

Nicolas Morales <nm529@cornell.edu> * †¹
Alex C. Ogonna <aco46@cornell.edu> * †¹
Bryan J. Ellerbrock <bjie24@cornell.edu> *¹
Guillaume J. Bauchet <gjb99@cornell.edu> *^{1 2}
Titima Tantikanjana <tt15@cornell.edu> *
Isaak Y. Teclé <isaak.yosief@cornell.edu> *
Adrian F. Powell <afp43@cornell.edu> *
David Lyon <dal333@cornell.edu> *³
Naama Menda <nm249@cornell.edu> *
Christiano C. Simoes <ccs263@cornell.edu> *
Surya Saha <ss2489@cornell.edu> *
Prashant Hosmani <psh65@cornell.edu> *⁴
Mirella Flores <mrf252@cornell.edu> *
Naftali Panitz <np298@cornell.edu> *
Ryan S. Preble <rsp98@cornell.edu> *
Afolabi Agbona <A.Afolabi@cgiar.org> ‡
Ismail Rabbi <i.rabbi@cgiar.org> ‡
Peter Kulakow <p.kulakow@cgiar.org> ‡
Prasad Peteti <p.prasad@cgiar.org> ‡
Robert Kawuki <kawukisezirobert@gmail.com> §
Williams Esuma <esumawilliams@yahoo.com> §
Micheal Kanaabi <kanaabimichael@gmail.com> §
Doreen M. Chelangat <dmurenju@gmail.com> §

- 26 Ezenwanyi Uba <ubaezenwanyi@yahoo.com> **
- 27 Adeyemi Olojede <yemiolojede@gmail.com> **
- 28 Joseph Onyeka <jonyeka@yahoo.com> **
- 29 Trushar Shah <tm.shah@cgiar.org> ††
- 30 Margaret Karanja <M.Karanja@cgiar.org> ††
- 31 Chiedozie Egesi <C.Egesi@cgiar.org> † ‡ **
- 32 Hale Tufan <hat36@cornell.edu> †
- 33 Agre Paterne <P.Agre@cgiar.org> ‡
- 34 Asrat Asfaw <A.Amele@cgiar.org> ††
- 35 Jean-Luc Jannink <jeanluc.jannink@usda.gov> §§ †
- 36 Marnin Wolfe <mw489@cornell.edu> †
- 37 Clay L. Birkett <clb343@cornell.edu> §§ †
- 38 David J. Waring <djw64@cornell.edu> §§ †
- 39 Jenna M. Hershberger <jmh579@cornell.edu> †
- 40 Michael A. Gore <mag87@cornell.edu> †
- 41 Kelly R. Robbins <krr73@cornell.edu> †
- 42 Trevor Rife <trife@ksu.edu> ***
- 43 Chaney Courtney <chaneylc@gmail.com> ***
- 44 Jesse Poland <jpoland@ksu.edu> ***
- 45 Elizabeth Arnaud <e.arnaud@cgiar.org> †††
- 46 Marie-Angélique Laporte <m.a.laporte@cgiar.org> †††
- 47 Heneriko Kulembeka <kulembeka@yahoo.com> §§§
- 48 Kasele Salum <kaselesalum@gmail.com> §§§
- 49 Emmanuel Mrema <emmanuel.oroaya@gmail.com> §§§
- 50 Allan Brown <a.brown@cgiar.org> ‡
- 51 Stanley Bayo <s.bayo@cgiar.org> ‡

- 52 Brigitte Uwimana <b.uwimana@cgiar.org> ‡
- 53 Violet Akech <v.akesh@cgiar.org> ‡
- 54 Craig Yencho <Craig_Yencho@ncsu.edu> ††††
- 55 Bert de Boeck <b.deboeck@cgiar.org> ****
- 56 Hugo Campos <h.campos@cgiar.org> ****
- 57 Rony Swennen <r.swennen@cgiar.org> ††††
- 58 Jeremy D. Edwards <jeremy.edwards@usda.gov> †††
- 59 Lukas A. Mueller <lam87@cornell.edu> *
- 60
- 61
- 62
- 63 * Boyce Thompson Institute, Ithaca NY 14853, USA
- 64 † Cornell University, Ithaca, NY 14853, USA
- 65 ‡ IITA Ibadan, 200001 Ibadan, Nigeria
- 66 § NaCCRI, Namulonge, Uganda
- 67 †† IITA Abuja, 901101 Abuja, Nigeria
- 68 †† IITA Nairobi, 30709-00100 Nairobi, Kenya
- 69 ** National Root Crops Research Institute (NRCRI), 463109 Umudike, Nigeria
- 70 §§ USDA-ARS, Ithaca NY 14853, USA
- 71 *** Kansas State University; Manhattan, KS 66506, USA
- 72 ††† USDA-ARS, Stuttgart AR 72160, USA
- 73 ††† Bioversity-CIAT Alliance, 34397 Montpellier, France
- 74 §§§ TARI, 33518 Ukiriguru, Tanzania
- 75 **** CIP, 15000 Lima, Peru
- 76 †††† KU Leuven, 3000 Leuven, Belgium
- 77 †††† North Carolina State University (NCSSU), Raleigh NC 27695, USA

78

79 ¹ These authors contributed equally to this work

80 ² Present address: Terre de Lin, Saint-Pierre-le-Viger, France

81 ³ Present address: Lawrence Berkeley National Laboratory, Berkeley, CA, USA

82 ⁴ Present address: NRGene, Morristown, NC, USA

83

84 Breedbase, a digital ecosystem for breeding

85

86 Keywords:

87 Database, breeding, phenotyping, genotyping, predictive breeding, genomic selection,

88 genome-based breeding, digital ecosystem, digital agriculture, web-based software,

89 open source breeding software

90

91 Corresponding author:

92 Lukas A. Mueller

93 Boyce Thompson Institute

94 533, Tower Road

95 Ithaca NY 14853

96 (607) 255 6557

97 LAM87@cornell.edu

98

99

100

101

102 ABSTRACT

103 Modern breeding methods integrate next-generation sequencing (NGS) and phenomics
104 to identify plants with the best characteristics and greatest genetic merit for use as
105 parents in subsequent breeding cycles to ultimately create improved cultivars able to
106 sustain high adoption rates by farmers. This data-driven approach hinges on strong
107 foundations in data management, quality control, and analytics. Of crucial importance is
108 a central database able to 1) track breeding materials, 2) store experimental
109 evaluations, 3) record phenotypic measurements using consistent ontologies, 4) store
110 genotypic information, and 5) implement algorithms for analysis, prediction and
111 selection decisions. Because of the complexity of the breeding process, breeding
112 databases also tend to be complex, difficult, and expensive to implement and maintain.
113 Here, we present a breeding database system, Breedbase (<https://breedbase.org/>).
114 Originally initiated as Cassavabase (<https://cassavabase.org/>) with the NextGen
115 Cassava project (<https://www.nextgencassava.org/>), and later developed into a crop-
116 agnostic system, it is presently used by dozens of different crops and projects. The
117 system is web-based and is available as open source software. It is available on GitHub
118 (<https://github.com/solgenomics/>) and packaged in a Docker image for deployment
119 (<https://dockerhub.com/breedbase/>). The Breedbase system enables breeding
120 programs to better manage and leverage their data for decision making within a fully
121 integrated digital ecosystem.

122

123 Availability

124 <https://github.com/solgenomics>

125 <https://hub.docker.com/r/breedbase/breedbase>

126

127 License - MIT License.

128

129

130

131 INTRODUCTION

132 Modern plant breeding is a data intensive process requiring multiple diverse datasets to
133 be integrated and assessed in decision making. In classical plant breeding, promising
134 individuals are intentionally interbred to generate a diverse population of progeny, from
135 which individuals with the best phenotypic characteristics are selected to be used as
136 elite parents in subsequent breeding cycles or released as improved cultivars
137 (Breseghello and Coelho 2013). Modern plant breeding extends classical breeding with
138 the use of marker assisted selection (MAS) and genomic selections (GS) to augment
139 phenotypic selection (Ribaut and Hoisington 1998). Furthermore, with the emergence of
140 high-throughput phenotyping technologies as tools for breeding, the number of potential
141 phenotypes to be tracked has vastly increased (Andrade-Sanchez et al. 2014; White et
142 al. 2012).

143

144 The development of inexpensive genotyping technologies allow even small breeding
145 programs to acquire high-density genotyping data for a large portion of their germplasm.
146 The availability of this genomic data has enabled more efficient approaches to evaluate
147 important and complex traits in the breeding process (VanRaden 2008). One such
148 approach is genomic selection (GS), which combines genomic and phenomic data to
149 develop a predictive model that can be used to estimate genotypic or breeding values
150 (Meuwissen and Goddard 2001). Since genotyping is both less expensive and faster
151 than phenotypic selection, genomic selection can result in significant acceleration of the
152 breeding cycle with concomitant faster increases in gain. A challenge for genome-based
153 breeding methods is the establishment of an adequate data management infrastructure

154 to integrate the complex datasets spanning the breeding process (Volk et al. 2021). This
155 represents a severe constraint to mainstreaming predictive breeding to small breeding
156 programs, particularly in developing countries.

