

Lawrence Berkeley National Laboratory

LBL Publications

Title

A primer on artificial intelligence in plant digital phenomics: embarking on the data to insights journey

Permalink

<https://escholarship.org/uc/item/6kt1d51f>

Journal

Trends in Plant Science, 28(2)

ISSN

1360-1385

Authors

Harfouche, Antoine L

Nakhle, Farid

Harfouche, Antoine H

et al.

Publication Date

2023-02-01

DOI

10.1016/j.tplants.2022.08.021

Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

1 A primer on artificial intelligence in plant digital phenomics:
2 embarking on the data to insights journey

3 Antoine L. Harfouche,^{1,*} Farid Nakhle,¹ Antoine H. Harfouche,² Orlando G.
4 Sardella,¹ Eli Dart,³ and Daniel Jacobson⁴

5 ¹Department for Innovation in Biological, Agro-food and Forest systems, University of
6 Tuscia, Viterbo, VT 01100, Italy

7 ²Unité de Formation et de Recherche en Sciences Économiques, Gestion,
8 Mathématiques et Informatique, Université Paris Nanterre, 92001 Nanterre, France

9 ³ESnet, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

10 ⁴Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

11 *Correspondence: aharfouche@unitus.it (A. L. Harfouche)

12 **Keywords**

13 AI system architecture, black box models, data analytics, digital phenomics, explainable
14 artificial intelligence, interpretable by design models

15

16 **Artificial intelligence (AI) has emerged as a fundamental component of global**
17 **agricultural research poised to impact many aspects of plant science. In digital**
18 **phenomics, AI is capable of learning intricate structure and patterns in large**
19 **datasets. Here, we provide a perspective and primer on AI applications to**
20 **phenome research. We propose a novel human-centric explainable AI (X-AI)**
21 **system architecture, consisting of data architecture, technology infrastructure,**
22 **and AI architecture design. We clarify the difference between post-hoc models**
23 **and interpretable by design models. We include guidance for effectively using an**
24 **interpretable by design model in phenomics analysis. We also provide a direction**
25 **to sources of tools and resources for making data analytics increasingly**
26 **accessible. This primer is accompanied by an interactive online tutorial.**

27 **Approaching plant phenomics from different angles**

28 Crop breeding relies heavily on phenotypic information, which remains a bottleneck for
29 realizing its full potential. The advent of plant phenomics (topic reviewed in [1–9]), which
30 broadly can be considered the systematic study of phenotypes, however, marks a
31 turning point. Phenomics platforms equipped with novel imaging sensors promise to
32 make it possible to perform phenotyping of a wide range of plant traits, organs, and
33 environmental situations at scale (Figure 1). These new technological developments
34 have opened up avenues for automated data acquisition, evolving phenomics to
35 become a thriving research field of its own [10].

36 Embracing digital technology holds tremendous potential for driving transformative
37 changes in plant phenomics by improving the collection of, access to, and analysis of

38 phenomic big data. As such, digital phenomics furnishes tools and resources that aid in
39 the digitization of plant phenomics. It uses phenome data and metadata to guide
40 decision-making along the entire data analytics cycle [11]. As data grows in volume and
41 varies in sources, effective management strategies must be put into place (Figure 2).
42 Digital phenomics gives rise to computational phenomics, which allows the assembly of
43 a broad array of methods that aid in the discovery of intricate structure and patterns
44 from phenotypic data, using technology infrastructure (Figure 3) and artificial
45 intelligence (AI) architecture design (Figure 4).

46 In the past decade, AI – the science of studying, designing, and developing intelligent
47 computer systems that can perform tasks that normally require human intelligence –
48 finally began to reveal its remarkable power and disruptive potential. Driven mainly by
49 the advent of machine learning (ML) – a particular approach to AI in which intelligent
50 systems learn and derive models from training datasets – and deep learning (DL) – a
51 specialized branch of ML that leverages neural networks to spot patterns in complex
52 data – AI flexed its muscles by achieving predictive successes in phenomics. For
53 example, in red–green–blue (RGB) image analysis, convolutional neural networks
54 (CNNs) were used to predict the yield of individual plants of barley and wheat [12], to
55 classify and quantify biotic and abiotic stresses in leaves of various fruit and vegetable
56 crops [13], to segment roots of chicory, wheat, and rapeseed [14–16], and to count
57 tobacco leaves [17]. A CNN was also used to provide explainable classifications of
58 biotic and abiotic stresses in soybean leaves by isolating the top-k feature maps learned
59 by the model [18]. A random forest (RF), a neural network (NN), a k-nearest neighbor
60 (KNN), a partial least squares (PLS), and a support vector machine (SVM) were

61 employed to estimate nitrogen nutrition index for improving nitrogen use efficiency in
62 rice [19,20]. In wood anatomical images at a microscopic level, a mask region-based
63 CNN (Mask R-CNN) was used to analyze the intrinsic variability of wood anatomical
64 features in conifers, alder, beech, and oak [21]. In analyzing multispectral and
65 hyperspectral data, a CNN was used to provide explainable identification of biotic stress
66 in individual soybean plants by incorporating saliency maps [22]; a PLS and a RF to
67 estimate above ground biomass in maize [23]; an SVM, a KNN, and a linear
68 discriminant analysis (LDA) to detect and segment root decay in wheat [24]; an LDA
69 and a PLS to detect response to drought stress in bell pepper, courgette, sunflower,
70 radish, foxtail millet, and sorghum [25]; and a data mining sharpener to guarantee
71 consistent spatial resolution among heterogeneous remote sensing image datasets to
72 dissect the latent heat flux signature of poplar in response to drought [26]. In three-
73 dimensional (3D) point cloud analysis, an SVM was applied to estimate yield and
74 canopy geometric characterization in apple [27]. In thermal infrared (TIR) image
75 analysis, SVMs and Gaussian processes were used to identify drought stress in spinach
76 [28] and rotation forests were used to predict plant water status in grapevine [29]. In
77 chlorophyll fluorescence image analysis, a CNN was used to identify abnormalities in
78 organelle morphology in Arabidopsis [30]. In X-ray computed tomography (X-ray CT)
79 analysis, an encoder-decoder network was used to segment wheat roots [31].

80 Analyzing data coming from different sensors and imaging techniques of the same
81 biological sample (i.e., plant) simultaneously can further improve phenotypic trait
82 predictions. Recent studies demonstrated that fusion of multiple data sources
83 originating from the same plants (i.e., paired data) perform better than a single source.

84 For example, using a deep NN (DNN), the fusion of RGB, TIR, and multispectral data
85 delivered superior performance over single sensor data analytics for yield prediction in
86 soybean [32]; similar results were reported for the estimation of soybean chlorophyll
87 content, nitrogen concentration, leaf area index, and above ground biomass by
88 employing an extreme learning machine [33]. Now and in the future, rigorous data
89 integration of phenomics and other different omics datasets that were not originally set
90 out to be integrated and are of distinct biological samples (i.e., unpaired data) may help
91 dissecting biological mechanisms that underlie desirable traits and shed light on the flow
92 of information that underpins plant responses to environmental stresses [10,34–37].
93 And because the phenotype of a plant is the result of interaction between its genotype
94 and the environment (G x E) in which it grows [38], integration efforts should also
95 include environmental data such as climatypes; this will be crucial for designing new
96 crop ideotypes that are optimized for niche environments in a world with a rapidly
97 changing climate [34].

98 Importantly, data management strategies should incorporate the findable, accessible,
99 interoperable, and reusable (FAIR) guiding principles [39] to put those phenome and
100 envirome data to their most effective use. This requires standards to ensure that
101 necessary metadata are recorded about data generation methods and the experimental
102 and environmental conditions in which they were acquired [40]. In this regard, the
103 minimum information about a plant phenotyping experiment (MIAPPE) standard has
104 been a great step forward to harmonize data from phenotyping experiments with
105 controlled vocabulary and ontologies [41,42]. Accordingly, the development of tools for
106 capturing the complete set of metadata is poised to have high impact on the support of

107 FAIR data. As standards and tools become more widely disseminated and explored, we
108 envisage metadata becoming commonly annotated by users, expected by referees, and
109 required by journals and data repositories. Furthermore, the workflows that are used to
110 analyze data should themselves be FAIR [43]. When data sharing is not viable due to
111 possible privacy or security concerns, **federated learning** (FL; see Glossary) and
112 communication-efficient FL offer an unprecedented opportunity to train AI models
113 without sharing data [44,45]. FL gained traction in medical imaging applications [46–48]
114 and carries great promise for overcoming data sharing challenges in plant phenomics.

115 This primer provides suggestions on how to use AI effectively in plant phenomics, on
116 how to ensure that human-centric explainable AI (X-AI) can benefit all, and discusses
117 various X-AI approaches and techniques. We have created a central directory of all
118 publicly available plant imaging datasets, and report their sources, accessibility, and a
119 summary of species and organ systems represented (Table 1). This review is
120 accompanied by an interactive tutorial to train an interpretable by design model to
121 deliver predictive and prescriptive analytics to users. Our primer is intended as an
122 educational resource for phenomicists who are interested in applying X-AI approaches
123 and techniques, and plant scientists who seek a high-level understanding of this rapidly
124 evolving field. Data scientists and information systems (IS) scientists may also use this
125 primer as an introduction to the promising applications of X-AI in phenomics.

126 **How to use AI effectively: ménage-à-trois between plant science, data science,**
127 **and IS**

128 Plant science has the potential to provide innovative solutions for the world's most
129 pressing challenges; however, recent advances in discovery methods have greatly
130 accelerated our ability to collect data, leaving us with the challenge of analyzing,
131 interpreting, and integrating the plethora of data [49]. To handle such data, data science
132 has attracted a lot of attention, promising to turn data into useful predictions and insights
133 [50]. To do that, data science needs supporting resources including algorithms,
134 software, and hardware infrastructure. IS combines those resources to create AI
135 architecture designs, and to transform, store, and distribute data for analysis. While the
136 relationship between these disciplines has not been reinforced repeatedly in history,
137 today with the depth of data analysis, the scale and dimension of the data, and the
138 nature of the scientific questions, an interaction in a ménage-à-trois fashion is highly
139 needed.

140 AI and the bias cascade

141 Multiple sources of **bias** can affect the performance of AI systems used in phenomics
142 and can occur across the different development steps of AI applications: data collection
143 or selection, data preprocessing, model development, model evaluation, and
144 deployment. Introduced bias can have a domino effect as it propagates from its entry
145 point to the succeeding development steps, creating a bias cascade.

146 The bias cascade starts with the data collection or selection step, where experimental
147 data are collected or selected from publicly available datasets (Table 1). Here, bias can
148 occur for a number of reasons: (i) 'measurement' bias, when data contains faulty
149 measurements originating from instrumentation malfunctions, wrong values from

150 miscalibrated sensors, or errors of precision that result in data distortion [51]; (ii) 'label'
151 bias, when data is laden with subjective judgments of human experts and thus
152 inconsistently or wrongly labeled [52]; (iii) 'sample selection' bias, when the training data
153 does not represent a random sample from the entire dataset [53], causing a model to
154 ignore data belonging to classes that were not represented during the data selection
155 process; and (iv) 'group attribution' bias, when a data sample is selected from an
156 incorrect target population [53], where a model can fail to distinguish between some
157 classes and consider them the same.

