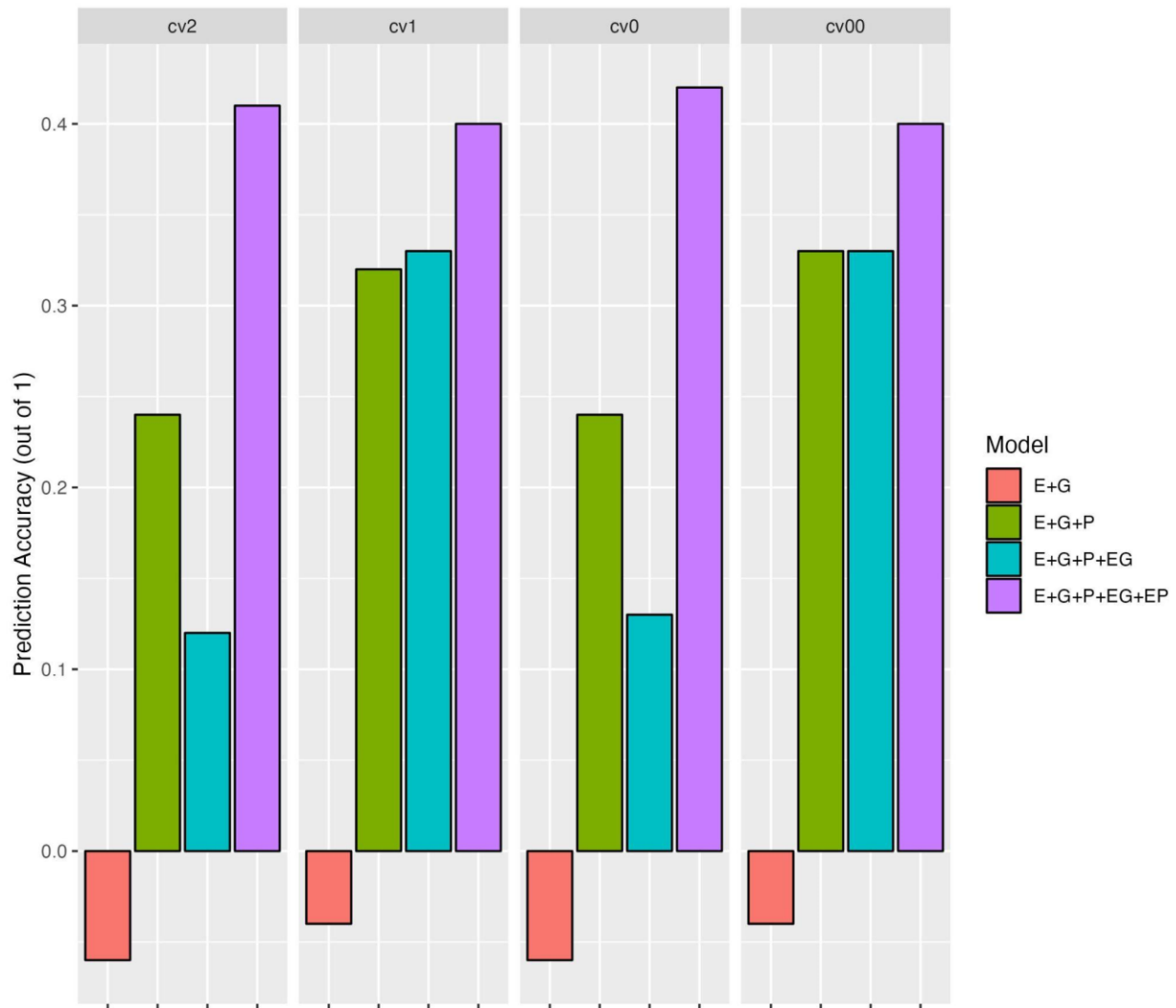
813

814



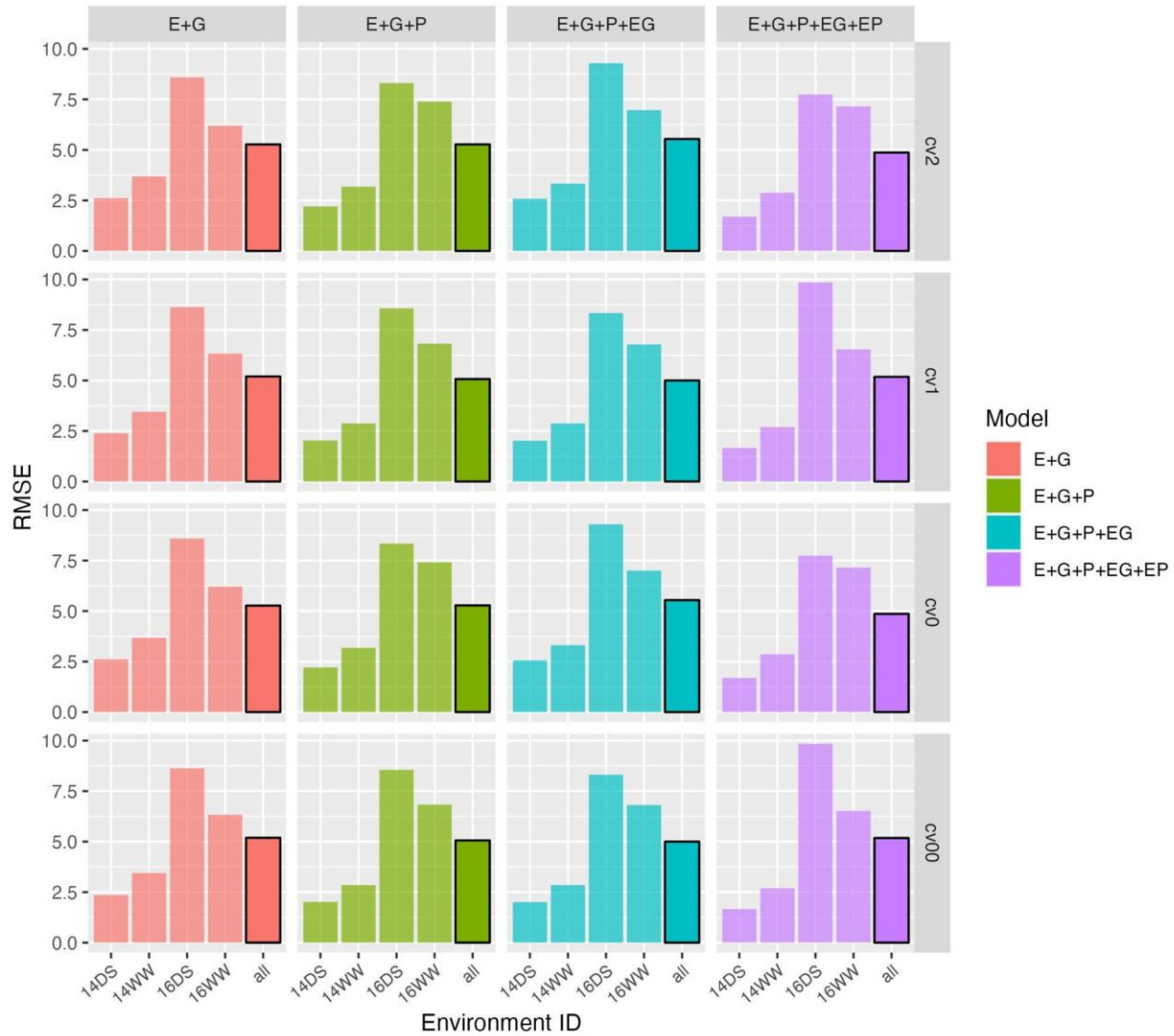
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832

Figure 3. Workflow diagram for CHiDO. This diagram demonstrates the logic implemented in the application to create and train linear models using arguments and data provided by the user.



833
834
835
836
837
838
839
840
841
842
843
844
845

Figure 4. Example of prediction accuracy results by model, and by cross-validation scheme.



846
847
848
849

Figure 5. Example of the model root-mean-square error (RMSE) by environment, and by cross-validation scheme.