157
158 To address these data management challenges, we initiated a system called
159 Cassavabase (<https://cassavabase.org/>) for the NextGen Cassava project building on a
160 genomics codebase developed for many years for the Solanaceae called SGN
161 (<https://solgenomics.net/>) (Mueller et al. 2005a; Menda et al. 2008; Bombarely et al.
162 2011; Fernandez-Pozo et al. 2015a). With an initial focus on tomato and sequencing its
163 genome (Mueller et al. 2005b; Tomato Genome Consortium 2012), SGN already
164 contained a comprehensive genomics database with a strong phenotype management
165 component (Menda et al. 2008), a number of genomics-centric tools (Mueller et al.
166 2008; Tecle et al. 2010; Fernandez-Pozo et al. 2015b), and a rudimentary version of a
167 genotyping storage backend (Fernandez-Pozo et al. 2015a). Cassavabase is an open-
168 source, web-based breeding data management and analysis system built with the ability
169 to manage the genomic selection process (Tecle et al. 2014). As more instances of the
170 software were deployed for other crops, the system expanded to better meet each
171 project's needs by adding further breeding-related tools, such as image-based or near-
172 infrared spectroscopy (NIRS)-based phenotyping tools (Hershberger et al. 2021).. To
173 reflect that the underlying software and database are amenable to any crop and to
174 promote adoption by new communities, we named the system "Breedbase"
175 (<https://breedbase.org/>). Major clonal crops using Breedbase currently are cassava
176 (<https://cassavabase.org/>), yam (<https://yambase.org/>), banana (<https://musabase.org/>),
177 and sweetpotato (<https://sweetpotatobase.org/>), collectively known as the RTBbases
178 (<https://rtbbase.org/>); however, major non-clonal crops using Breedbase include wheat
179 (<https://wheat.triticeaetoolbox.org/>) and rice (<https://ricebase.org/>). Breeding and

180 research groups have adopted the system as well, such as the Gore Lab at Cornell
181 University (<https://gorelabbase.sgn.cornell.edu/>).

182
183 The purpose of Breedbase is to enable a '*digital ecosystem*' that contains an integrated
184 breeding workflow. Processes and data comprising germplasm banks, parental
185 selection, crossing design, experimental design, data collection, analyses, and decision
186 making tools are aggregated into a single system. This improves efficiency and reduces
187 data errors that can happen when using disjointed informatics tools, for instance when
188 transferring and restructuring data for analyses (Cobb et al. 2019). When data are
189 loaded into a database, many checks can be performed to make sure the data are
190 consistent and in line with specified quality control criteria.

191
192 Many breeders, especially in smaller programs that cannot allocate resources to data
193 management tools, maintain their data in spreadsheets. While spreadsheets provide a
194 straightforward way to manage data and analyses, they suffer from a number of
195 drawbacks, even with relatively small volumes of data. For example, it is difficult to
196 precisely merge data across different spreadsheets, often resulting in errors and data
197 quality issues, or to visualize or analyze data across spreadsheets. Data in
198 spreadsheets are typically not normalized, resulting in typographical issues, inconsistent
199 identifiers, liberal use of synonyms, and similar issues that make the data hard to
200 aggregate. Nevertheless, the largest problem with spreadsheets is that their storage is
201 not centralized; in fact, they are often stored on personal computers and laptops, often
202 in multiple inconsistent versions, with potentially limited backup strategies and little
203 recourse if accidental data loss occurs or if a person leaves the breeding program,
204 taking all the breeding data with them. Breeding programs can be very large,
205 encompassing many locations with many collaborators; as such, spreadsheets hinder

206 collaboration because data cannot be accessed in a consistent state by many people at
207 once. Furthermore, with genome-based breeding, spreadsheets become unworkable,
208 as it is difficult to maintain and analyze potentially very large genotypic data sets in
209 spreadsheets in any useful way. It is important to note that using a database is not
210 sufficient for managing a modern breeding cycle - the entire breeding process needs to
211 be integrated around the database to create an efficient digital ecosystem.

212

213 Breedbase implements a robust system of breeding workflows, data management
214 procedures and analysis tools to address breeder informatics problems. Here we
215 present the rationale, design, implementation and major use cases for Breedbase.

216 MATERIALS AND METHODS

217

218 Implementation

219 The Breedbase data architecture is built around a Postgres (<https://postgresql.org/>)
220 relational database with a schema that is mainly derived from Chado (Jung et al. 2011),
221 with some historic, pre-Chado tables from SGN, as well as minor customizations
222 (Fernandez-Pozo, Menda, et al. 2015) (**Figure 1a**). In relational databases, information
223 is systematically structured into concepts represented as tables (“normalization”), a
224 format that facilitates many aspects of data management. The information in the
225 different tables can be joined based on primary and foreign keys, which are usually
226 numeric values assigned to every row in a table. For some data types, such as
227 genotypic data, Breedbase uses non-SQL extensions built into Postgres, such as
228 JSONb-based data structures (Morales, Bauchet, et al. 2020). The application layer is
229 implemented in Perl, using the Moose object system, based on the Model-View-
230 Controller (MVC) Catalyst web framework (<http://www.catalystframework.org/>), with

231 Mason as the templating toolkit (<https://metacpan.org/pod/Mason>). The system uses an
232 object-relational layer based on DBIx::Class, with the main Chado classes organized in
233 the Bio::Chado::Schema namespace. For statistical analyses and some of the data
234 visualizations, the R language and add-on R packages (<https://r-project.org/>) are used.
235 Image analyses and machine learning models are implemented in Python TensorFlow
236 (<https://www.tensorflow.org/>) and OpenCV (<https://opencv.org/>) (Morales, Kaczmar, et
237 al. 2020). The frontend graphical user interface (GUI) development has recently
238 transitioned away from Mason components to JavaScript, with a heavy reliance on
239 asynchronous JavaScript requests. Almost all functionalities are implemented as
240 RESTful services, allowing for a more interactive user experience and reusable
241 codebase. JavaScript frameworks used for the GUI include JQuery (<https://jquery.org/>),
242 D3.js (<https://d3js.org/>), Bootstrap (<https://getbootstrap.com/>) and Brapi.js
243 (<https://brapi.org/>). The entire Breedbase system is built on open source software and is
244 packaged in a Docker image for deployment (<https://docker.com/>). For interoperability
245 with other breeding database and tools, Breedbase implements the BrAPI 2.0
246 specification (Selby et al. 2019).

247
248 In terms of user interface, the goal of Breedbase is to provide a standard, modern web
249 interface for all breeding tools. Breedbase is essentially a cloud-based app, obviating
250 the need for the user to install any software. For anyone with web-browsing experience,
251 the interface should be intuitive and straightforward, and it is continuously improved
252 based on user driven feedback. In Breedbase, processes are presented in an
253 interactive workflow system, providing step-by-step guidance to breeders and users in
254 accomplishing specific tasks. A few of the widely-used interfaces include the Wizard,
255 Lists, and Datasets tools, which will be described in more detail later.

256

257 Use cases

258 The initial development of BreedBase focused on addressing the data collection and
259 management stages necessary to facilitate genomic selection within a breeding
260 program, including:

261

- 262 ● Manage accessions and pedigrees in the database, with ontology-based
263 descriptions and support for rich metadata including images
- 264 ● Design field layouts and track all field metadata
- 265 ● Load historical data from breeding programs
- 266 ● Collect phenotypic data on tablets in the field and upload the subsequent
267 phenotypes
- 268 ● Manage genotypic data associated with the accessions
- 269 ● Enable genome-based predictive breeding by calculating correlations between
270 phenotypes and genotypes, and predict phenotypes from genotypes [the solGS
271 tool (Teclé et al. 2014), <https://cassavabase.org/solgs/search>]
- 272 ● Support controlled crossing using customized tracking tools

273

274 More recently, a number of other use cases were pursued:

275

- 276 ● Advanced statistical analyses including Principal Component Analysis (PCA),
277 stability analysis (AMMI) (Duarte and Pinto 2002) heritability calculations
278 (Holland, Nyquist, and Cervantes-Martínez 2010), mixed model analysis, and
279 genome-wide association studies (GWAS)
- 280 ● Marker-assisted breeding
- 281 ● Processing and analysis of unoccupied aerial vehicle (UAV) image data
- 282 ● Image analysis

283 ● NIRS data storage and analysis

284

285 Plant breeding operations requiring decision support within a growing season include
286 three broad activities: crossing, evaluations, and selections. These activities typically
287 include setup of crossing and trial experiments (design, labeling), data and seed
288 collection, genotyping, and subsequent statistical analysis. Breedbase offers support for
289 each of these components through online tools. To streamline accessibility and usage
290 for key routine activities, Breedbase has established workflow components. Each
291 workflow offers the user a guided process for a targeted activity. For example, the trial
292 creation workflow comprises trial creation, planting material and checklist creation,
293 randomization and statistical design selection, field visualization and storage. During
294 this process, field trial experiment parameters (see Phenotyping Trials section) are input
295 into Breedbase and the relevant experimental design is calculated using open source R
296 libraries such as Agricolae (Mendiburu et al.) or Digger (Coombes 2009). The
297 experimental layout is calculated and displayed, and can be reviewed and potentially
298 improved by re-running randomization before the trial design is stored in Breedbase.
299 Additional parameters such as field management factors (ie: agronomic management or
300 fertilizer application) can also be entered. Similar workflows exist for other activities,
301 such as phenotyping and genotyping.

302 Development Process

303 The development process can be broadly described as agile (Shore, Chromatic, and
304 Warden 2008; Beck and Andres 2004), in which shorter-term goals are defined and
305 implemented, and subsequently further refined based on new feedback from users;
306 agile teams provide for short release cycles and continuous improvement to the
307 software (**Figure 1b**). Progress is tracked using a version control system with built-in
308 issue tracking software (GitHub, <https://github.com/>). New features are discussed with

309 breeders and other stakeholders. Issues and bugs discovered in Breedbase are tracked
310 on the public GitHub issue tracker. A programmer is then assigned to a ticket, and will
311 create an issue-specific topical git branch in the relevant code repositories, and
312 implements the required changes in the branch, including tests and edits to the user
313 documentation. When the implementation is ready for release, a pull request is
314 generated on GitHub and a reviewer is assigned. In the review, the code is verified for
315 errors, programming style, tests, and documentation. If the reviewer approves the pull
316 request, the code is merged into the master branch. The test-driven software
317 development approach is tightly integrated with our development process, consisting of
318 unit and integration tests. A ticket meeting is held once a week and all open pull
319 requests and important tickets are discussed. If all the pull requests were merged
320 successfully, and no issues are discovered with tests or other checks, a new release tag
321 is created, the new version is deployed in production, and a new Docker image is
322 released. Since Breedbase is open source, programmers outside of the core
323 development team are able to make contributions to the code base via the same
324 process. The Breedbase project has had 40+ contributors addressing various issues
325 and improvements (<https://github.com/solgenomics/sgn/graphs/contributors>).