158 Data preprocessing is performed to eliminate noisy (e.g., blurred images, images with
159 unfavorable lighting conditions, images that do not represent the object of interest),
160 incomplete (e.g., unannotated images), duplicate data, and to normalize datasets as
161 needed to account for batch effects (e.g., groups of images taken under different
162 lighting conditions or with different camera settings) or systematic experimental artifacts
163 (e.g., reflections in images). In this step, even if the training set was representative of
164 the entire dataset, data can be intrinsically **unbalanced** where certain plant species,
165 genotypes, or even stresses are underrepresented. Such cases can introduce the 'class
166 imbalance' bias.

167 Bias may also arise during the model development and evaluation steps, where a model
168 is trained and its ability to generalize beyond the training set, on new, previously unseen
169 data is evaluated. As most AI algorithms identify correlations between variables in the
170 underlying data but without being able to detect causal relations, two biases are likely to
171 arise: (i) the 'correlation fallacy' that confuses correlation with causation [53] where a

172 model wrongly deduces a cause-and-effect relationship between correlated variables;
173 and (ii) the ‘apophenia’ when a model sees patterns while none actually exist [54].
174 These two biases can be amplified when a massive quantity of training data is used,
175 mistakenly offering connections that radiate in all directions [54], and producing
176 probable yet uncertain predictions. Further, training complex models (i.e., models with
177 many trainable parameters) can capture noise-generated patterns, tricking them into
178 thinking that the noise encodes real information [55]. This problem introduces the
179 ‘overfitting’ bias and causes a steep drop-off in predictive performance at the evaluation
180 step. Similarly, such performance drop-offs also occur when models are unable to
181 accurately capture relationships between variables and thus introducing the
182 ‘underfitting’ bias [56].

183 Finally, at the deployment step, bias can occur in situations where data used in practice
184 differs from training data (e.g., different weed or crop species), which is known as the
185 ‘domain shift’ bias.

186 Creating a human-centric X-AI

187 It is therefore crucial to mitigate bias to increase the success probability of the AI
188 algorithm for the task at hand. Let alone that bias mitigation serves as a building block
189 towards AI **trustworthiness** [57]. So, what can be done to mitigate detrimental biases
190 in AI in plant phenomics? There is a consensus on the need to develop a human-centric
191 X-AI system that will not just aspire to meet human requirements regarding
192 explainability and trustworthiness, but, more importantly, will actively aim to keep a
193 **human-in-the-loop** (HITL) for a harmonious human and AI system symbiosis. We

194 believe that such a system should not only put humans at its center, but also integrate
195 their knowledge into its predictive process.

196 Designing human-centric X-AI for plant phenomics is not without challenges; it requires
197 a dedicated and multidisciplinary team effort, involving plant scientists, data scientists,
198 and IS scientists to bring AI to its most feasible, desirable, viable, and responsible state.
199 This novel multidisciplinary knowledge is clearly imperative to identify and reduce AI
200 biases, and to facilitate explainability and accountability. We advocate that such a
201 system architecture is required to constantly realign data architecture and technology
202 infrastructure to serve novel AI architecture designs. Phenotyping complex traits
203 demands the integration of data on different morphological, physiological, temporal,
204 geospatial, and environmental variables [35,58]. While large datasets are vital for
205 creating accurate AI models and validating their results, storing them in a FAIR manner
206 can be challenging. Data architecture plays a fundamental role in meeting these
207 requirements. It consists of a set of standards that govern which data is collected,
208 whether it should be transformed (e.g., data cleaning, deduplication, format conversion,
209 structuring, validation, etc.) before or after storage using extract, transform, load (ETL)
210 or extract, load, transform (ELT) processes, and where (data warehouses or data lakes)
211 and how (matrices, cubes, polytopes, or distributed in-memory) it is stored (Figure 2).
212 Without AI, these data streams would be overwhelming and chaotic [35], but reaching
213 the full potential of AI-based analysis of large phenomic datasets comes down to the
214 right technology infrastructure which defines the components that serve as a foundation
215 for the data life cycle, including hardware infrastructure, network flow, software
216 frameworks, and programming languages (Figure 3). High performance computing

217 (HPC), like pre-exascale supercomputers, is boosting both the accuracy and predictive
218 power of these approaches. While central processing units (CPUs) maximize the
219 performance of an algorithm, graphics processing units (GPUs) can dramatically
220 increase AI training speed thanks to their processing cores initially designed to process
221 visual data such as images and videos [11]. For example, to take advantage of GPUs,
222 the compute unified device architecture (CUDA) software framework provides a
223 development environment for creating and optimizing AI applications on **GPU-**
224 **accelerated** local computers or supercomputers. However, CUDA works exclusively on
225 Nvidia GPUs; alternatively, the open computing language (OpenCL) and openACC
226 frameworks work on multiple types of GPUs [59]. Another option is to translate
227 automatically CUDA source code into portable heterogeneous-computing interface for
228 portability (HIP) using source-to-source translators such as HIPify, so that non-Nvidia
229 GPUs can benefit from the rapid development of CUDA applications. Additionally,
230 software libraries such as kokkos, RAJA, open multi-processing (OpenMP), and one
231 application programming interface (oneAPI) can be leveraged to unlock the promise of
232 heterogeneous computing where **compute nodes** employ more than one type of
233 processors including CPUs, GPUs, and tensor processing units (TPUs), among others.
234 This enables the development of scalable AI-based applications in a hardware agnostic
235 way. With the advent of exascale computing, supercomputers will deliver higher
236 performance in pattern searching in phenomic big data, and thus, will boost AI abilities
237 in digital phenomics, speeding up crop design (Figure 3A). But, building powerful
238 supercomputers is a never-ending race, and as new ones get launched, the number of
239 compute nodes they comprise increases. For example, the first supercomputer to break

240 the exascale barrier, Summit, comprises 4,608 compute nodes, while the most powerful
241 exascale supercomputer that tops the latest TOP500 list¹, Frontier, contains 9,472. This
242 makes it harder to exploit supercomputers efficiently because of their need to transmit
243 data back and forth between their nodes, running huge numbers of computations at the
244 same time [60]. Implementing AI algorithms (Figure 4A) for such **parallel computing** is
245 not easy. Luckily, emerging free and open-source software frameworks such as
246 Tensorflow Keras, PyTorch, scikit-learn, and XGBoost, among others, and software
247 libraries such as cuNumeric are enabling scalability on parallel computing. As more
248 powerful exascale supercomputers are being anticipated [61], researchers may start to
249 utilize quantum computers at some point in the future [62]. This will ultimately drive
250 digital phenomics towards designing faster, better crops and providing sustainability-
251 friendly solutions (Figure 3A). Beside the hardware infrastructure, properly designed
252 network flows (Figure 3B), such as the 'science demilitarized zone (DMZ)' that includes
253 network architecture and performance tools [63], enable high-throughput access to
254 datasets in a secure and timely manner while conforming with the FAIR data principles
255 [39]. Software frameworks (Figure 3C) provide a working environment that helps
256 researchers achieve higher productivity in designing AI algorithms; they support more
257 than one programming language (Figure 3D), enabling fast and efficient implementation
258 of algorithms without compromising code quality.

259 Because data are only as good as the tools and algorithms available to analyze them,
260 solving complex biological questions requires a creative process during which efficient
261 AI algorithm architectures are designed and developed. Customized algorithms and
262 architectures can leverage currently available AI architecture designs (Figure 4) to come

263 up with new architecture designs tailor-made to find answers to the questions at hand.
264 Such promising designs should combine knowledge-based AI, to represent human
265 expert knowledge, with data-driven AI to discover connections and correlations
266 automatically in big data. This combination will result in an informed AI that acquires
267 both **tacit** and **explicit knowledge** of its designers (e.g., the interaction between data
268 scientists, IS scientists, and plant scientists) and users (e.g., breeders and farmers),
269 and integrates that tacit and explicit knowledge with knowledge discovered from data
270 and metadata.

271 It is noteworthy, however, that new architecture designs should also integrate
272 knowledge into X-AI to enable the monitoring of the inputs and outputs of the
273 algorithms, provide more human-comprehensible explanations for their decisions,
274 deliver superior performance, mitigate bias, and aid in verifying models' adherence to
275 ethical and socio-legal values. Ensemble methods can, for example, be leveraged to
276 design new AI algorithms that are both informed and explainable (Figure 4B).
277 Ultimately, improvements in informed X-AI would help develop novel interpretable
278 algorithms and are likely to be crucial to enable human-centric X-AI in phenomics.

279 Mitigating bias in human-centric X-AI

280 A human-centric X-AI system is emerging, whereby plant scientists, data scientists, and
281 IS scientists must work together to seize this opportunity to help identify and mitigate
282 bias by using a number of strategies.

283 Starting from the top of the bias cascade, at the data collection or selection step, the
284 minority of data that do not conform to the general characteristics of a given dataset,

285 known as outliers, should be removed during data cleaning to mitigate the
286 'measurement' bias. As for the 'label' bias, data annotators should be supplied with
287 detailed instructions containing visual examples of the correct output for a given input to
288 be able to reduce ambiguities and avoid mistakes that result from incorrect or
289 incomplete knowledge. For example, when labeling weed species, in addition to their
290 morphological descriptions, a visual representation of each species could be helpful for
291 annotators. Next, the 'sample selection' and 'group attribution' biases can be mitigated
292 by establishing random sample selection and statistical correction processes [64,65].

293 When preprocessing data, intrinsically unbalanced datasets can be balanced by means
294 of oversampling (i.e., augmenting the number of training examples within the minority
295 class to be equivalent to other classes) and/or undersampling (i.e., reducing the number
296 of training examples within the majority class to be equivalent to other classes) [11] to
297 eliminate the 'class imbalance' bias [66].

298 Properly sampled and preprocessed data mitigate the risk of 'correlation fallacy' and
299 'apophenia' biases, which can occur during the model development step. When training,
300 'overfitting' can be debiased by either increasing the size of training data, decreasing
301 the model complexity, or ignoring the less important features in a process called
302 regularization [67]. Whereas 'underfitting' bias can be resolved by increasing the
303 complexity of the model to capture nonlinear relationships in data. During the evaluation
304 step, models yielding incorrect predictions such as misclassifying crops as weeds or
305 vice versa should be inspected carefully; X-AI can be leveraged to better understand
306 how the model reached its predictions which helps identify previously unknown bias.
307 However, post-hoc approaches to explainability are not necessarily transparent (i.e.,

308 because they only approximate models' prediction procedure), and thus, it might be
309 better to employ interpretable by design models (see next section).

310 Notably, to avoid the 'domain shift' bias and identify unwanted biases in deployment, it
311 is crucial that the model is carefully monitored to assess whether the data being used in
312 practice are representative of those used during training. It is important to note that risk-
313 based regulations of AI are on the horizon in the USⁱⁱ and Europe^{iii,iv}. When new
314 regulations enter into force, post-authorization monitoring of AI applications becomes
315 crucial to ensure that the performance of models does not degrade in practice. Once a
316 model has passed regulatory authorization and is implemented in phenomics, it needs
317 to be retrained periodically using new datasets to prevent it from becoming outdated,
318 ensuring 'domain shift' bias mitigation.