326 Ontologies

327 A key aspect of data integration is the necessity of standardization. Breedbase is based
328 on the Chado database schema, which relies heavily on controlled vocabularies and
329 ontologies to describe its data, and requires numerous ontologies for its internal
330 functioning. In many ways, it can be described as an ontology-based database. For the
331 breeding application, data standardization in the form of trait catalogs is especially
332 important when several sites or breeding programs share data in the database. Without
333 standardization, the data would not be comparable, limiting the utility of an integrated
334 database. The creation and maintenance of trait ontologies is a considerable task. The

335 Crop Ontology (CO) project was developed by CGIAR to define and maintain relevant
336 breeding ontologies (Shrestha et al. 2012). All the RTBbases use the Crop Ontology
337 vocabularies and collaborate with Crop Ontology and breeders to improve and expand
338 these vocabularies (Arnaud et al. 2020). If no ontologies are available, they have to be
339 created, which can be a lengthy and arduous task. The Protégé tool
340 (<https://protege.stanford.edu/>) (Musen 2015) is commonly used by curators for editing
341 ontologies before upload to Crop Ontology and Breedbase. The Trait Dictionary
342 Template along with the Guidelines (Pietragalla et al, 2020), available in the CO
343 website, remain useful to collect the trait details from the research community and reach
344 consensus. Each species is allocated a code by the CO coordination team to identify
345 the ontology and crop repositories are created in the Planteome Github to secure the
346 ontology version management. An online term submission form is accessible in
347 Breedbase for users wishing to suggest missing traits or modifications to the Crop
348 Ontology (<https://submit.rtbbase.org>).

349 Interoperability and BrAPPs

350 Databases must interoperate with a variety of tools to perform their functions in data
351 acquisition, analysis, and data export. Recently, a standard called the Breeding
352 Application Programming Interface (BrAPI; <https://brapi.org/>) was developed to
353 exchange breeding data (Selby et al. 2019), which breeding databases can implement
354 to provide a standard interoperability layer. Standardized application programming
355 interfaces (APIs) allow Breedbase to integrate and interface with a broader set of BrAPI
356 enabled applications, or BrAPPs, that can be written across diverse programming
357 languages including Android, R, and Javascript. The BrAPI R package allows data
358 retrieval from Breedbase for further statistical processing within the R environment.
359 Javascript based BrAPPs provide dynamic visualization of plant breeding data, such as
360 pedigrees exploration, experimental field maps, and data from multiple trials. BrAPPs

361 can interact with data from any BrAPI compliant database, such as Breedbase or the
362 Breeding Management System (BMS) (Figure 1a). Activities such as dynamic data
363 filtering, trial comparison, box plotting, and a comparative genetic map viewer are also
364 implemented with BrAPPs on Breedbase. Breedbase fully supports BrAPI version 2.0
365 and is committed to updating the system for future versions of this essential
366 infrastructure.

367 Querying Breedbase

368 Breedbase has a number of query options, which are grouped in the “Search” menu.
369 The most important data types each have a search (“Accessions and plots”, “Trials”,
370 “Organisms”, “Crosses”, etc). A powerful combined search is available in the form of the
371 Search Wizard (**Figure 2**).

372 The Search Wizard and Datasets

373 The Search Wizard allows users to slice their data in different dimensions, such as
374 breeding programs, locations, years, and so forth. The data in the database can be
375 thought of as a multi-dimensional cube which is cut along different dimensions,
376 providing an intersection that represents the data of interest. This approach is
377 conceptually related to a query method called Online Analytical Processing (OLAP)
378 (Celko 2006). The current Wizard presents four boxes, for four different dimensions,
379 which can be selected using pull down menus (**Figure 2**). For example, a user who is
380 interested in the performance of cassava clones evaluated by IITA in 2017 and 2018 at
381 the Mokwa station in Nigeria can use the wizard to find this information. Working from
382 left to right, the user selects as the first dimension “Breeding Programs”, which displays
383 all the breeding programs in the database in the first box. The user then selects “IITA”
384 from the individual breeding programs listed in the box. When the user selects “Years”
385 in the second box, all the years for which data for IITA exist are listed. In this example,
386 the user selects 2017 and 2018. Finally, after selecting locations in the third box, the

387 user specifies "Mokwa". When trials or accessions are selected, phenotypic and
388 genotypic data corresponding to the selection can be downloaded using buttons below
389 the Wizard boxes. The Wizard also allows the combination of current selections to be
390 stored in the database under a user-given name, representing an intersect of data of
391 interest in the database. This stored selection is called a "dataset". Datasets are used
392 across Breedbase to efficiently reference a complex query with a simple, assigned
393 name. Tools that support the dataset concept in Breedbase include solGS, GWAS, the
394 heritability tool, the stability analysis, and the general mixed model tool.

395

396 Quick Search

397 A quick search is provided in the upper right corner of the menu bar that searches a
398 keyword across all data types in the database, and is a fast way to retrieve named
399 objects such as stocks and genes.

400 Special searches

401 Topic specific searches are available from the Search menu, including a trial search, a
402 trait search, searches for genotyping data (including genotyping protocols, projects and
403 plates), an image search that searches image descriptions and associated tags, and a
404 user search that searches the users of the database. All these searches work in a
405 straightforward and consistent way: a search form is filled in with search criteria, and the
406 search is submitted to the database. A list with matched search results is displayed,
407 from which links are provided to the corresponding detail pages.

408 Analysis Tools

409 Breedbase is more than a static collection of data, as it enables users to explore and
410 analyze data in the database. Once data is uploaded to the database, users can view
411 summary statistics, evaluate phenotypic variances, and identify observations with
412 missing or outlier data. They can filter observations in a trial based on a range of trait or

413 traits values. For an experiment phenotyped in multiple environments, they can evaluate
414 trait performance across environments using pairwise comparison scatter plots and
415 histograms.

416

417 Breedbase also has tools for ANOVA, correlation, principal component analysis, data
418 partitioning using K-means clustering, genomic prediction, genome wide association
419 study, selection index calculation, genetic gain visualization, and linear mixed models.

420 With the Search Wizard, as explained above, users can construct datasets that can be
421 used as inputs to various tools. Most tools follow a similar blueprint in terms of user

422 interface: (1) select the dataset of interest from a drop-down menu of all available

423 datasets, (2) adjust parameters for the tool, (3) submit the calculation for analysis, and

424 (4) display the results. For some tools that require heavy computation, an email can be

425 optionally sent to the user with a link to the results. Query implementation is a relatively

426 complex task in the programming of a tool, but the Wizard enables the modularization of

427 algorithms into Breedbase with relatively little glue-code, facilitating tool coverage

428 expansion. Results such as predictions from solGS and adjusted means from mixed

429 models can be saved in the database as analysis results. These results can be used

430 like primary data in downstream analyses such as the selection index tool to help

431 identify favorable germplasm.

432

433

434

435 Managing a Breeding Program Using Breedbase - Accessions,
436 Phenotyping, Crossing, and Genotyping data

437 General principles

438 Plant breeding involves the collection of a wide variety of data types at different time

439 points and locations, and across different scenarios (e.g., field, laboratory, seed

440 storage). To give users flexibility and mobility in data collection, smartphone-based

441 applications are often required. Android applications, such as PhenoApps

442 (<http://phenoapps.org/>), are developed with this perspective (Rife and Poland 2014).

443 Breedbase has adopted the PhenoApps tool suite created by Kansas State University

444 (KSU).

445

446 PhenoApps include applications for phenotyping (Field Book), cross management

447 (Intercross), sample collection (Coordinate), and inventory management (Inventory).

448 Breedbase has worked to build in native support for these applications and integrate

449 them into best practices workflows. Since internet access is not available at all field

450 sites, the functionality has been developed to allow configuration of these applications

451 prior to field data collection. Field layouts, plant accessions, and traits to be measured

452 can be loaded onto mobile devices through special interfaces in Breedbase. Following

453 collection, data is imported back into Breedbase. Because all the trial information in the

454 collection device was initially downloaded from Breedbase, required identifiers can

455 easily be matched with the existing data in the database. This process is called “round-

456 tripping”, and is a crucially important concept for high quality data management.

457 List management

458 Breeding activities often require the maintenance of lists of various types - for example,
459 a list of accessions to plant, traits to measure, or trials to evaluate - and, consistent with
460 digital ecosystem principles, these lists should be managed entirely through the
461 database. Accordingly, Breedbase implements comprehensive list management
462 functions. By default, lists are associated with the user that creates the list. The main list
463 interface can be reached by clicking on the Lists link on the top right of the toolbar,
464 which appears when logged in. A dialog appears that allows users to view, create and
465 edit new lists. Each list has a data type from an internal ontology called 'list_type', which
466 includes terms for 'accessions', 'trials', 'traits', 'years', etc. Lists are collections of text
467 elements that correspond to names of database objects. Lists can be validated against
468 names that are already present in the database. A validated list can then be used to
469 submit data to various tools, including the Wizard, right on the website. Sometimes, it
470 can be useful to share a list with other users, and this can be achieved by making a list
471 public by clicking the appropriate checkbox in the list detail view. Public lists are shown
472 in a separate section, and become visible to all users. They can be "unshared" if
473 needed.

474

475 Germplasm management

476 Germplasm is the foundation of a breeding program and plays a similarly important role
477 in a database such as Breedbase. In plant breeding programs, tracking and
478 characterization of germplasm is a major challenge. Germplasm in this context includes
479 accessions, stocks, varieties, or, in clonal crops, clones. Breedbase commonly uses the
480 term "accession". Breedbase is pre-populated with the complete plant section of the
481 NCBI taxonomy database, defining all known species with their associated genus,
482 abbreviation, common name, and GenBank taxon identifier. Researchers using

483 Breedbase can usually find their crops of interest within the 100,000+ organisms
484 available. Accessions are always created in association with one of these organisms.

485

486 Some instances of Breedbase, such as Cassavabase and Sweetpotatobase, are
487 designed to only contain germplasm of their respective species; however, it is possible
488 for a single instance of Breedbase to be used for a variety of crop species. Combining
489 many crop species into a single instance can complicate the search interfaces and lead
490 to bloated databases; however, aggregating all data allows for more consistent and
491 queryable data. Alternatively, separating instances can lead to potentially duplicated
492 and inconsistent data, but can be beneficial for fostering communities.

493

494 In Breedbase, there are two distinct concepts that describe accessions: (1) an
495 accession that can be ordered from a seed bank, which may have been selfed and
496 could be genetically quite pure, or landraces. These are “long-term use” accessions (ie:
497 historical germplasm, parental inbred lines) , which may be actively maintained and can
498 be obtained easily; whereas (2), are “short-term use” accessions (ie: intermediate
499 generations) that are produced in a breeding program and may go through a few rounds
500 of selection, but most of which will be discarded in the process. These accessions may
501 also not be genetically pure, as they may result from crosses between relatively distant
502 parents.

503

504 To create an accession in Breedbase, only a unique name and the organism species
505 name are required. As with all objects stored in a relational database, Postgres will
506 create a primary key identifier for each object, using a data structure called a sequence,
507 which is used to link the accession to other objects in the database using a foreign key.
508 This means that even if the accession name is modified, it will still retain all the

509 connections to other objects of the original entry. Germplasm can be further annotated
510 with configurable properties from the Multi-Crop Passport Descriptors (MCPD)
511 standards (Food and Agriculture Organization of the United Nations 2018) and BrAPI
512 standards (Selby et al. 2019); these properties include 'variety', 'donor', 'donor institute',
513 'donor PUI', 'country of origin', 'institute code', 'institute name', 'notes', 'accession
514 number', and 'PUI'. Germplasm can be added to the database using the interactive list
515 tool (see previous section) or an Excel file upload; the Excel file upload also allows for
516 storing and updating of all attributes listed above. The first step in the initiation of a
517 breeding program is to load relevant accessions into the database. This is critical, as
518 the naming of accessions is often not uniform between breeding programs and the
519 community at large. In some cases, a single name can refer to several different
520 accessions or a single accession may have many different names or synonyms, often
521 the result of historical transcription error or case inconsistency. Before the first upload it
522 is therefore essential to define a standard unique name and set of possible synonyms
523 for each accession. Though Breedbase allows for synonyms of accession names, they
524 should also be unique. It is best practice to use synonyms only to find accessions and
525 not when performing routine tasks with the database during the breeding
526 process. Whenever new accession names are encountered, Breedbase provides a
527 workflow to compare new names to all existing accessions in the database. In this
528 workflow, a user can consolidate synonyms, for instance to add 'Tx 303' as a synonym
529 of 'TX303'.

530
531 After initial accession upload, it is often necessary to add more accessions, increasing
532 the chance of generating duplicated accessions in the database, or other upload issues.
533 As is the case with synonyms, many of these problems result from poorly defined
534 accession identifiers with capitalization inconsistencies and special characters such as

535 slashes, dots, dashes, underlines, and spaces. Although we recommend avoiding such
536 special characters, especially in primary identifiers, it is not always feasible, notably with
537 legacy data. To ease upload and tracking of such cases, Breedbase has a fuzzy search
538 (also called approximate string matching search) component, enabling an accurate
539 quality control of existing similar germplasm names in the database.

540

541 Phenotyping Trials

542 Phenotyping trials are a core activity of plant breeding programs, and must be carefully
543 designed. Trial designs can either be generated directly in Breedbase using the
544 integrated, comprehensive trial design tool or uploaded using Excel files formatted with
545 a Breedbase-provided template. Trial metadata fields include breeding program,
546 location, name, trial type, year, plot dimensions, field size, and trial design type.
547 Supported statistical trial design types currently include alpha lattice, lattice, augmented,
548 split plot, partially replicated, and Wescott designs. Designs should also include the
549 ordinal row and column positions of each plot as it is planted in the field, so Breedbase
550 allows this information to be added either during or after design storage. Once a trial
551 design is finalized, it is stored in the Breedbase schema. Within Breedbase, a field trial
552 links phenotypic observations to the experimental layout under a specific statistical
553 design.

554

555 Row crops usually use the concept of plot as the minimal entity for data collection, but
556 many specialty crops (*i.e.*, vegetables) require data collection on a per plant or per
557 tissue basis. Breedbase allows plant- and tissue-level entry creation for each plot in a
558 trial, resulting in database entries and identifiers at each level, which can also be
559 encoded in barcode labels for data collection.

560 Crossing

561 To collect data from crosses, Breedbase requires the creation of a top-level crossing
562 experiment; the crossing experiment is defined with a unique name, a breeding
563 program, a location, a year, and a description. The individual crosses performed are
564 then stored under the crossing experiment and defined by a cross unique id, parents,
565 and a cross type. The cross type can be one of the following: biparental, self, sib, open-
566 pollinated, bulk, bulk selfed, bulk and open-pollinated, doubled haploid, polycross,
567 reciprocal, or multicross. Depending on the type of cross performed, different metadata
568 must be provided; for example, in a biparental cross, information from both the male
569 and female parent is required, whereas in an open-pollinated cross, information on only
570 the female is required. In the case of an open-pollinated cross, a population name
571 representing a group of male germplasm can be given as the male parent. In addition to
572 cross unique id, which captures specific details of each cross, users have the option to
573 group crosses having the same parental genotypes via family name for downstream
574 progeny analysis.

575

576 Breedbase tracks parental information from crosses in two ways: (1) through the
577 accession names of the female and male parents, allowing for simple ancestry tracking
578 of AxB pedigrees for the progeny from a cross. When a cross is created in Breedbase,
579 the pedigree between progeny and parental germplasm is automatically created as well.
580 This first form of parental tracking is applied in all cases when a cross is created in
581 Breedbase. (2) through the plot or plant names of the male and female parents. The plot
582 or plant names of the parents are related to the field trial in which they are planted, as is
583 described in the above field trial section. This approach allows detailed tracking of
584 female and male parents used in crossing, but is optional in Breedbase because of the
585 difficulty in recording this information in many cases.

586
587 Recording information on parental plots is facilitated by mobile data collection platforms.
588 Of note are customized Open Data Kit (ODK) Android applications, such as BTract and
589 the PhenoApps app Intercross. BTract assigns and prints a unique cross barcode label
590 after scanning barcodes to track the precise male and female plots or plants involved in
591 the pollination. Through ODK data synchronization, the cross information can be
592 uploaded into Breedbase. Intercross can be used to scan parental barcodes and
593 associate a unique cross id to the performed cross. The output from Intercross can also
594 be uploaded directly into Breedbase.

595
596 In crossing experiments that include evaluation of crosses, Breedbase can store
597 annotations regarding properties of the cross. Default properties include pollination
598 date, tag number, number of flowers, number of bags, number of fruits, and number of
599 seeds; however, these properties are set in the configuration file for the Breedbase
600 instance, allowing researchers flexibility in defining these terms. Breedbase also
601 supports tracking of tissue culture samples.

602
603 Crosses can be created individually using an interactive interface on Breedbase or can
604 be uploaded in bulk using an Excel spreadsheet by providing cross unique ids, cross
605 types, and parents involved. Once each cross unique id is saved in Breedbase,
606 additional data can be added or uploaded using the cross unique id as an identifier.
607 Progeny of the cross can be saved as new germplasm in the database, automatically
608 creating pedigrees for the new germplasm.

609
610 Genotyping Data
611

612 High-density genotyping data are a complex data type that have become an important
613 resource in modern breeding programs due to the advent of low-cost next-generation
614 sequencing (NGS) and genotyping technologies (Thomson 2014). Breedbase offers
615 simple laboratory information management functionalities from field tissue sampling to
616 SNP data storage. Functions include tissue samples collection and tracking via plot
617 barcodes and PCR plate formats (ie: 96 or 384 wells), genotyping protocol definition,
618 data storage and subsequent analytics (Morales, Bauchet, et al. 2020; Teclé et al.
619 2014).

620
621 The primary means of organizing genotyping data between sequencing events is the
622 ‘genotyping protocol’ in Breedbase. A ‘genotyping protocol’ consists of a specific set of
623 genotypic markers and records all metadata about how the genotypes were produced,
624 including the reference genome and specifics about, analytical platform and related
625 variant calling software. The ‘genotyping protocols’ can be grouped in Breedbase under
626 a ‘genotyping project’ which displays all relevant genotyping data and provides an
627 overview, which is especially useful for very active genotyping programs.

628
629 Multiple genotyping technologies can be stored in Breedbase from low density
630 genotyping (ie: Kompetitive allele-specific PCR, KASP) to high density genotyping such
631 as Genotyping-by-sequencing (GBS) or DArT-seq (Elshire et al. 2011; Semagn et al.
632 2014; Kilian et al. 2012). The preferred method for uploading high-density genotyping
633 data to Breedbase is through variant call format (VCF) files. VCF provides for compact
634 representation of genotypic scores for large numbers of samples and markers (Danecek
635 et al. 2011). PostgreSQL non-relational functionalities allow Breedbase to store high-
636 density genotyping data in JavaScript Object Notation (JSON) structures within the
637 larger relational database schema (“ISO/IEC TR 19075-6:2017” 2018). Breedbase

638 particularly relies on the binary JSON (JSONb) data type for compressed data storage
639 and faster retrieval (Morales, Bauchet, et al. 2020).

640

641 Genotyping data can be queried alongside relationally stored phenotypic and
642 experimental information for analyses, including computation of a genomic relationship
643 matrix (GRM) for user specified germplasm and computation of a genome-wide
644 association study (GWAS) for user specified germplasm and phenotypic traits
645 (VanRaden 2008). Queries spanning specific markers or marker sets and experimental
646 information can be readily constructed. Genotyping data results can be downloaded as
647 VCF files from the Search Wizard web-interface. The genotyping data are also used in
648 the Genomic Selection tool, solGS, to predict GEBVs of genotyped lines.