319 Reducing the risk of bias in AI models requires continuous human attention across the
320 five development steps, keeping HITL. Studies have shown that human-computer
321 interaction in HITL AI has improved the predictive performance of AI-based image
322 analysis and reduced biases [68–70]. HITL can make a significant impact in phenomic
323 data collection, data preprocessing, model development, evaluation, and deployment. It
324 plays a critical role in the collection and preparation of data to be used for training an AI
325 model. As such, model training is often a HITL iterative process that identifies biases or
326 weaknesses of the model (e.g., images on which the model fails due to incomplete
327 training sets or inappropriate parameterization) and adjusts the training set and
328 parameters to reduce any biases and ensure the best model performance. It is
329 recommended to start each training step with small iterations and plan on how the
330 feedback of the team of humans can be collected and propagated to other steps, relying

331 on their intelligence to perform complex tasks. This paradigm allows leveraging the
332 advantages of AI while having humans at various checkpoints to fill gaps where models
333 are not confident in their predictions or where they may fall short due to underlying
334 biases [71]. HITL may also offer advantages to evaluating the accuracy of AI predictions
335 and interpreting their decisions by interacting with explainable models. The benefits of
336 HITL extend to deployment by monitoring the model for possible biases and ensuring
337 the reliability of the AI system. HITL can feedback into itself to respond to changes in
338 the real-world environment. For example, after data collection and preprocessing, in
339 each training iteration, plant scientists are shown a list of misclassified images with the
340 outputs of the AI algorithm to hand-verify predictions and assess false positives and
341 false negatives. For instance, the model might misclassify crops as weed; but this could
342 be due to an algorithmic or learning error, or to mislabeled images. They then correct
343 the wrong labels, if any, to ensure high-quality data. Data scientists evaluate the model,
344 tune its **hyperparameters**, and retrain it. Such iterations between humans and AI are
345 effective to generate training data based on human judgment to increase learning
346 efficiency and enhance model performance [71]. Data scientists can provide the
347 expertise necessary to help IS scientists design AI architectures with explanatory
348 capacity supported by theoretical underpinnings. Finally, HITL monitors the model
349 outcomes post deployment to ensure that all biases are identified and mitigated.
350 Furthermore, ensembles of models can be used in the HITL process. An intrinsically
351 interpretable model such as an iterative RF [72,73], can be used initially for feature
352 engineering to determine the variables (e.g., wavelet decomposition in RGB and
353 hyperspectral images) that will then be used in a deep learning predictive model.

354 This three-way collaboration can amplify knowledge about domain-specific feature
355 engineering and selection to reach a level of augmented intelligence that can help
356 discovering new ways to make AI more efficient, less biased, and explainable. It also
357 creates new opportunities for human-centric X-AI to predict desirable phenotypic traits
358 and aid efforts to breed climate-proof crops fast enough.

359 **How to move from data inputs to outcomes: opening the black box or designing a**
360 **transparent glass box for explainability**

361 AI continues to permeate plant phenomics as recently reviewed in [74–76]. However,
362 complex AI models are difficult to explain even among data scientists; they operate as
363 black boxes and require a leap of faith to believe their predictions [35]. Explainability of
364 AI models would not only increase the trust of users in why and how predictions were
365 made but also help data scientists enable better diagnostics and enhance their
366 performance. Although these desirable properties of explainability have led to a recent
367 growing interest in X-AI research [77], its origin traces back to the early 1970s when
368 Edward Shortliffe introduced the AI-based antimicrobial therapy consultation system for
369 assisting physicians who need advice about appropriate therapy. The system made use
370 of a set of decision rules coded, categorized, and hand-entered into it to give advice and
371 explain the reasons behind its predictions [78]. In 1979, Jon Doyle introduced the truth
372 maintenance systems (TMS), an independent module that constructs explanations of
373 predictions by recording and maintaining a representation of the knowledge acquired by
374 an expert system [79]. TMS research and development continued until the 1990s, when
375 researchers began to study the possibility of extracting meaningful explanations from
376 non-hand-coded rules that are generated by trained models such as NN [80].

377 The rise of DL in the 2010s [11] increased the complexity of AI models and
378 consequently, the demand for X-AI algorithms. To address this issue, researchers have
379 been developing new approaches and techniques to make these models explainable.
380 Unfortunately, the rush in X-AI development has caused confusion on its various
381 approaches in the literature, where they are not accurately described and are often
382 confused together [81]. While all those approaches revolve around allowing humans to
383 observe how predictions of an AI model came to be, we can technically distinguish
384 between research involving post-hoc models and interpretable by design models.

385 As current AI models are often developed with only predictive performance in mind,
386 post-hoc algorithms can be used to explain them. They are employed after a black box
387 model is trained and are not connected to its internal design; they can either be model-
388 specific or model-agnostic [82]. In principle, model-specific algorithms are limited to
389 certain black box models. For example, DL important features (DeepLift) is a model-
390 specific algorithm that can explain DNNs and does not work for any other algorithm. On
391 the other hand, model-agnostic algorithms such as the local interpretable model-
392 agnostic explanations (LIME) [83] and Shapley additive explanation (SHAP) [84] are
393 more general and can be applied to any black box model. Commonly, post-hoc
394 algorithms work by: (i) probing or inspecting the trained parameters to understand what
395 has the black box model learned; (ii) employing data perturbation strategies which
396 involve modifying the input data and observing the changes in the black box model
397 predictions; or (iii) using a more interpretable model (e.g., decision tree) referred to as a
398 surrogate model to approximate and provide explanations of predictions made by the
399 black box model. Recently, researchers started applying post-hoc algorithms in plant

400 phenomics to identify, classify, and quantify plant stresses [11,85–89] and to count
401 leaves [90]. However, as post-hoc algorithms approximate the inner workings of black
402 box models, it is possible that their generated explanations do not provide enough detail
403 to understand what the black box model is actually doing [91]. On the other hand,
404 interpretable by design algorithms do not need an additional (post-hoc) algorithm to be
405 explainable; they provide their own explanations, which are faithful to what the model
406 actually computes [91].

407 These algorithms have existed since the development of expert systems in the 1970s.
408 They have, however, been labeled as less accurate because scientists argue that there
409 is a tradeoff between accuracy and explainability in a way that, the highest performing
410 algorithms are the least explainable, and the most explainable ones are less accurate^v.
411 This belief proved to be imprecise, especially when analyzing structured data with
412 meaningful features [91]. This also depends on the algorithms being compared. For
413 example, according to [91] it would not be fair to compare the 1984 decision tree
414 algorithm to a more recent DL one and conclude that interpretable by design models are
415 not as accurate. Indeed, the recently developed interpretable by design ‘this looks like
416 that’ algorithm, derived from a CNN, proved to be as accurate as the non-explainable
417 CNN [92]. Figure 5 highlights the two categories of X-AI, their corresponding
418 representative algorithms, and the explainable outcomes associated with their
419 implementation.

420 Regardless of whether a post-hoc or an interpretable by design algorithm is used,
421 model explanations can occur on a global or local level. While the former describes the

422 overall extracted relationships based on the entire model behavior, the latter reveals the
423 rationale behind a specific prediction [93].

424 Finally, it is worth noting that, just as different X-AI techniques exist, there exists a range
425 of approaches to explainability since different contexts give rise to different explainability
426 needs [94]. For example, when training and evaluating an AI model, plant scientists
427 might want to understand which data features are being used for prediction and how
428 they are correlated together, while data scientists might require technical details about
429 how the model functions to help in its testing, debugging, bias identification and
430 mitigation, hyperparameter tuning, and evaluation; IS scientists can leverage details
431 about the model training process to help optimize the architecture of the algorithm using
432 suitable design approaches and methods (Figure 4B). At the model deployment step,
433 regulators might require assurance about how data is being processed to assess its risk
434 level by inspecting its reliability, as well as the impact of its predictions on its users to
435 ultimately decide whether or not it requires authorization and regulation. Similarly,
436 farmers and breeders might require explanations to understand why and how the model
437 came to a prediction and to ensure its trustworthiness.

438 Presently, the hope for human-comprehensible explanations for black-box algorithms to
439 increase technical confidence, generate trust, and make better informed choices
440 remains an open challenge. In light of this challenge, we strongly recommend that a
441 single prediction might therefore need to be explained in various ways, reflecting the
442 requirements of all stakeholders.

443 **How to devise X-AI-driven analytics for phenomics questions**

444 Interpretable by design models

445 X-AI bears great potential for the analysis and interpretation of phenomic data. In what
446 follows, we provide an example of an X-AI workflow design and describe for the first
447 time, the steps needed to foster practical applicability of interpretable by design
448 algorithms in phenomics image analysis (Figure 6). We have also accompanied this
449 review by an interactive online tutorial that acts as an educational resource, intended for
450 readers with little to no knowledge of X-AI algorithms; it also serves as a good starting
451 point for self-learning and raises an early awareness that computational phenomics
452 need not be intimidating. In addition, we have created a set of self-test quizzes and
453 hands-on practice exercises to provide users with opportunities to augment their
454 learning by practically applying the concepts explained in order to assess their acquired
455 knowledge. The code and computational **notebooks** are open source and freely
456 accessible through our GitHub repository^{vi}. Collectively, this will accelerate the rate of
457 discovery and move toward open science and AI ethics in digital phenomics.

458 In our tutorial, we train ‘this looks like that’ algorithm to classify diseases using the
459 **crowdsourced** cassava disease classification dataset. This choice is motivated by the
460 importance of cassava, being a key food security crop grown by smallholder farmers in
461 Africa, Asia, and South America. However, diseases that plague the crop are a major
462 cause of poor yield [95]. Existing methods to identify diseases require governmental
463 agricultural experts to visually inspect and diagnose the plants [96]. This labor-intensive
464 process makes it difficult to monitor and treat disease progression. With the help of X-
465 AI, we can identify and classify cassava diseases and monitor their progression rapidly
466 enough to address these current limitations in disease surveillance. Our model,

467 initialized with **transfer learning** (TL), was trained to predict five classes, and provide
468 corresponding prototypical explanations by marking activated patches with bounding
469 boxes and generating heatmaps to show which parts of the image are similar to the
470 prototypes. The resulting **confusion matrix** illustrates the percentage of correctly
471 classified images in each class. The overall accuracy was 88.7% after cycling through
472 the training set 240 times (Figure 6).

473 A more detailed description of all steps of the analysis, including the computational
474 notebook and code to train, validate, and test/replicate our models, is provided on the
475 tutorial website. This description covers, as relevant, data preprocessing, image classes
476 and format, architecture of the model, model training and evaluation, prototypical
477 explanations, and the **computer cluster** used for training.

478 Dealing with small datasets

479 As some phenotyping experiments generate small amounts of data, X-AI models get
480 fewer training examples to learn from. But how can models learn well from small
481 datasets? Using low complexity models that have a small number of trainable
482 parameters can perform better than complex ones as they are less prone to overfitting
483 and generalize better [97]. Additionally, TL can be employed to transfer knowledge
484 acquired while learning a different but related task from a model trained on a large
485 dataset to fit a new model using a small dataset [98]. When TL is not powerful enough
486 due to the absence of large datasets, cumulative learning (CL) can be used to train a
487 model over various small datasets and accumulate knowledge in the resulting network
488 representation (i.e., model weights) [99]. Even when pretraining a model with TL or CL

489 is not possible, the cosine loss function can substantially improve the predictive
490 performance of the model [100]. While the loss function of the model measures the error
491 between the input and predicted output, the cosine loss function maximizes the cosine
492 similarity between them. One-shot or few-shot learning can also be used to train a
493 model from one or a handful of training image data by basing predictions on a similarity
494 metric (e.g., cosine similarity) that compares training data to new inputs [101]. Most
495 recently, with the embedding of human knowledge into AI, it will be possible to
496 supplement small training datasets. Representation of such knowledge can be
497 incorporated into AI by means of changes to the input data and loss function [102], to
498 the architecture of the algorithm [103], or to a combination thereof. Alternatively,
499 oversampling can be another workaround that produces new sample data to augment
500 small datasets. These new data, however, should be meaningful, sufficient, and
501 realistic, and should contribute for better performance of predictive models [104].
502 Oversampling can be achieved by: (i) performing geometric transformations on existing
503 images using primitive data manipulation techniques, including flipping, rotation,
504 shearing, cropping, and translation [105]; (ii) generating new **synthetic data** with
505 generative adversarial networks (GANs), which are powerful models for learning
506 complex distributions to synthesize semantically meaningful samples from an actual
507 training set [104]. GANs can be employed for image-to-image translation, fusion image
508 generation, label-to-image mapping, and text-to-image translation [104]; (iii) simulating
509 real-world scenarios, by making use of virtual reality [106] or other extended reality
510 technologies, including augmented and mixed reality, to create immersive 3D virtual
511 environments, in which cameras can automatically collect photorealistic synthetic

512 images; and (iv) pairing existing images using methods such as cut-and-paste [107] or
513 CutMix [108], which automatically ‘cut’ objects of interest from training images and
514 ‘paste’ them on random backgrounds or on other training images, respectively.