649 Authentication and Authorization

650 During breeding processes, a potentially large number of people will need to access the
651 database to download, upload, modify or delete data. This requires a fine-tuned layer of
652 authentication and authorization management in the database. Breedbase requires a
653 user to login for most functionalities (authentication). Every user account is associated
654 with “roles” that determine what the user will be allowed to do in the system
655 (authorization). Currently, there are three major roles: user, submitter and curator. The
656 user role allows read-only access. With the submitter role, a user can upload data, and
657 can modify or delete data that they themselves uploaded. The curator role allows a user
658 to modify any type of data. In addition, every breeding program in the database has a
659 corresponding role that controls authorization over specific breeding program activities,
660 such as creating and uploading trial data.

661

662 Cassavabase, the flagship Breedbase database
 663 Cassavabase (<https://cassavabase.org/>) is the breeding database for the NextGen
 664 Cassava project (<https://nextgencassava.org/>). The NextGen Cassava partners, IITA
 665 (Ibadan, Nigeria), NRCRI (Umudike, Nigeria), NaCRRRI (Namulonge, Uganda), TARI
 666 (Ukiriguru, Tanzania), Embrapa (Cruz das Almas, Brazil) and CIAT (Cali, Colombia) use
 667 Cassavabase for their breeding programs, starting as early as 2014. To date,
 668 Cassavabase has accumulated an immense amount of cassava breeding data (**Figure**
 669 **1c**), consisting of information on more than 500,000 cassava accessions, characterized
 670 by over 19 million phenotypic measurements in over 4,000 trials, and nearly 35,000
 671 genotyping experiments. This shows that the Breedbase system can scale to fairly large
 672 datasets and large, multi-institute and multi-national programs.

673 Other instances of Breedbase

674 In addition to Cassavabase, Breedbase has been deployed for various crops, notably
 675 for other Roots, Tuber and Banana (RTB) crops (<https://rtbbase.org/>) in the CGIAR:
 676 banana, (<https://musabase.org/>), sweetpotato (<https://sweetpotatobase.org/>) and yam
 677 (<https://yambase.org/>). In addition, several dozen Breedbase instances are currently
 678 deployed for other crops, such as rice (<https://ricebase.org/>), wheat
 679 (<https://wheat.triticeaetoolbox.org/>), oat (<https://oat.triticeaetoolbox.org/>), kelp
 680 (<https://sugarkelpbase.org/>), potato and maize. While the afore-mentioned projects use
 681 Breedbase for mainly breeding informatics purposes, other Breedbase instances focus
 682 on genomics. These include SGN (<https://solgenomics.net/>, (Fernandez-Pozo, Menda,
 683 et al. 2015), which focuses on tomato and other *Solanaceae*, fern
 684 (<https://fernabase.org/>, (Li et al. 2018), *Erysimum* (<https://erysimum.org/>, (Züst et al.
 685 2020)) and milkweed (<https://milkweedbase.org/>). In addition, a Breedbase instance has
 686 been deployed to characterize a tritrophic vector-borne disease system, the citrus
 687 greening disease (<https://citrusgreening.org/>) (Saha et al. 2017). An instance named

688 ImageBreed has been deployed for high-throughput imaging of maize and alfalfa field
689 experiments (<https://imagebreed.org/>) (Morales, Kaczmar, et al. 2020). A number of
690 academic labs and breeding companies also use Breedbase for data management
691 within their programs. The Breeding Insight project (<https://breedinginsight.org/>), which
692 creates breeding databases for USDA breeding programs, has also adopted the
693 Breedbase system as a foundation for their breeding solutions.

694

695

696

Box 1.

697

698 **Providing data management tools for small grains breeders: The Triticeae**

699 **Toolbox adaptation of Breedbase**

700

701 As documented in this article, Breedbase provides many features for working breeding
702 programs. The mission of The Triticeae Toolbox (T3) is to provide these features to a
703 diverse audience of small grains breeding programs, by mandate in the United States,
704 and by extension globally.

705

706 The development of T3 is motivated by the belief that larger datasets provide greater
707 power to identify genetic effects that are relevant to all breeders. Across wheat, oat, and
708 barley, T3 stores 5,600 trials, comprising over 1,800,000 phenotypic data points on over
709 30,000 lines with genotype data. From there, T3 seeks to provide breeders with results
710 from analyses that tap into these data, in the hope that this will help breeders gain
711 insights from their own data. The primary example we have in this area is a function to
712 show marker trait associations identified among all trials submitted to T3 with adequate
713 marker density, and meta-analyzed to determine robust associations across trials. The

714 next milestone on the roadmap of this function is to develop marker imputation
715 functionality on T3 that will present genotype trials with uniform high-density marker
716 scores, enabling meta-analysis over more trials. Indeed, marker data are a critical
717 rationale for T3's mission: the database contains data on many lines that now are
718 connected to current populations primarily through the marker alleles segregating.

719
720 An important advantage of a web-based data management platform is that it links the
721 data to the world of knowledge available on the web. T3 provides that connectivity by
722 providing links to external information on markers, traits, and germplasm. Our primary
723 partners in that regard are GrainGenes, Wheat Expression Browser, and the Wheat
724 KnetMiner (Hassani-Pak et al. 2021). For example, a marker trait association close to a
725 gene can be used to connect that trait to [JBrowse](https://jbrowse.org/) (https://jbrowse.org/), to gene
726 expression data ([expVIP](#) and [EMBL-EBI](#)) or to a knowledge network, [KnetMiner](#). Traits
727 in Breedbase are defined using collaborative ontologies crucial to forging these links:
728 the ontologies represent agreements on naming traits and gene functions that enable
729 meaningful bridges across knowledge platforms.

730
731 The diversity of T3 users means that they will not operate together as an integrated
732 breeding organization. Rather each breeding program submitting data to T3 will want
733 data privacy and ease in determining what data becomes incorporated into the public
734 production database. Currently, all data on the production database is available to
735 anyone. We plan on implementing privacy settings specifying data visibility as public or
736 restricted. Absent this feature, we now work with a few users by providing them with
737 separate instances of T3 that are not publicly visible but can easily transmit datasets to
738 the T3 production database when ready.

739

740 The wide range of T3 users also means that we expect them to have varying degrees of
741 familiarity with the Breedbase platform. To allow users to test the addition and
742 modification of datasets without modifying curated data by mistake, T3 has created
743 sandbox instances for each crop. Users can freely upload data to the sandbox, ensure
744 that the uploaded data added to the database is correct, and then easily publish the
745 data to the production instance. A data curator checks the submitted data before adding
746 it to the production database. The Breedbase system was crucial in establishing these
747 features and reduced duplication of effort.

748

749

750

751 Box 2. Usage Example

752

753 Recently, Obgonna and colleagues leveraged legacy breeding data to investigate the
754 genetic architecture of cyanide content in cassava, a key trait in food safety (Ogbonna
755 et al. 2020). Authors performed a retrospective analysis, mining historical cyanide data
756 from the African IITA breeding program (18 locations, 23 years and 393 trials) and
757 Colombian CIAT (41 locations, 11 years and 155 trials) program from the Breedbase
758 instance cassavabase.org. Recycling open source, standardized, breeding data in
759 conjunction with novel genotypic data provided a high statistical power and allowed the
760 detection of key loci controlling cassava root cyanide content using GWAS. Such loci
761 would otherwise have gone undetected, and was identified only because of the
762 availability of the Breedbase digital ecosystem.

763

764

765 DISCUSSION

766 Breeding is a complex process involving many different types of data, especially
767 considering genome-based breeding methods at the current state of the art. Creating
768 and maintaining breeding databases is therefore generally considered to be time-
769 consuming and expensive. Many large breeding companies maintain their own
770 databases and software for managing breeding processes and selection, but this is not
771 an option for smaller programs. The lack of bespoke databases is especially true in
772 resource poor areas of the world, where the need for plant improvement is often the
773 greatest. A free, user-driven and open source platform such as Breedbase that
774 integrates a complete digital ecosystem for breeding will help close the gap for these
775 programs as well as many smaller to mid-sized organizations. Still, Breedbase
776 databases can scale significantly to large breeding programs with hundreds of
777 thousands of accessions and millions of phenotypic scores.

778 Integration in breeding programs

779 Even the best breeding data management tools will fail to deliver if breeding programs
780 do not use them or use them incorrectly. A significant effort is required to integrate a
781 breeding database into the workflow of a breeding organization, as data management is
782 central to the work of modern breeding programs but remains a shortcoming. Breeding
783 activities need to be closely tracked; to ensure complete integration, all materials,
784 operations, and operators need to be systematically recorded and reviewed throughout
785 the process. This is important to enable analyses, improve data quality, and to identify
786 sources of errors in real time and *post hoc*.

787

788 It is important for breeding programs to work closely with groups that have significant
789 experience in data management, which can also help the breeding programs to
790 understand their needs, and to train staff better in the use of the database. In the RTB

791 breeding programs, we found it to be helpful to designate specific staff as Data
792 Managers, who receive extensive database training. Data Managers have spent time at
793 the BTI to learn more about the database developments, and can provide additional
794 training and help on the ground in the breeding programs. They also provide timely
795 feedback on the tools and features based on their first hand experiences, which is vital
796 for the improvement of the database. We have put a significant effort in user training
797 through in-person workshops, reciprocal visits, and training materials, such as a
798 complete on-line manual, slideshows, and most recently a youtube channel with
799 recorded workshops.