515 **How to ensure that human-centric X-AI benefits all: team science, open science,**
516 **open education, and embedded ethics**

517 AI in phenomics can potentially impact many aspects of plant science, from basic
518 research discovery to translational research. It is critical that these advances in
519 technology broadly benefit society as a whole.

520 So, how do we effectively ensure that human-centric X-AI benefits and does not harm
521 individuals and communities? This can be done in several ways.

522 First, we suggest that pivoting toward multidisciplinary team science is necessary to
523 tackle the most pressing scientific, societal, and ethical problems of plant digital
524 phenomics. Over the last decade, funding agencies across the US and Europe
525 dedicated resources to facilitating team science. This work is evidenced by
526 interdisciplinary and multidisciplinary team requirements in funding announcements and
527 programs. For example, addressing the problem of bias in phenomics AI requires the
528 integrated knowledge of socially and intellectually diverse researchers who specialize in
529 plant science, plant phenomics, plant pathology, data science, computer science, IS,
530 social science, and bioethics, just to name a few.

531 Second, we emphasize the crucial importance of an open science system that aspires
532 to open access not only to research outputs, but the whole research process, and posit

533 that all phenomics and data centers should participate in these practices. Promotion of
534 open science and team science are synergistic goals, both of which are essential for
535 improving our knowledge and scientific rigor.

536 Third, we call for mobilizing open educational resources relevant to AI in phenomics that
537 advocate digitized materials offered freely and openly for educators, students, and
538 interested learners worldwide, including developing countries to use and reuse for
539 teaching, learning, training, and research. Open education holds great promise to create
540 knowledge and put it to use, promote content quality through sharing of materials for
541 feedback and continuous improvement, and achieve competencies.

542 Fourth, we propose the development of socially and ethically responsible AI in
543 phenomics by reforming curricula and embedding bioethicists into the technology
544 development team. Ethical concerns around AI regarding handling of data, data bias,
545 transparency, explainability, and responsibility, have prompted us to consider how AI
546 technology can be designed, implemented, deployed, and monitored post deployment in
547 an ethical manner. Embedding bioethicists into the AI development team can ensure
548 that developers be practically assisted in anticipating, identifying, and addressing ethical
549 issues through critical ethical reasoning and bioethical decision-making. Universities
550 across the US and Europe have recently joined the effort to develop socially responsible
551 AI by reforming curricula. For example, Harvard University initiated an ‘Embedded
552 EthiCS’ curriculum that integrates ethical issues into the core computer science
553 curriculum. We advocate that universities around the world implement similar
554 approaches to empower students and early-career scientists to think ethically as they
555 develop algorithms and build AI systems, in their studies, in their new business

556 ventures, and as they pursue technical work in their careers. These free and open
557 courses should be taught by interdisciplinary teams of computer scientists, social
558 scientists, and bioethicists.

559 We encourage the scientific community to embrace a growth mindset regarding team
560 science, open science, open education, and embedded ethics, which altogether can be
561 harnessed to create extraordinary phenomic resources that benefit all. The rewards to
562 these efforts come from investments of energy, time, and action.

563 **Concluding remarks and future perspectives**

564 We are experiencing an unprecedented time where the availability of vast amounts of
565 phenomic data, combined with advances in AI, is providing the opportunity to
566 turbocharge the data to insight journey. This opportunity is an incentive to not only
567 design and implement effective and reliable data management strategies but also to
568 improve visibility, accessibility, and usability of publicly available datasets that can
569 support research and innovation in plant digital phenomics.

570 Although AI has demonstrated impressive potential in phenomics, risks due to bias and
571 lack of transparency of models should be considered. Reducing these risks entails
572 multidisciplinary science and technology teams working together. The involvement of
573 plant scientists, data scientists, and IS scientists during the complete lifecycle of AI
574 analysis is integral to ensure explainability and to identify bias in the predictive models.
575 Interpretable by design models can potentially be leveraged to mitigate bias and provide
576 transparency into the decision-making process.

577 In the past, AI research focused on a one-way interaction, from AI to humans; today,
578 human-centric X-AI aims to enable bidirectional interaction so that human intelligence
579 and AI are brought together to collectively achieve superior results and continuously
580 improve by learning from each other. Human-centric X-AI will have an extraordinary
581 impact on phenomics in the near future, and we should do all we can to ensure that it is
582 designed, implemented, deployed, and regulated in a way that maximizes benefits for
583 breeders, farmers, and consumers. In this regard, the academic and AI communities
584 should ensure that computational phenomics, in addition to social and ethical analysis,
585 are integrated into plant science curriculum as a step toward this goal (see Outstanding
586 Questions).

587 **Acknowledgements**

588 The authors are grateful to the editor and three anonymous reviewers for their
589 constructive and insightful comments which greatly helped improve the manuscript.
590 Partial support for this work was provided by the EU FP7 project WATBIO, grant no.
591 311929; the EU H2020 project EMPHASIS-PREP and its Italian node, PHEN-ITALY,
592 grant no. 739514; the Italian Ministry of University and Research Brain Gain
593 Professorship to A.L.H; the Center for Bioenergy Innovation, a U.S. Department of
594 Energy (DOE) Bioenergy Research Center, the Plant Microbe Interface SFA, and the
595 integrated Pennycress Resilience Project, all supported by the Biological and
596 Environmental Research in the DOE Office of Science; the Oak Ridge Leadership
597 Computing Facility, a DOE Office of Science User Facility supported under Contract
598 DE-AC05-00OR22725; and the DOE, Laboratory Directed Research and Development
599 funding ORNL AI Initiative ProjectID 10875 at the Oak Ridge National Laboratory.

600 **Resources**

601 ⁱ<https://www.top500.org/lists/top500/2022/06/>

602 ⁱⁱ[https://www.whitehouse.gov/wp-content/uploads/2020/01/Draft-OMB-Memo-on-](https://www.whitehouse.gov/wp-content/uploads/2020/01/Draft-OMB-Memo-on-Regulation-of-AI-1-7-19.pdf)
603 [Regulation-of-AI-1-7-19.pdf](https://www.whitehouse.gov/wp-content/uploads/2020/01/Draft-OMB-Memo-on-Regulation-of-AI-1-7-19.pdf)

604 ⁱⁱⁱ[https://eur-lex.europa.eu/legal-](https://eur-lex.europa.eu/legal-content/EN/TXT/?qid=1623335154975&uri=CELEX%3A52021PC0206)
605 [content/EN/TXT/?qid=1623335154975&uri=CELEX%3A52021PC0206](https://eur-lex.europa.eu/legal-content/EN/TXT/?qid=1623335154975&uri=CELEX%3A52021PC0206)

606 ^{iv}[https://eur-lex.europa.eu/legal-](https://eur-lex.europa.eu/legal-content/en/TXT/?qid=1593079180383&uri=CELEX%3A52020DC0064)
607 [content/en/TXT/?qid=1593079180383&uri=CELEX%3A52020DC0064](https://eur-lex.europa.eu/legal-content/en/TXT/?qid=1593079180383&uri=CELEX%3A52020DC0064)

608 ^v<https://www.darpa.mil/attachments/DARPA-BAA-16-53.pdf>

609 ^{vi}<https://github.com/HarfoucheLab/a-primer-on-AI-in-plant-digital-phenomics>

610 **References**

- 611 1. Yang, W. *et al.* (2020) Crop phenomics and high-throughput phenotyping: past
612 decades, current challenges, and future perspectives. *Mol. Plant* 13, 187–214
- 613 2. Houle, D. *et al.* (2010) Phenomics: the next challenge. *Nat. Rev. Genet.* 11, 855–
614 866
- 615 3. Tardieu, F. *et al.* (2017) Plant phenomics, from sensors to knowledge. *Curr. Biol.*
616 27, 770–783
- 617 4. Dhondt, S. *et al.* (2013) Cell to whole-plant phenotyping: the best is yet to come.
618 *Trends Plant Sci.* 18, 428–439

- 619 5. Fiorani, F. and Schurr, U. (2013) Future scenarios for plant phenotyping. *Annu.*
620 *Rev. Plant Biol.* 64, 267–291
- 621 6. Furbank, R.T. and Tester, M. (2011) Phenomics – technologies to relieve the
622 phenotyping bottleneck. *Trends Plant Sci.* 16, 635–644
- 623 7. Song, P. *et al.* (2021) High-throughput phenotyping: breaking through the
624 bottleneck in future crop breeding. *Crop J.* 9, 633–645
- 625 8. Huang, Y. *et al.* (2020) Phenotypic techniques and applications in fruit trees: a
626 review. *Plant Methods* 16, 107
- 627 9. Feng, L. *et al.* (2021) A comprehensive review on recent applications of
628 unmanned aerial vehicle remote sensing with various sensors for high-throughput
629 plant phenotyping. *Comput. Electron. Agric.* 182, 106033
- 630 10. Coppens, F. *et al.* (2017) Unlocking the potential of plant phenotyping data
631 through integration and data-driven approaches. *Curr. Opin. Syst. Biol.* 4, 58–63
- 632 11. Nakhle, F. and Harfouche, A.L. (2021) Ready, steady, go AI: a practical tutorial on
633 fundamentals of artificial intelligence and its applications in phenomics image
634 analysis. *Patterns* 2, 100323
- 635 12. Nevavuori, P. *et al.* (2019) Crop yield prediction with deep convolutional neural
636 networks. *Comput. Electron. Agric.* 163, 104859
- 637 13. Gopal, G. *et al.* (2019) Identification of plant leaf diseases using a nine-layer deep
638 convolutional neural network. *Comput. Electr. Eng.* 76, 323–338
- 639 14. Smith, A.G. *et al.* (2020) Segmentation of roots in soil with U-Net. *Plant Methods*

- 640 16, 13
- 641 15. Yasrab, R. *et al.* (2019) RootNav 2.0: deep learning for automatic navigation of
642 complex plant root architectures. *Gigascience* 8, giz123
- 643 16. Gaggion, N. *et al.* (2021) ChronoRoot: high-throughput phenotyping by deep
644 segmentation networks reveals novel temporal parameters of plant root system
645 architecture. *Gigascience* 10, giab052
- 646 17. Ubbens, J.R. and Stavness, I. (2017) Deep plant phenomics: a deep learning
647 platform for complex plant phenotyping tasks. *Front. Plant Sci.* 8, 1190
- 648 18. Ghosal, S. *et al.* (2018) An explainable deep machine vision framework for plant
649 stress phenotyping. *Proc. Natl. Acad. Sci.* 115, 4613–4618
- 650 19. Qiu, Z. *et al.* (2021) Estimation of nitrogen nutrition index in rice from UAV RGB
651 images coupled with machine learning algorithms. *Comput. Electron. Agric.* 189,
652 106421
- 653 20. Shi, P. *et al.* (2021) Rice nitrogen nutrition estimation with RGB images and
654 machine learning methods. *Comput. Electron. Agric.* 180, 105860
- 655 21. Resente, G. *et al.* (2021) Mask, train, repeat! Artificial intelligence for quantitative
656 wood anatomy. *Front. Plant Sci.* 12, 767400
- 657 22. Nagasubramanian, K. *et al.* (2019) Plant disease identification using explainable
658 3D deep learning on hyperspectral images. *Plant Methods* 15, 98
- 659 23. Shu, M. *et al.* (2021) The application of UAV-based hyperspectral imaging to
660 estimate crop traits in maize inbred lines. *Plant Phenomics* 2021, 1–14

- 661 24. Bodner, G. *et al.* (2018) Hyperspectral imaging: a novel approach for plant root
662 phenotyping. *Plant Methods* 14, 84
- 663 25. Burnett, A.C. *et al.* (2021) Detection of the metabolic response to drought stress
664 using hyperspectral reflectance. *J. Exp. Bot.* 72, 6474–6489
- 665 26. Tauro, F. *et al.* (2022) Latent heat flux variability and response to drought stress
666 of black poplar: A multi-platform multi-sensor remote and proximal sensing
667 approach to relieve the data scarcity bottleneck. *Remote Sens. Environ.* 268,
668 112771
- 669 27. Gené-Mola, J. *et al.* (2020) Fruit detection, yield prediction and canopy geometric
670 characterization using LiDAR with forced air flow. *Comput. Electron. Agric.* 168,
671 105121
- 672 28. Raza, S.-A. *et al.* (2014) Automatic detection of regions in spinach canopies
673 responding to soil moisture deficit using combined visible and thermal imagery.
674 *PLoS One* 9, e97612
- 675 29. Gutiérrez, S. *et al.* (2018) Vineyard water status assessment using on-the-go
676 thermal imaging and machine learning. *PLoS One* 13, e0192037
- 677 30. Li, J. *et al.* (2021) DeepLearnMOR: a deep-learning framework for fluorescence
678 image-based classification of organelle morphology. *Plant Physiol.* 186, 1786–
679 1799
- 680 31. Soltaninejad, M. *et al.* (2020) Three dimensional root CT segmentation using
681 multi-resolution encoder-decoder networks. *IEEE Trans. Image Process.* 29,
682 6667–6679

- 683 32. Maimaitijiang, M. *et al.* (2020) Soybean yield prediction from UAV using
684 multimodal data fusion and deep learning. *Remote Sens. Environ.* 237, 111599
- 685 33. Maimaitijiang, M. *et al.* (2017) Unmanned aerial system (UAS)-based phenotyping
686 of soybean using multi-sensor data fusion and extreme learning machine. *ISPRS*
687 *J. Photogramm. Remote Sens.* 134, 43–58
- 688 34. Streich, J. *et al.* (2020) Can exascale computing and explainable artificial
689 intelligence applied to plant biology deliver on the United Nations sustainable
690 development goals? *Curr. Opin. Biotechnol.* 61, 217–225
- 691 35. Harfouche, A.L. *et al.* (2019) Accelerating climate resilient plant breeding by
692 applying next-generation artificial intelligence. *Trends Biotechnol.* 37, 1217–1235
- 693 36. Chhetri, H.B. *et al.* (2020) Genome-wide association study of wood anatomical
694 and morphological traits in *Populus trichocarpa*. *Front. Plant Sci.* 11, 545748
- 695 37. Chhetri, H.B. *et al.* (2019) Multitrait genome-wide association analysis of *Populus*
696 *trichocarpa* identifies key polymorphisms controlling morphological and
697 physiological traits. *New Phytol.* 223, 293–309
- 698 38. Poorter, H. *et al.* (2016) Pampered inside, pestered outside? differences and
699 similarities between plants growing in controlled conditions and in the field. *New*
700 *Phytol.* 212, 838–855
- 701 39. Wilkinson, M.D. *et al.* (2016) The FAIR guiding principles for scientific data
702 management and stewardship. *Sci. Data* 3, 160018
- 703 40. Williamson, H.F. *et al.* (2021) Data management challenges for artificial

- 704 intelligence in plant and agricultural research. *F1000Research* 10, 324
- 705 41. Papoutsoglou, E.A. *et al.* (2020) Enabling reusability of plant phenomic datasets
706 with MIAPPE 1.1. *New Phytol.* 227, 260–273
- 707 42. Ówiek-Kupczyńska, H. *et al.* (2016) Measures for interoperability of phenotypic
708 data: minimum information requirements and formatting. *Plant Methods* 12, 44
- 709 43. Wolf, M. *et al.* (2021) Reusability first: toward FAIR workflows. In *2021 IEEE*
710 *International Conference on Cluster Computing (CLUSTER)*, pp. 444–455
- 711 44. McMahan, B. *et al.* (2017) Communication-efficient learning of deep networks
712 from decentralized data. In *Proceedings of the 20th International Conference on*
713 *Artificial Intelligence and Statistics*, 54, pp. 1273–1282
- 714 45. Chen, M. *et al.* (2021) Communication-efficient federated learning. *Proc. Natl.*
715 *Acad. Sci.* 118, e2024789118
- 716 46. Kaissis, G.A. *et al.* (2020) Secure, privacy-preserving and federated machine
717 learning in medical imaging. *Nat. Mach. Intell.* 2, 305–311
- 718 47. Kaissis, G. *et al.* (2021) End-to-end privacy preserving deep learning on multi-
719 institutional medical imaging. *Nat. Mach. Intell.* 3, 473–484
- 720 48. Dayan, I. *et al.* (2021) Federated learning for predicting clinical outcomes in
721 patients with COVID-19. *Nat. Med.* 27, 1735–1743
- 722 49. Brink, S.C. (2021) 25 years of trends in plant science: we should all be plant
723 worshippers. *Trends Plant Sci.* 26, 527–529
- 724 50. Blei, D.M. and Smyth, P. (2017) Science and data science. *Proc. Natl. Acad. Sci.*

- 725 114, 8689–8692
- 726 51. Xu, L. *et al.* (2015) Plant photosynthesis phenomics data quality control.
727 *Bioinformatics* 31, 1796–1804
- 728 52. Miceli, M. *et al.* (2020) Between subjectivity and imposition: power dynamics in
729 data annotation for computer vision. *Proc. ACM Human-Computer Interact.* 4, 1–
730 25
- 731 53. Akter, S. *et al.* (2021) Algorithmic bias in data-driven innovation in the age of AI.
732 *Int. J. Inf. Manage.* 60, 102387
- 733 54. Boyd, D. and Crawford, K. (2012) Critical questions for big data. *Information,*
734 *Commun. Soc.* 15, 662–679
- 735 55. Mehta, P. *et al.* (2019) A high-bias, low-variance introduction to Machine Learning
736 for physicists. *Phys. Rep.* 810, 1–124
- 737 56. Greener, J.G. *et al.* (2022) A guide to machine learning for biologists. *Nat. Rev.*
738 *Mol. Cell Biol.* 23, 40–55
- 739 57. Schwartz, R. *et al.* (2021) A proposal for identifying and managing bias in artificial
740 intelligence. *NIST Spec. Publ.* Doi: 10.6028/NIST.SP.1270-draft
- 741 58. Großkinsky, D.K. *et al.* (2015) Plant phenomics and the need for physiological
742 phenotyping across scales to narrow the genotype-to-phenotype knowledge gap.
743 *J. Exp. Bot.* 66, 5429–5440
- 744 59. Matthews, D. (2018) Supercharge your data wrangling with a graphics card.
745 *Nature* 562, 151–152

- 746 60. Mann, A. (2020) Core concept: nascent exascale supercomputers offer promise,
747 present challenges. *Proc. Natl. Acad. Sci.* 117, 22623–22625
- 748 61. Skibba, R. (2021) Japan’s fugaku supercomputer crushes competition, but likely
749 not for long. *Engineering* 7, 6–7
- 750 62. Marx, V. (2021) Biology begins to tangle with quantum computing. *Nat. Methods*
751 18, 715–719
- 752 63. Dart, E. *et al.* (2014) The science DMZ: a network design pattern for data-
753 intensive science. *Sci. Program.* 22, 701405
- 754 64. Phillips, S.J. *et al.* (2009) Sample selection bias and presence-only distribution
755 models: implications for background and pseudo-absence data. *Ecol. Appl.* 19,
756 181–197
- 757 65. Castro, D.C. *et al.* (2020) Causality matters in medical imaging. *Nat. Commun.*
758 11, 3673
- 759 66. Buda, M. *et al.* (2018) A systematic study of the class imbalance problem in
760 convolutional neural networks. *Neural Networks* 106, 249–259
- 761 67. Lever, J. *et al.* (2016) Regularization. *Nat. Methods* 13, 803–804
- 762 68. Jin, L. *et al.* (2020) Deep-learning-assisted detection and segmentation of rib
763 fractures from CT scans: Development and validation of FracNet. *EBioMedicine*
764 62, 103106
- 765 69. Sreeram, M. and Nof, S.Y. (2021) Human-in-the-loop: role in cyber physical
766 agricultural systems. *Int. J. Comput. Commun. Control* 16, 4166

- 767 70. Budd, S. *et al.* (2021) A survey on active learning and human-in-the-loop deep
768 learning for medical image analysis. *Med. Image Anal.* 71, 102062
- 769 71. Patel, B.N. *et al.* (2019) Human–machine partnership with artificial intelligence for
770 chest radiograph diagnosis. *npj Digit. Med.* 2, 111
- 771 72. Basu, S. *et al.* (2018) Iterative random forests to discover predictive and stable
772 high-order interactions. *Proc. Natl. Acad. Sci.* 115, 1943–1948
- 773 73. Cliff, A. *et al.* (2019) A high-performance computing implementation of iterative
774 random forest for the creation of predictive expression networks. *Genes (Basel)*.
775 10, 996
- 776 74. Jiang, Y. and Li, C. (2020) Convolutional neural networks for image-based high-
777 throughput plant phenotyping: a review. *Plant Phenomics* 2020, 1–22
- 778 75. Danilevicz, M.F. *et al.* (2021) Resources for image-based high-throughput
779 phenotyping in crops and data sharing challenges. *Plant Physiol.* 187, 699–715
- 780 76. Mochida, K. *et al.* (2019) Computer vision-based phenotyping for improvement of
781 plant productivity: a machine learning perspective. *Gigascience* 8, giy153
- 782 77. Vilone, G. and Longo, L. (2020) Explainable artificial intelligence: a systematic
783 review. *arXiv* 2006.00093
- 784 78. Shortliffe, E.H. *et al.* (1973) An artificial intelligence program to advise physicians
785 regarding antimicrobial therapy. *Comput. Biomed. Res.* 6, 544–560
- 786 79. Doyle, J. (1979) A truth maintenance system. *Artif. Intell.* 12, 231–272
- 787 80. Tickle, A.B. *et al.* (1998) The truth will come to light: directions and challenges in