800
801 In our experience, one of the bottlenecks in implementing a breeding database is the
802 availability of standardized trait ontologies for the crop in question. Especially in larger
803 projects, it can be difficult for all breeders to agree on a common ontology, including
804 common sample preparation and measurement protocols, as well as measurement
805 units. Without this standardization, a database loses much of its appeal as it becomes
806 impossible to aggregate and reconcile disparate data. This challenge cannot be
807 understated as it is a major obstacle especially when phenotypic data is collected
808 across different locations for a variety of crops and has to be stored in a single
809 integrated system. We have focused on developing ontologies and common
810 vocabularies to address this issue but it can be harder than expected, as there are often
811 diverging and strong opinions on these matters. In addition, breeding programs
812 introduce new traits to be measured, for example, quality traits, and there needs to be a
813 process to integrate such new terms into the ontology. Fortunately, the Crop Ontology
814 project (Shrestha et al. 2012) has created trait ontologies for a wide range of crops,
815 which we contribute to and many Breedbase instances rely on. Crop Ontology has also

816 defined processes for updating and developing the ontologies, which allows new traits
817 and methods to be introduced to breeding programs with relative ease.

818 Future developments

819 Progress in the last few years in digital agriculture has been enormous and will continue
820 to be so in the foreseeable future. New genotyping and phenotyping technologies, such
821 as near-infrared spectroscopy, are constantly being developed or improved. Breeding
822 databases must co-evolve with the technological advances to remain relevant, requiring
823 significant effort in refactoring and implementation. Systems that easily adapt to new
824 technologies will have a distinct advantage; in terms of software development
825 strategies, agile software development will be more efficient than older waterfall type
826 models. Another area of improvement is that of algorithms and other aspects of
827 methodology. With a strong connection to the R programming language, it is relatively
828 easy to implement new algorithms in Breedbase, as they often require little modification
829 from standalone scripts to work within Breedbase. At its core, Breedbase uses a
830 relational database with integrated JSON data storage, which provides a healthy
831 balance between highly structured, normalized data and flexibility. However, other
832 systems, such as graph databases and highly parallelized solutions like Hadoop, or a
833 combination thereof, are becoming popular and may be integrated into Breedbase in the
834 future.

835

836 All of the Breedbase code is open source and readily available on the code sharing site
837 GitHub (<https://github.com/solgenomics>).

838

839 Conclusions

840 Breedbase provides a fully open-source, scalable and feature-rich breeding digital
841 ecosystem that has been in use at the RTB crops breeding centers of the CGIAR for

842 many years, starting with the NextGen Cassava database, Cassavabase
843 (<https://cassavabase.org/>). The system has now been adopted by various breeding
844 programs including vegetable and grain crops and maintains an open and collaborative
845 approach to software development, allowing database customization for each research
846 community while sustaining a common framework. Our hope is that Breedbase, and the
847 digital ecosystem that it provides, can contribute, in a small way, to solving the world's
848 big problems with food scarcity and food quality, and thus contribute to improving
849 subsistence farmers' lives around the world.

850 Web Resources

851 <https://github.com/solgenomics/> Github repositories for Breedbase code

852 <https://hub.docker.com/r/breedbase/breedbase#> Docker image for Breedbase server

853 <https://breedbase.org/> Breedbase demo site

854 <https://cassavabase.org/> Cassavabase, the flagship Breedbase site

855 <https://musabase.org/> Breedbase site for banana breeding

856 <https://yambase.org/> Breedbase site for yam breeding

857 <https://sweetpotatobase.org/> Breedbase site for sweet potato breeding

858 <https://www.youtube.com/channel/UC3jrvvzGKKEHzOriDBgnj0A> YouTube channel for

859 Breedbase

860

861 Data availability statement

862 All code is available from Github (<https://github.com/solgenomics>) and docker hub

863 (<https://hub.docker.com/r/breedbase/breedbase#>).

864

865 Acknowledgments

866

867 Funding: This work was partially supported by the NEXTGEN Cassava project,
868 through a grant to Cornell University by the Bill & Melinda Gates Foundation
869 (Grant INV-007637 <http://www.gatesfoundation.org>) and the UK's Foreign,
870 Commonwealth & Development Office (FCDO). Other funding was provided by BMGF
871 through the Africa Yam project, Better Breeding Bananas, and the sweetpotato-focused
872 GT4SP and SweetGAINS projects.

873 JLJ, CLB, and DJW received partial support from the US Wheat and Barley Scab
874 Initiative and the National Research Initiative Competitive Grants 2017- 67007- 25939
875 and 2017- 67007- 25929 from the National Institute of Food and Agriculture, U.S.
876 Department of Agriculture.

877 The Crop Ontology is financially supported through the CGIAR Platform for Big Data in
878 Agriculture and the CGIAR Agrifood Research Programmes, by the CGIAR Trust Fund
879 (<https://www.cgiar.org/funders/>), and UKAID. The Crop Ontology was additionally
880 supported through the Planteome Project, led by Pankaj Jaiswal (Oregon State
881 University), by the National Science Foundation, USA (IOS:1340112).

882 Contributions: We would like to thank the Breeding Insight Project, including Moira
883 Sheehan, Tim Parsons, Nick Palladino, and Chris Tucker, for providing code to the
884 project. We would also like to thank many of the previous members of the breeding
885 leads and staff, including Racheal Mukisa, Dorcus Gemenet, Edward Carey, Jolien
886 Swanckaert, Jan Low, Robert Mwanga, and many others. Thanks to Thomas Hickey,
887 Keo Corak and Julie Dawson for their many bug reports and suggestions. For the
888 development of food quality related features, we would like to thank the RTBFoods
889 project, especially Karima Meghar, Thierry Tran and Dominique Dufour. Thanks to Lynn
890 Johnson, Chris Hernandez, and Peter Selby for their help and suggestions. Thanks to
891 Luka Wanjohi, Reinhard Simon, Ciro Rosales, Ivan Perez and Elisa Salas for comments
892 and suggestions. Special thanks to Kathy Kahn, Jim Lorenzen and the entire BMGF for

893 their tireless support of the African RTB breeding programs. This paper is dedicated to
 894 the memory of Martha Hamblin.

895

896 Conflict of Interest

897 None declared.

898 Literature Cited

- 899 Arnaud, Elizabeth, Marie-Angélique Laporte, Soonho Kim, Céline Aubert, Sabina Leonelli, Berta
 900 Miro, Laurel Cooper, et al. 2020. "The Ontologies Community of Practice: A CGIAR
 901 Initiative for Big Data in Agrifood Systems." *Patterns (New York, N.Y.)* 1 (7): 100105.
- 902 Beck, Kent, and Cynthia Andres. 2004. *Extreme Programming Explained: Embrace Change*.
 903 Pearson Education.
- 904 Bombarely, Aureliano, Naama Menda, Isaak Y. Teclé, Robert M. Buels, Susan Strickler,
 905 Thomas Fischer-York, Anuradha Pujar, Jonathan Leto, Joseph Gosselin, and Lukas A.
 906 Mueller. 2011. "The Sol Genomics Network (solgenomics.net): Growing Tomatoes Using
 907 Perl." *Nucleic Acids Research* 39 (Database issue): D1149–55.
- 908 Celko, Joe. 2006. "OLAP Basics." *Joe Celko's Analytics and OLAP in SQL*.
 909 <https://doi.org/10.1016/b978-012369512-3/50028-7>.
- 910 Cobb, Joshua N., Roselyne U. Juma, Partha S. Biswas, Juan D. Arbelaez, Jessica Rutkoski,
 911 Gary Atlin, Tom Hagen, Michael Quinn, and Eng Hwa Ng. 2019. "Enhancing the Rate of
 912 Genetic Gain in Public-Sector Plant Breeding Programs: Lessons from the Breeder's
 913 Equation." *TAG. Theoretical and Applied Genetics. Theoretische Und Angewandte Genetik*
 914 132 (3): 627–45.
- 915 Coombes, N. E. 2009. "DiGGeR, a Spatial Design Program." *Biometric Bulletin* NSW DPI.
- 916 Danecek, Petr, Adam Auton, Goncalo Abecasis, Cornelis A. Albers, Eric Banks, Mark A.
 917 DePristo, Robert E. Handsaker, et al. 2011. "The Variant Call Format and VCFtools."
 918 *Bioinformatics* 27 (15): 2156–58.
- 919 Duarte, J. B., and R. M. C. Pinto. 2002. "Biplot AMMI Graphic Representation of Specific
 920 Combining Ability." *Cropp Breeding and Applied Biotechnology*.
 921 <https://doi.org/10.12702/1984-7033.v02n02a01>.
- 922 Elshire, Robert J., Jeffrey C. Glaubitz, Qi Sun, Jesse A. Poland, Ken Kawamoto, Edward S.
 923 Buckler, and Sharon E. Mitchell. 2011. "A Robust, Simple Genotyping-by-Sequencing
 924 (GBS) Approach for High Diversity Species." *PLoS One* 6 (5): e19379.
- 925 Fernandez-Pozo, Noe, Naama Menda, Jeremy D. Edwards, Surya Saha, Isaak Y. Teclé, Susan
 926 R. Strickler, Aureliano Bombarely, et al. 2015. "The Sol Genomics Network (SGN)—from
 927 Genotype to Phenotype to Breeding." *Nucleic Acids Research*.
 928 <https://doi.org/10.1093/nar/gku1195>.
- 929 Fernandez-Pozo, Noe, Hernan G. Rosli, Gregory B. Martin, and Lukas A. Mueller. 2015. "The
 930 SGN VIGS Tool: User-Friendly Software to Design Virus-Induced Gene Silencing (VIGS)
 931 Constructs for Functional Genomics." *Molecular Plant* 8 (3): 486–88.
- 932 Food and Agriculture Organization of the United Nations. 2018. *Genebank Standards for Plant
 933 Genetic Resources for Food and Agriculture*. Food & Agriculture Org.
- 934 Hassani-Pak, Keywan, Ajit Singh, Marco Brandizi, Joseph Hearnshaw, Jeremy D. Parsons,
 935 Sandeep Amberkar, Andrew L. Phillips, John H. Doonan, and Chris Rawlings. 2021.
 936 "KnetMiner: A Comprehensive Approach for Supporting Evidence-Based Gene Discovery
 937 and Complex Trait Analysis across Species." *Plant Biotechnology Journal*, March.
 938 <https://doi.org/10.1111/pbi.13583>.
- 939 Hershberger, J., N. Morales, C. C. Simoes, and B. Ellerbrock. 2020. "Making WAVES in
 940 Breedbase: An Integrated Spectral Data Storage and Analysis Pipeline for Plant Breeding