- 788 extracting the knowledge embedded within trained artificial neural networks. *IEEE*
789 *Trans. Neural Networks* 9, 1057–1068
- 790 81. Lipton, Z.C. (2018) The mythos of model interpretability. *Queue* 16, 31–57
- 791 82. Molnar, C. (2019) *Interpretable machine learning. A guide for making black box*
792 *models explainable*, (1st edition), leanpub.com.
- 793 83. Ribeiro, M.T. *et al.* (2016) “Why should I trust you?”: explaining the predictions of
794 any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference*
795 *on knowledge discovery and data mining*, pp. 1135–1144
- 796 84. Lundberg, S.M. and Lee, S.I. (2017) A unified approach to interpreting model
797 predictions. In *Advances in Neural Information Processing Systems*, 30, pp.
798 4768–4777
- 799 85. Ghosal, S. *et al.* (2018) An explainable deep machine vision framework for plant
800 stress phenotyping. *Proc. Natl. Acad. Sci.* 115, 4613–4618
- 801 86. Toda, Y. and Okura, F. (2019) How convolutional neural networks diagnose plant
802 disease. *Plant Phenomics* 2019, 1–14
- 803 87. Nagasubramanian, K. *et al.* (2019) Plant disease identification using explainable
804 3D deep learning on hyperspectral images. *Plant Methods* 15, 98
- 805 88. Nagasubramanian, K. *et al.* (2020) Usefulness of interpretability methods to
806 explain deep learning based plant stress phenotyping. *arXiv* 2007.05729
- 807 89. Mostafa, S. *et al.* (2021) Visualizing feature maps for model selection in
808 convolutional neural networks. In *Proceedings of the IEEE/CVF International*

- 809 *Conference on Computer Vision (ICCV) Workshops*, pp. 1362–1371
- 810 90. Dobrescu, A. *et al.* (2019) Understanding deep neural networks for regression in
811 leaf counting. In *2019 IEEE/CVF Conference on Computer Vision and Pattern
812 Recognition Workshops (CVPRW)*, pp. 2600–2608
- 813 91. Rudin, C. (2019) Stop explaining black box machine learning models for high
814 stakes decisions and use interpretable models instead. *Nat. Mach. Intell.* 1, 206–
815 215
- 816 92. Chen, C. *et al.* (2019) This looks like that: deep learning for interpretable image
817 recognition. In *Advances in Neural Information Processing Systems*, 32, pp.
818 8930–8941
- 819 93. Barredo Arrieta, A. *et al.* (2020) Explainable artificial intelligence (XAI): concepts,
820 taxonomies, opportunities and challenges toward responsible AI. *Inf. Fusion* 58,
821 82–115
- 822 94. Preece, A. *et al.* (2018) Stakeholders in explainable AI. *arXiv* 1810.00184
- 823 95. Otim-Napri, G.W. *et al.* (2001) Changes in the incidence and severity of Cassava
824 mosaic virus disease, varietal diversity and cassava production in Uganda. *Ann.
825 Appl. Biol.* 138, 313–327
- 826 96. Mwebaze, E. *et al.* (2019) iCassava 2019 fine-grained visual categorization
827 challenge. *arXiv* 1908.02900
- 828 97. Brigato, L. and Iocchi, L. (2021) A close look at deep learning with small data. In
829 *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 2490–

- 830 2497
- 831 98. Gupta, V. *et al.* (2021) Cross-property deep transfer learning framework for
832 enhanced predictive analytics on small materials data. *Nat. Commun.* 12, 6595
- 833 99. Seddiki, K. *et al.* (2020) Cumulative learning enables convolutional neural network
834 representations for small mass spectrometry data classification. *Nat. Commun.*
835 11, 5595
- 836 100. Barz, B. and Denzler, J. (2020) Deep learning on small datasets without pre-
837 training using cosine loss. In *2020 IEEE Winter Conference on Applications of*
838 *Computer Vision (WACV)*, pp. 1360–1369
- 839 101. Wang, Y. *et al.* (2021) Generalizing from a few examples: a survey on few-shot
840 learning. *ACM Comput. Surv.* 53, 1–34
- 841 102. Childs, C.M. and Washburn, N.R. (2019) Embedding domain knowledge for
842 machine learning of complex material systems. *MRS Commun.* 9, 806–820
- 843 103. Hasebe, T. (2021) Knowledge-embedded message-passing neural networks:
844 improving molecular property prediction with human knowledge. *ACS Omega* 6,
845 27955–27967
- 846 104. Shamsolmoali, P. *et al.* (2021) Image synthesis with adversarial networks: a
847 comprehensive survey and case studies. *Inf. Fusion* 72, 126–146
- 848 105. Khalifa, N.E. *et al.* (2022) A comprehensive survey of recent trends in deep
849 learning for digital images augmentation. *Artif. Intell. Rev.* 55, 2351–2377
- 850 106. Martinez-Gonzalez, P. *et al.* (2020) UnrealROX: an extremely photorealistic virtual

- 851 reality environment for robotics simulations and synthetic data generation. *Virtual*
852 *Real.* 24, 271–288
- 853 107. Dwibedi, D. *et al.* (2017) Cut, paste and learn: surprisingly easy synthesis for
854 instance detection. In *2017 IEEE/CVF International Conference on Computer*
855 *Vision (ICCV)*, pp. 1310–1319
- 856 108. Yun, S. *et al.* (2019) CutMix: regularization strategy to train strong classifiers with
857 localizable features. In *2019 IEEE/CVF International Conference on Computer*
858 *Vision (ICCV)*, pp. 6022–6031
- 859 109. Pound, M.P. *et al.* (2017) Deep machine learning provides state-of-the-art
860 performance in image-based plant phenotyping. *Gigascience* 6, gix083
- 861 110. Pound, M.P. *et al.* (2017) Deep learning for multi-task plant phenotyping. In *2017*
862 *IEEE International Conference on Computer Vision Workshops (ICCVW)*, pp.
863 2055–2063
- 864 111. David, E. *et al.* (2020) Global wheat head detection (GWHD) dataset: a large and
865 diverse dataset of high-resolution RGB-labelled images to develop and
866 benchmark wheat head detection methods. *Plant Phenomics* 2020, 1–12
- 867 112. David, E. *et al.* (2021) Global wheat head detection 2021: an improved dataset for
868 benchmarking wheat head detection methods. *Plant Phenomics* 2021, 1–9
- 869 113. Wang, X. *et al.* (2019) High-throughput phenotyping with deep learning gives
870 insight into the genetic architecture of flowering time in wheat. *Gigascience* 8, 1–
871 11

- 872 114. Quiñones, R. *et al.* (2021) Multi-feature data repository development and analytics
873 for image cosegmentation in high-throughput plant phenotyping. *PLoS One* 16,
874 e0257001
- 875 115. Taghavi Namin, S. *et al.* (2018) Deep phenotyping: Deep learning for temporal
876 phenotype/genotype classification. *Plant Methods* 14, 66
- 877 116. Minervini, M. *et al.* (2016) Finely-grained annotated datasets for image-based
878 plant phenotyping. *Pattern Recognit. Lett.* 81, 80–89
- 879 117. Cruz, J.A. *et al.* (2016) Multi-modality imagery database for plant phenotyping.
880 *Mach. Vis. Appl.* 27, 735–749
- 881 118. Dobos, O. *et al.* (2019) A deep learning-based approach for high-throughput
882 hypocotyl phenotyping. *Plant Physiol.* 181, 1415–1424
- 883 119. Khanna, R. *et al.* (2019) A spatio temporal spectral framework for plant stress
884 phenotyping. *Plant Methods* 15, 13
- 885 120. Sa, I. *et al.* (2018) WeedMap: a large-scale semantic weed mapping framework
886 using aerial multispectral imaging and deep neural network for precision farming.
887 *Remote Sens.* 10, 1423
- 888 121. Chebrolu, N. *et al.* (2017) Agricultural robot dataset for plant classification,
889 localization and mapping on sugar beet fields. *Int. J. Rob. Res.* 36, 1045–1052
- 890 122. Mignoni, M.E. *et al.* (2022) Soybean images dataset for caterpillar and *Diabrotica*
891 *speciosa* pest detection and classification. *Data Br.* 40, 107756
- 892 123. dos Santos Ferreira, A. *et al.* (2017) Weed detection in soybean crops using

- 893 ConvNets. *Comput. Electron. Agric.* 143, 314–324
- 894 124. Bosilj, P. *et al.* (2020) Transfer learning between crop types for semantic
895 segmentation of crops versus weeds in precision agriculture. *J. F. Robot.* 37, 7–
896 19
- 897 125. Haug, S. and Ostermann, J. (2015) A crop/weed field image dataset for the
898 evaluation of computer vision based precision agriculture tasks. In *Computer*
899 *Vision - ECCV 2014 Workshops*, pp. 105–116
- 900 126. Nakatumba-Nabende, J. *et al.* (2020) A dataset of necrotized cassava root cross-
901 section images. *Data Br.* 32, 106170
- 902 127. Thapa, R. *et al.* (2020) The plant pathology challenge 2020 data set to classify
903 foliar disease of apples. *Appl. Plant Sci.* 8, e11390
- 904 128. Gené-Mola, J. *et al.* (2020) LFuji-air dataset: annotated 3D LiDAR point clouds of
905 Fuji apple trees for fruit detection scanned under different forced air flow
906 conditions. *Data Br.* 29, 105248
- 907 129. Hani, N. *et al.* (2020) MinneApple: a benchmark dataset for apple detection and
908 segmentation. *IEEE Robot. Autom. Lett.* 5, 852–858
- 909 130. Dias, P.A. *et al.* (2018) Multispecies fruit flower detection using a refined semantic
910 segmentation network. *IEEE Robot. Autom. Lett.* 3, 3003–3010
- 911 131. Fenu, G. and Mallocci, F.M. (2021) DiaMOS plant: a dataset for diagnosis and
912 monitoring plant disease. *Agronomy* 11, 2107
- 913 132. Bargoti, S. and Underwood, J. (2017) Deep fruit detection in orchards. In *2017*

- 914 *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3626–
915 3633
- 916 133. Kusriani, K. *et al.* (2020) Data augmentation for automated pest classification in
917 Mango farms. *Comput. Electron. Agric.* 179, 105842
- 918 134. Kicherer, A. *et al.* (2015) An automated field phenotyping pipeline for application
919 in grapevine research. *Sensors* 15, 4823–4836
- 920 135. Navarro, P.J. *et al.* (2022) A novel ground truth multispectral image dataset with
921 weight, anthocyanins, and Brix index measures of grape berries tested for its
922 utility in machine learning pipelines. *Gigascience* 11, giac052
- 923 136. Alessandrini, M. *et al.* (2021) A grapevine leaves dataset for early detection and
924 classification of esca disease in vineyards through machine learning. *Data Br.* 35,
925 106809
- 926 137. Abdelghafour, F. *et al.* (2021) An annotated image dataset of downy mildew
927 symptoms on Merlot grape variety. *Data Br.* 37, 107250
- 928 138. Kestur, R. *et al.* (2019) MangoNet: a deep semantic segmentation architecture for
929 a method to detect and count mangoes in an open orchard. *Eng. Appl. Artif. Intell.*
930 77, 59–69
- 931 139. Chouhan, S.S. *et al.* (2019) A data repository of leaf images: practice towards
932 plant conservation with plant pathology. In *2019 4th International Conference on*
933 *Information Systems and Computer Networks (ISCON)*, pp. 700–707
- 934 140. Ali, H. *et al.* (2017) Symptom based automated detection of citrus diseases using

- 935 color histogram and textural descriptors. *Comput. Electron. Agric.* 138, 92–104
- 936 141. Sharif, M. *et al.* (2018) Detection and classification of citrus diseases in agriculture
937 based on optimized weighted segmentation and feature selection. *Comput.*
938 *Electron. Agric.* 150, 220–234
- 939 142. Chitwood, D.H. and Otoni, W.C. (2017) Morphometric analysis of *Passiflora*
940 leaves: the relationship between landmarks of the vasculature and elliptical
941 Fourier descriptors of the blade. *Gigascience* 6, giw008
- 942 143. Chitwood, D.H. and Otoni, W.C. (2017) Divergent leaf shapes among *Passiflora*
943 species arise from a shared juvenile morphology. *Plant Direct* 1, e00028
- 944 144. Kaufmane, E. *et al.* (2022) QuinceSet: dataset of annotated Japanese quince
945 images for object detection. *Data Br.* 42, 108332
- 946 145. Altaheri, H. *et al.* (2019) Date fruit dataset for intelligent harvesting. *Data Br.* 26,
947 104514
- 948 146. Altaheri, H. *et al.* (2019) Date fruit classification for robotic harvesting in a natural
949 environment using deep learning. *IEEE Access* 7, 117115–117133
- 950 147. Wiesner-Hanks, T. *et al.* (2018) Image set for deep learning: field images of maize
951 annotated with disease symptoms. *BMC Res. Notes* 11, 440
- 952 148. Lac, L. *et al.* (2022) An annotated image dataset of vegetable crops at an early
953 stage of growth for proximal sensing applications. *Data Br.* 42, 108035
- 954 149. Schunck, D. *et al.* (2021) Pheno4D: a spatio-temporal dataset of maize and
955 tomato plant point clouds for phenotyping and advanced plant analysis. *PLoS*

956 *One* 16, e025634

957 150. Jepkoech, J. *et al.* (2021) Arabica coffee leaf images dataset for coffee leaf
958 disease detection and classification. *Data Br.* 36, 107142

959 151. Manso, G.L. *et al.* (2019) A smartphone application to detection and classification
960 of coffee leaf miner and coffee leaf rust. *arXiv* 1904.00742

961 152. Parraga-Alava, J. *et al.* (2019) RoCoLe: a robusta coffee leaf images dataset for
962 evaluation of machine learning based methods in plant diseases recognition. *Data*
963 *Br.* 25, 104414

964 153. Uchiyama, H. *et al.* (2017) An easy-to-setup 3D phenotyping platform for
965 KOMATSUNA dataset. In *2017 IEEE International Conference on Computer*
966 *Vision Workshops (ICCVW)*, pp. 2038–2045

967 154. Seidel, D. *et al.* (2021) Predicting tree species from 3D laser scanning point
968 clouds using deep learning. *Front. Plant Sci.* 12, 1–12

969 155. Espejo-Garcia, B. *et al.* (2020) Towards weeds identification assistance through
970 transfer learning. *Comput. Electron. Agric.* 171, 105306

971 156. Olsen, A. *et al.* (2019) DeepWeeds: a multiclass weed species image dataset for
972 deep learning. *Sci. Rep.* 9, 2058

973 **Figure legends**

974 Figure 1. Phenomics platforms and sensors for high-throughput plant phenotyping in
975 controlled environments and field conditions: collecting relevant data from a wide range
976 of sources. A network of comprehensive automated weather stations collects hourly

977 weather and soil data, including, among others, rainfall, air temperature, solar radiation,
978 relative humidity, and soil moisture and temperature. Varying phenotyping scales allow
979 for precise and consistent monitoring of individual plants, plots, and fields. Ground-
980 based and aerial platforms can mount a variety of cameras and sensors for non-
981 invasive, high-throughput (HTP) phenotyping: visible light camera for RGB imaging;
982 LiDAR sensor and 3D laser scanners for 3D imaging; multispectral cameras and
983 hyperspectral sensors for spectral imaging; TIR cameras for thermal imaging; and
984 chlorophyll fluorescence sensor for chlorophyll fluorescence imaging. Automated and
985 environmentally controlled platforms, growth chambers, and multifunction printers can
986 be used for HTP in controlled environments. Root phenotyping in the field can be
987 invasive (e.g., shovelomics and its automation with root excavating robots); minimally
988 invasive (e.g., minirhizotrons); or non-invasive (e.g., ERT, electrical capacitance, GPR
989 mapping, and electromagnetic inductance mapping). Field deployable linear X-ray CT
990 cart, and handheld X-ray fluorescence elemental mapping are being explored for non-
991 invasive field root phenotyping. Multispectral, hyperspectral, RGB, and EIT imaging can
992 be used to phenotype roots in soil-filled rhizotrons (rhizoboxes) in controlled
993 environments. Similarly, NMR, X-ray CT, and PET imaging can be used to phenotype
994 roots in soil-filled pots. RhizoTubes, which are cylindrical rhizotrons, allow full
995 visualization of the root system of a single or up to six plants simultaneously. The
996 RhizoCab is designed to take images of the entire root systems of plants growing in
997 RhizoTubes. These platforms and sensing technologies are generating a massive
998 amount of data, which creates a need for proper data management and processing –
999 the first step of the data life cycle in digital phenomics (Figure 2). Abbreviations: EIT,

1000 electrical impedance tomography; ERT, electrical resistance tomography; GPR, ground
1001 penetrating radar; LiDAR, light detection and ranging; NMR, nuclear magnetic
1002 resonance; PET, positron emission tomography; RGB, red–green–blue; TIR, thermal
1003 infrared; UAVs, unmanned aerial vehicles; X-ray CT, X-ray computed tomography.

1004 Figure 2. Data architecture blueprint to drive human-centric explainable artificial
1005 intelligence (X-AI) innovation. Phenomics data can be structured, semistructured, or
1006 unstructured. Structured (e.g., spreadsheet files) data (blue line) are typically ‘at-rest’,
1007 transformed into rows and columns, and loaded into relational databases in data
1008 warehouses using a process known as ETL. Semistructured (e.g., extensible markup
1009 language files) and unstructured (e.g., flat files) data (red line) are streamed ‘in-motion’,
1010 loaded into non-relational databases, and stored in data lakes in their raw form; their
1011 transformation occurs on-demand using a process known as ELT. As ELT loads data
1012 immediately, it prevents any slowdown that often occurs at the transformation step, and
1013 thus, enables near real-time analytics for fast and practical decision-making. Whether
1014 ETL or ELT is used, data warehouses and data lakes store data as matrices, cubes,
1015 polytopes, or distributed in memory. A well-designed data architecture results in higher-
1016 quality phenomic datasets that allow plant scientists to ask biological questions and to
1017 devise data-driven analytics, searching for answers. Abbreviations: ELT, extract, load,
1018 transform; ETL, extract, transform, load.

1019 Figure 3. Technology infrastructure to support human-centric explainable artificial
1020 intelligence (X-AI). The technology infrastructure consists of the hardware, network flow,
1021 software frameworks, and programming languages that enable data transmission,
1022 transformation, storage, access, and analysis. (A) Computing hardware supporting data

1023 analysis. Pre-exascale supercomputers (e.g., University of Waterloo's Graham, and
1024 Lawrence Livermore National Laboratory's Sierra) reach a performance of a million
1025 billion FLOPS. With a similar hardware architecture but an increased number of CPUs
1026 and GPUs, exascale supercomputers (e.g., Oak Ridge National Laboratory's Summit
1027 and Frontier) reach a billion billion FLOPS and can deliver higher performance in
1028 pattern searching in phenomic big data, and thus, speeding up crop design. Quantum
1029 computers (e.g., International Business Machines' System One and Quantinuum's H1-
1030 2) represent a new paradigm in computation that leverages the fundamental principles
1031 of quantum mechanics to perform calculations. They employ **quantum bits** (qubits) that
1032 can be entangled, giving them the ability to manipulate vast amounts of data with few
1033 operations, and thus, the capacity to solve problems polynomially faster than classical
1034 computers (i.e., pre-exascale and exascale supercomputers) to ultimately design faster,
1035 better crops. Researchers can simulate quantum circuits on classical computers using
1036 free and open-source software development kits such as Cirq or Qiskit, and the
1037 cuQuantum software library to leverage the power of GPUs and parallel computing to
1038 perform faster calculations. Examples of classical and quantum computers are
1039 compared based on their peak performance that is the theoretical highest processing
1040 power they can reach. For classical computers, the LINPACK benchmark tests the
1041 performance in double precision (64-bit) compute capabilities while HPL-AI scores
1042 performance based on mixed precision (16- and 32-bit). As quantum computers use
1043 **QPUs** to manipulate the quantum states of qubits to perform computations, their
1044 performance is measured using **QV**. (B) Network flow to enable high-throughput access
1045 to and sharing of phenomic datasets. Requests coming from the wide area network are

1046 forwarded through a router to one of two paths: (i) the data query and browse path (red
1047 line) where requests to browse or search phenomic datasets are filtered through a
1048 firewall and processed by the hosting server; and (ii) the data transfer path (green line)
1049 where requests to download or upload phenomic datasets are inspected in the DMZ for
1050 access control, and are forwarded to the transfer nodes (typically Linux servers) to
1051 reach the filesystem where data can be transformed before or after storage using ETL
1052 or ELT, respectively (see Figure 2). (C, D) Representative free and open-source
1053 software frameworks and their supported programming languages used to implement AI
1054 algorithms. Abbreviations: CPU, central processing unit; ELT, extract, load, transform;
1055 ETL, extract, transform, load; FLOPS, floating-point operations per second; GPU,
1056 graphics processing unit; HPL-AI, high performance LINPACK for accelerator
1057 introspection; LINPACK, linear equations software package; PB, petabyte; QPU,
1058 quantum processing unit; Qubit, quantum bit; QV, quantum volume.

1059 Figure 4. Artificial intelligence (AI) architecture design to unleash the power of human-
1060 centric explainable AI (X-AI). (A) Representative AI algorithms that are used for AI tasks
1061 in digital phenomics including classification and regression (supervised learning), and
1062 clustering and dimensionality reduction (unsupervised learning). Reinforcement learning
1063 algorithms can be applied to search optimal architecture designs and improve their
1064 performance. (B) Representative AI algorithm design approaches and methods, a
1065 higher level of abstraction that help scientists in their efforts to design and implement
1066 novel AI algorithms to answer complex biological questions. The knowledge-based AI
1067 approach represents human expert knowledge as a collection of rules to form a
1068 knowledge base that is applied to solve a specific problem. It offers a consistent answer

1069 for a repetitive problem and its decisions are explainable. It can be implemented using
1070 rule-based methods. The Data-driven AI approach discovers connections and
1071 correlations automatically in a large amount of data and learns a black box model. It can
1072 be implemented using various methods including CNN, ensemble, or statistical
1073 methods, among others. The Informed AI approach combines knowledge-based AI with
1074 data-driven AI by leveraging human knowledge with knowledge acquired from data to
1075 make faster, more accurate decisions. It can be implemented using ensemble or rule-
1076 based methods. Finally, X-AI approaches provide meaningful explanations of decisions
1077 made by X-AI models to humans through a decipherable decision-making process.
1078 They allow the monitoring of inputs and outputs with the purpose of verifying X-AI
1079 models' adherence to ethical and socio-legal values by: (i) opening the black box of
1080 data-driven or informed AI models using ensemble methods; or (ii) designing new,
1081 transparent glass box algorithms that are interpretable by design using ensemble or
1082 CNN methods. Abbreviations: CNN, convolutional neural network; DBSCAN, density-
1083 based spatial clustering of applications with noise; DNN, deep neural network; GAN,
1084 generative adversarial network; GMM, Gaussian mixture model; HMM, hidden Markov
1085 model; KNN, k-nearest neighbors; NN, neural network; PCA, principal component
1086 analysis; RNN, recurrent neural network; RF, random forest; SAE, sparse autoencoder;
1087 SARSA, state–action–reward–state–action; SSAE, stacked SAE; SVM, support vector
1088 machine.