- 941 Programs.” *bioRxiv*.
 942 <https://www.biorxiv.org/content/10.1101/2020.09.18.278549v1.abstract>.
- 943 Holland, James B., Wyman E. Nyquist, and Cuauhtemoc T. Cervantes-Martínez. 2010.
 944 “Estimating and Interpreting Heritability for Plant Breeding: An Update.” *Plant Breeding*
 945 *Reviews*. <https://doi.org/10.1002/9780470650202.ch2>.
- 946 “ISO/IEC TR 19075-6:2017.” 2018. ISO. 2018. <https://www.iso.org/standard/67367.html>.
- 947 Jung, Sook, Naama Menda, Seth Redmond, Robert M. Buels, Maren Friesen, Yuri Bendana,
 948 Lacey-Anne Sanderson, et al. 2011. “The Chado Natural Diversity Module: A New Generic
 949 Database Schema for Large-Scale Phenotyping and Genotyping Data.” *Database: The*
 950 *Journal of Biological Databases and Curation* 2011 (November): bar051.
- 951 Kilian, Andrzej, Peter Wenzl, Eric Huttner, Jason Carling, Ling Xia, Hélène Blois, Vanessa Caig,
 952 et al. 2012. “Diversity Arrays Technology: A Generic Genome Profiling Technology on
 953 Open Platforms.” *Methods in Molecular Biology* 888: 67–89.
- 954 Li, Fay-Wei, Paul Brouwer, Lorenzo Carretero-Paulet, Shifeng Cheng, Jan de Vries, Pierre-Marc
 955 Delaux, Ariana Eily, et al. 2018. “Fern Genomes Elucidate Land Plant Evolution and
 956 Cyanobacterial Symbioses.” *Nature Plants* 4 (7): 460–72.
- 957 Menda, Naama, Robert M. Buels, Isaak Teclé, and Lukas A. Mueller. 2008. “A Community-
 958 Based Annotation Framework for Linking Solanaceae Genomes with Phenomes.” *Plant*
 959 *Physiology* 147 (4): 1788–99.
- 960 Mendiburu, Felipe De, Felipe De Mendiburu, and Reinhard Simon. n.d. “Agricolae - Ten Years
 961 of an Open Source Statistical Tool for Experiments in Breeding, Agriculture and Biology.”
 962 <https://doi.org/10.7287/peerj.preprints.1404v1>.
- 963 Meuwissen, T. H., and M. E. Goddard. 2001. “Prediction of Identity by Descent Probabilities
 964 from Marker-Haplotypes.” *Genetics, Selection, Evolution: GSE* 33 (6): 605–34.
- 965 Morales, Nicolas, Guillaume J. Bauchet, Titima Tantikanjana, Adrian F. Powell, Bryan J.
 966 Ellerbrock, Isaak Y. Teclé, and Lukas A. Mueller. 2020. “High Density Genotype Storage for
 967 Plant Breeding in the Chado Schema of Breedbase.” *PLOS ONE*.
 968 <https://doi.org/10.1371/journal.pone.0240059>.
- 969 Morales, Nicolas, Nicholas S. Kaczmar, Nicholas Santantonio, Michael A. Gore, Lukas A.
 970 Mueller, and Kelly R. Robbins. 2020. “ImageBreed: Open- access Plant Breeding Web-
 971 database for Image- based Phenotyping.” *The Plant Phenome Journal*.
 972 <https://doi.org/10.1002/ppj2.20004>.
- 973 Mueller, Lukas A., Adri A. Mills, Beth Skwarecki, Robert M. Buels, Naama Menda, and Steven
 974 D. Tanksley. 2008. “The SGN Comparative Map Viewer.” *Bioinformatics* 24 (3): 422–23.
- 975 Mueller, Lukas A., Teri H. Solow, Nicolas Taylor, Beth Skwarecki, Robert Buels, John Binns,
 976 Chenwei Lin, et al. 2005. “The SOL Genomics Network: A Comparative Resource for
 977 Solanaceae Biology and beyond.” *Plant Physiology* 138 (3): 1310–17.
- 978 Mueller, Lukas A., Steven D. Tanksley, Jim J. Giovannoni, Joyce van Eck, Stephen Stack, Doil
 979 Choi, Byung Dong Kim, et al. 2005. “The Tomato Sequencing Project, the First Cornerstone
 980 of the International Solanaceae Project (SOL).” *Comparative and Functional Genomics* 6
 981 (3): 153–58.
- 982 Musen, Mark A. 2015. “The Protégé Project.” *AI Matters*.
 983 <https://doi.org/10.1145/2757001.2757003>.
- 984 Ogonna, Alex C., Luciano Rogerio Braatz de Andrade, Ismail Y. Rabbi, Lukas A. Mueller, Eder
 985 Jorge de Oliveira, and Guillaume J. Bauchet. 2020. “Large-Scale GWAS Using Historical
 986 Data Identifies a Conserved Genetic Architecture of Cyanogenic Glucosides Content in
 987 Cassava (*Manihot Esculenta* Crantz.) Root.” *The Plant Journal: For Cell and Molecular*
 988 *Biology*.
 989 https://onlinelibrary.wiley.com/doi/abs/10.1111/tpj.15071?casa_token=iT4r8kmf2eMAAAA:q0QJGeKM0Db-rA-SLnV48jK9UIFnA02xJxq5AWcBUzUwESSOJiYJteaC-kyhCPQZTrgcq_dIFWPY3zY.
- 990
 991
- 992 Pietragalla J., Valette L., Shrestha R., Laporte M.-A., Hazekamp T., Arnaud E., 2020.
 993 Guidelines for creating crop-specific ontologies to annotate phenotypic data, version 2.0,
 994 Alliance Bioversity International-CIAT, December 2020
- 995 Rife, Trevor W., and Jesse A. Poland. 2014. “Field Book: An Open-Source Application for Field
 996 Data Collection on Android.” *Crop Science*. <https://doi.org/10.2135/cropsci2013.08.0579>.

- 997 Saha, Surya, Prashant S. Hosmani, Krystal Villalobos-Ayala, Sherry Miller, Teresa Shippy,
998 Mirella Flores, Andrew Rosendale, et al. 2017. "Improved Annotation of the Insect Vector of
999 Citrus Greening Disease: Biocuration by a Diverse Genomics Community." *Database: The*
1000 *Journal of Biological Databases and Curation* 2017 (January).
1001 <https://doi.org/10.1093/database/bax032>.
- 1002 Selby, Peter, Rafael Abbeloos, Jan Erik Backlund, Martin Basterrechea Salido, Guillaume
1003 Bauchet, Omar E. Benites-Alfaro, Clay Birkett, et al. 2019. "BrAPI-an Application
1004 Programming Interface for Plant Breeding Applications." *Bioinformatics* 35 (20): 4147–55.
- 1005 Semagn, Kassa, Raman Babu, Sarah Hearne, and Michael Olsen. 2014. "Single Nucleotide
1006 Polymorphism Genotyping Using Kompetitive Allele Specific PCR (KASP): Overview of the
1007 Technology and Its Application in Crop Improvement." *Molecular Breeding*.
1008 <https://doi.org/10.1007/s11032-013-9917-x>.
- 1009 Shore, James, Chromatic, and Shane Warden. 2008. *The Art of Agile Development*. "O'Reilly
1010 Media, Inc."
- 1011 Shrestha, Rosemary, Luca Matteis, Milko Skofic, Arlet Portugal, Graham McLaren, Glenn
1012 Hyman, and Elizabeth Arnaud. 2012. "Bridging the Phenotypic and Genetic Data Useful for
1013 Integrated Breeding through a Data Annotation Using the Crop Ontology Developed by the
1014 Crop Communities of Practice." *Frontiers in Physiology* 3 (August): 326.
- 1015 Teclé, Isaak Y., Jeremy D. Edwards, Naama Menda, Chiedozi Egesi, Ismail Y. Rabbi, Peter
1016 Kulakow, Robert Kawuki, Jean-Luc Jannink, and Lukas A. Mueller. 2014. "solGS: A Web-
1017 Based Tool for Genomic Selection." *BMC Bioinformatics* 15 (December): 398.
- 1018 Teclé, Isaak Y., Naama Menda, Robert M. Buels, Esther van der Knaap, and Lukas A. Mueller.
1019 2010. "solQTL: A Tool for QTL Analysis, Visualization and Linking to Genomes at SGN
1020 Database." *BMC Bioinformatics* 11 (October): 525.
- 1021 Thomson, M. J. 2014. "High-Throughput SNP Genotyping to Accelerate Crop Improvement."
1022 *Plant Breeding and Biotechnology*. <https://www.e-sciencecentral.org/articles/SC000009999>.
- 1023 Tomato Genome Consortium. 2012. "The Tomato Genome Sequence Provides Insights into
1024 Fleshy Fruit Evolution." *Nature* 485 (7400): 635–41.
- 1025 VanRaden, P. M. 2008. "Efficient Methods to Compute Genomic Predictions." *Journal of Dairy*
1026 *Science* 91 (11): 4414–23.
- 1027 Züst, Tobias, Susan R. Strickler, Adrian F. Powell, Makenzie E. Mabry, Hong An, Mahdieh
1028 Mirzaei, Thomas York, et al. 2020. "Independent Evolution of Ancestral and Novel
1029 Defenses in a Genus of Toxic Plants (, Brassicaceae)." *eLife* 9 (April).
1030 <https://doi.org/10.7554/eLife.51712>.

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042 **Figure 1a:** Breedbase platform architecture.

1043 User interface:

1044 To offer a dynamic, highly interactive user interface, several Javascript libraries are
1045 implemented including D3, JQuery and Bootstrap. RESTful APIs, including a full BrAPI
1046 2.0 implementation, handle the communication between the front and back end,
1047 allowing fast calculations without reloading the website. HTML5 for interactive graphical
1048 display, allowing instant reorganisation of visual elements. The Bootstrap framework is
1049 used for modern and dynamic page templating.

1050

1051 Middleware layer:

1052 A Perl software stack including Mason components to connect to the user interface, a
1053 Catalyst a web application framework, Moose an object oriented perl library and
1054 DBIX::Class an object-relational mapper to connect to SQL code. In addition, BrAPI
1055 libraries are used. Finally a job cluster scheduler, Slurm is implemented to allocate
1056 server resources and ensure scalability.

1057

1058 Data source layer:

1059 Breedbase operates on a relational database using Postgres. Postgres 12.0 offers “Big
1060 data” solutions including parallel query execution and optimized binary javascript object
1061 notation data type (JSON) handling. Binary JSON (JSONB) is a simple data structure
1062 designed to be storage space and scan-speed efficient. In Breedbase, JSONB is used
1063 in various data types including genotypic (marker) information. In addition to the
1064 relational database a standard file system space is available for flat files. Finally, other
1065 databases can communicate to a Breedbase instance to provide additional back-end for
1066 marker data (ie: Genomic Open Source Informatic Initiative (GOBii)) or to exchange
1067 germplasm information for example.

1068

1069 **Figure 1b:** Breedbase co-development process.

1070 User-developers interactions are promoted using various media. Users have online
1071 access to documentation (<https://solgenomics.github.io/sgn/>), video tutorials or through
1072 onsite training. Software development goals are extensively discussed between
1073 developers, data managers, breeders and other appropriate stakeholders. Agile
1074 development allows short term product release. Suggested improvements, issues and
1075 bugs discovered in Breedbase are submitted and tracked on the public GitHub issue
1076 tracking software (<https://github.com/>). Software development progress is tracked using
1077 a version control system and Docker releases.

1078

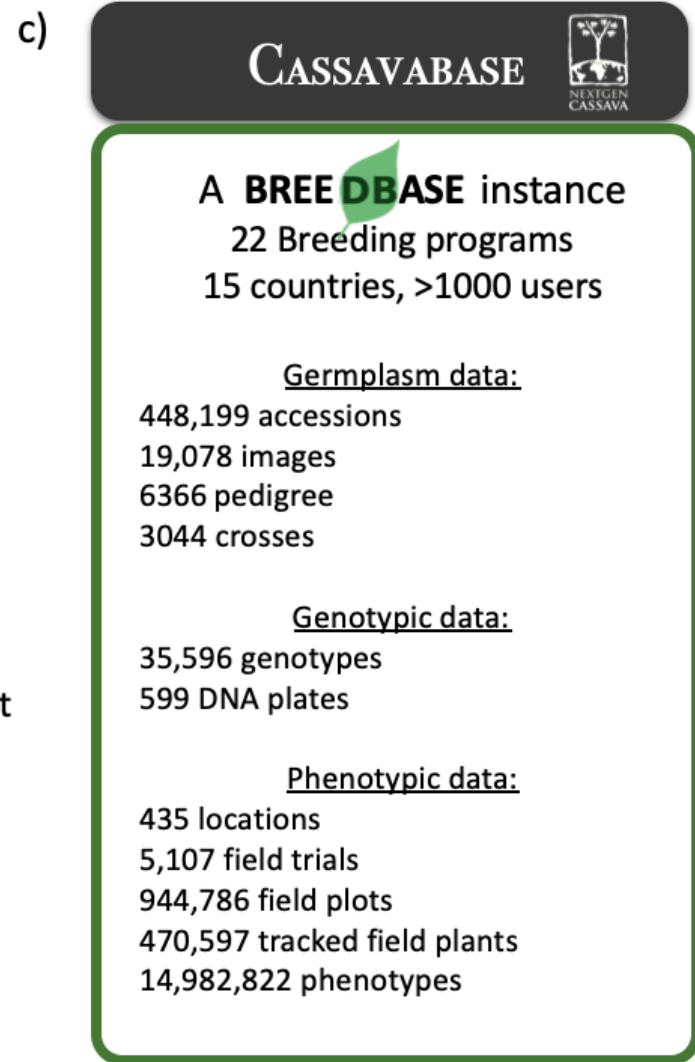
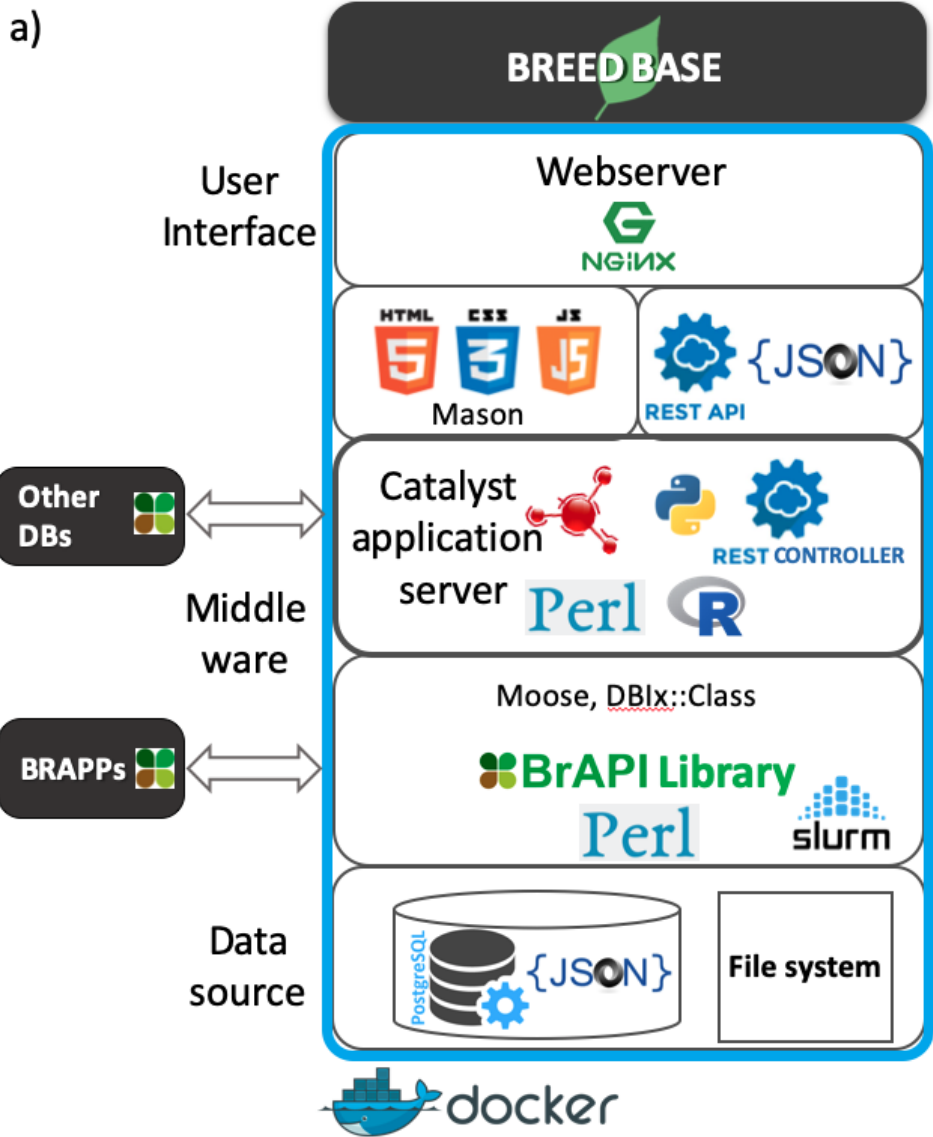
1079 **Figure 1c:** Cassavabase, a breedbase instance: data content overview.

1080 Cassavabase involves national and international breeding programs (22) from various
1081 African and South American countries (15) and currently has 1131 registered users.
1082 Cassavabase hosts various data types including high density and low density
1083 genotyping assays (35,000), plot based phenotypic data points (near 15 million), images
1084 from plants and plots from, trials (5107) and locations (435).

1085

1086

1087 Figure 2: Screenshot of the “Search Wizard” interface, a central query function on
1088 Breedbase. With the Search Wizard, the data in the database can be intersected by
1089 dimensions, such as locations, years, breeding programs, and traits. For each
1090 dimension, a number of elements can be selected. The individual selected dimensions
1091 can be stored in lists, and the combined selections can be saved as a dataset. Both lists
1092 and datasets can be used to feed data into various tools on Breedbase.



Search Wizard

Don't see your data?

Refresh Lists

Update Wizard

Years

Search

Select All 2/41 Clear

+ 2013
+ 2014
+ 2015
+ 2016
+ 2017
× 2018
× 2019

Match ANY ALL

Add to List... Add

Create New List... Create

Locations

Search

Select All 3/59 Clear

+ Ukerewe
+ Ukiriguru
+ UNIVASF Trial
+ unspecified
+ Wench-BA
× Umudike
× Uyo
× Zaria

Match ANY ALL

Add to List... Add

Create New List... Create

Traits

Search

Select All 4/74 Clear

+ cassava bacterial blight incidence
+ cassava bacterial blight severity
+ cassava brown streak disease
+ cassava green mite incidence
+ cassava green mite incidence
× cassava bacterial blight incidence
× cassava bacterial blight incidence
× cassava bacterial blight severity
× cassava bacterial blight severity

Match ANY ALL

Add to List... Add

Create New List... Create

Accessions

Search

Select All 0/1617 Clear

+ 50395
+ AR1410
+ AR144
+ AR1-82
+ AR311

Load/Create Datasets using Match Columns

Load Dataset

Load

Create New Dataset

Create

Related Genotype Data

Related Trial Metadata

Related Trial Phenotypes