1089 Figure 5. Cultivating conditions for explainable artificial intelligence (X-AI) to flourish in
1090 plant digital phenomics. Data preprocessing prepares input data for X-AI algorithms:
1091 descriptive data analysis provides statistical summaries about a dataset in order to spot

1092 anomalies; data annotation and standardization is done by labeling and adding relevant,
1093 structured information about the data such as its source and other details known as
1094 metadata; and feature engineering uses existing features to create new ones while
1095 feature selection extracts relevant features from the complete set of features in a
1096 dataset, increasing the predictive precision of learning algorithms. X-AI can be achieved
1097 by either opening the black box or designing a transparent glass box. X-AI can be
1098 interrogated to understand why a decision has been made, keeping human-in-the-loop
1099 (HITL) of such decision-making, and allowing a two-way transfer of knowledge where on
1100 the one hand, experts assist in the training of X-AI and on the other hand, explanations
1101 can be used to generate scientific hypotheses that can result in new discoveries. An X-
1102 AI that takes into account the requirements of all stakeholders interacting with it will
1103 drive successful adoption among agricultural technopreneurs, plant biologists,
1104 policymakers, and funders. This will help bridge the gap between science, policy,
1105 embedded ethics, and entrepreneurship, allowing for responsible TT, and leading to
1106 technological, regulatory, and social and ethical outcomes. Abbreviations: AI, artificial
1107 intelligence; CNN, convolutional neural network; DeconvNet, deconvolution network;
1108 DeepLift, deep learning important features; FAIR, findable, accessible, interoperable,
1109 reusable; IPP, intellectual property protection; LIME, local interpretable model-agnostic
1110 explanations; SHAP, Shapley additive explanation; TT, technology transfer.

1111 Figure 6. Planning, training, and interpreting an explainable artificial intelligence (X-AI)-
1112 based analysis in plant digital phenomics require careful consideration at each stage of
1113 the analysis. This figure sheds light on all the elements of designing such a workflow
1114 using cassava leaf disease classification task as an example. For data preparation, (i) a

1115 dataset shared on Kaggle by the AI lab at Makerere University was used for analysis;
1116 (ii) data cleaning was carried out to eliminate outliers and mislabeled images; (iii) the
1117 dataset was randomly split for training, validation, and testing; (iv) another shared
1118 version of the dataset with images cropped to leaf boundaries using a trained YOLO
1119 model was used to minimize noise in training images; and (v) the training dataset was
1120 augmented and balanced by oversampling, creating random transformations to image
1121 geometries. An alternative solution to oversampling is synthesizing leaf images;
1122 OpenCV can be used to segment leaves to train a deep convolutional generative
1123 adversarial network (DCGAN) to generate synthetic data. 'This looks like that'
1124 interpretable by design algorithm, implemented in Python and PyTorch was carefully
1125 chosen for the classification task; its training time was approximated and compared on
1126 different hardware, showing the advantages of GPUs over CPUs and exascale over
1127 pre-exascale supercomputers. However, increasing the number of GPUs comes at the
1128 price of increased network communication and input-output (I/O) operations to
1129 synchronize the model over cluster nodes. Such overheads can cause a delay in the
1130 training time. For example, while the algorithm is expected to complete 1000 training
1131 epochs in 31 hours using 26 Nvidia Tesla V100 GPUs, it is still expected to take an
1132 approximation of eight hours using the full power of Summit supercomputer (27,649
1133 GPUs). 'This looks like that' algorithm uses transfer learning to import convolutional
1134 layers from pre-trained models and during training, the prototype layer extracts parts of
1135 training images (prototypes) and learns a similarity metric between them; the final class
1136 prediction is based on the weighted sum of similarities between the input and
1137 prototypes. For some prototypes, the nearest image patches come from different

1138 classes, often corresponding to a background patch, and thus should be pruned. For
1139 interpretation, the model tries to find evidence for a test image to belong to a specific
1140 class, marking activated patches by bounding boxes. While heatmaps show which part
1141 of the image is similar to a prototype, the confusion matrix illustrates the percentage of
1142 images of a true class classified into the class indicated by the predicted class column,
1143 indicating an overall accuracy of 88.7% after 240 training epochs. Abbreviations: CPU,
1144 central processing unit; DenseNet, dense convolutional network; GPU, graphics
1145 processing unit; OpenCV, open source computer vision library; ResNet, residual
1146 network; VGG, visual geometry group; YOLO, you only look once.

1147 Table 1. Publicly available global datasets and their characteristics as valuable resources for plant digital phenomics research^{a,b}

Dataset	Country of origin	Plant species ^c	Plant organ systems	No. of images	Platform	Sensors	Image annotation types ^d	Potential applications	Access type ^e	Data access details	Data PID ^f	Link to dataset	File format ^g	Refs.
Supporting data for "deep machine learning provides state-of-the-art performance in image-based plant phenotyping"	UK	<i>Triticum aestivum</i>	Root	2697	Controlled environment stationary platform	RGB	Point annotations of root tips, leaf tips, leaf bases, ear tips, ear bases	Identification and localization of root tips, leaf, and wheat ear tips	OA	Downloadable tar file	http://doi.org/10.5524/100343	http://gigadb.org/dataset/100343	JPG	[109]
			Shoot	1664										
Wheat 2017	UK	<i>T. aestivum</i>	Shoot, spikes	520	Controlled environment stationary platform	RGB	Point annotations of spikelet, base and tip of each ear	Localization and counting of wheat spikes and spikelets	OA	Register for link to download zip file	–	https://plantimages.nottingham.ac.uk/	JPG	[110]
Global wheat head detection (GWHD)	Japan, France, Canada, UK, Switzerland, China, Australia	Wheat	Shoot	1094	Manned mobile platform	RGB	Bounding boxes	Detection and localization of wheat heads	OA	Downloadable zip file	http://doi.org/10.5281/zenodo.4298502	https://zenodo.org/record/4298502	PNG	[111]
				678	Handheld visible light camera in the field									
				447	Rail-based field automated gantry									
Global wheat head detection (GWHD) 2021	Japan, France, Canada, UK, Switzerland, China, Australia, USA, Mexico, Republic of Sudan, Norway, Belgium	Wheat	Shoot	2307	Manned mobile platform	RGB	Bounding boxes	Detection and localization of wheat heads	OA	Downloadable zip file	https://doi.org/10.5281/zenodo.5092309	https://zenodo.org/record/5092309	PNG	[112]
				2684	Handheld visible light camera in the field									
				1429	Rail-based field automated gantry									

Supporting data for "high throughput phenotyping with deep learning gives insight into the genetic architecture of flowering time in wheat"	KS - USA	Wheat	Shoot	>400000	Ground-based field robot	RGB	Image-level annotations	Estimation of plant morphology and developmental stages	OA	Downloadable zip file	http://dx.doi.org/10.5524/100566	http://gigadb.org/dataset/view/id/100566	JPEG	[113]
RootNav 2.0	UK	<i>T. aestivum</i>	Root	3630	Controlled environment stationary platform	RGB	Segmentation masks, image-level annotations	Root segmentation and species classification	OA	Download each image separately	http://doi.org/10.5524/100651	http://gigadb.org/dataset/100651	JPG	[15]
		<i>Brassica napus</i>		120		NIR							PNG	
		<i>Arabidopsis thaliana</i>		277										
Cosegmentation for plant phenotyping (CosegPP)	USA	Buckwheat	Shoot	56	Controlled environment stationary platform	RGB	Segmentation masks	Plant segmentation	OA	Downloadable zip file	https://doi.org/10.5281/zenodo.5117176	https://zenodo.org/record/5117176	PNG	[114]
				56		IR								
				56		CF								
		Sunflower		104		RGB								
				112		IR								
				112		CF								
Aberystwyth leaf evaluation dataset	UK	<i>A. thaliana</i>	Shoot	56	Controlled environment stationary platform	RGB	Semantic segmentation	Plant and leaf segmentation	OA	Downloadable zip file	http://doi.org/10.5281/zenodo.168158	https://zenodo.org/record/168158	PNG	
Deep phenotyping dataset	Australia	<i>A. thaliana</i>	Shoot	2134	Controlled environment stationary platform	RGB	Image-level annotations	Genotype classification	OA	Downloadable zip file	–	https://figshare.com/s/e18a978267675059578f	JPG	[115]
Supporting data for "ChronoRoot: high-throughput phenotyping by deep segmentation networks reveals novel temporal parameters of plant root system architecture"	France	<i>A. thaliana</i>	Root	331	Controlled environment stationary platform	RGB	Segmentation masks	Root segmentation	OA	Downloadable tar file	http://dx.doi.org/10.5524/100911	http://dx.doi.org/10.5524/100911	PNG	[16]
Plant phenotyping datasets	Italy	<i>A. thaliana</i>	Shoot	6287	GARNICS controlled environment robot gardener	RGB	Segmentation masks, bounding boxes, point annotations of	Plant and leaf segmentation, leaf counting, species classification	OA	Fill in a form and get access and download a zip file	–	https://www.plant-phenotyping.org/datasets-	HDF5	[116]

	Germany	<i>Nicotiana tabacum</i>		165120			leaf centers					home		
Multi-modality plant imagery database (MSU-PID)	MI - USA	<i>A. thaliana</i>	Shoot	576	Controlled environment stationary platform	RGB, RGB-depth, CF, NIR ^h	Polygon and point annotations of leaves and leaf tips	Leaf segmentation, counting, alignment, and tracking	OA	Downloadable zip file	–	http://cvlab.cse.msu.edu/multi-modality-imagery-database-msu-pid.html	PNG	[117]
		<i>Phaseolus vulgaris</i>		175										
Plant segmentation	Hungary	<i>A. thaliana</i>	Shoot	16 ⁱ	Computer scanner	RGB	Segmentation masks	Segmentation and length determination of hypocotyl	OA	Downloadable zip file	–	https://www.kaggle.com/tivadardanka/plant-segmentation	PNG	[118]
		<i>Brachypodium distachyon</i>		8 ⁱ	Handheld visible light camera in a controlled environment									
		<i>Sinapis alba</i>	Shoot and root	15 ⁱ										
Eschikon plant stress phenotyping dataset	Switzerland	<i>Beta vulgaris</i>	Shoot	496	Controlled environment stationary platform	HS (NIR)	Image-level annotations, bounding boxes	Classification of biotic and abiotic stress	OA	Downloadable zip file	–	https://projects.asl.ethz.ch/datasets/doku.php?id=2018plantstressphenotyping	PNG	[119]
				992		RGB-depth								
				496		RGB								
Remote sensing 2018 weed map dataset	Switzerland, Germany	<i>B. vulgaris</i> , weed	Shoot	18746	UAV	MS (NIR), 4 and 5 bands ^l	Semantic segmentation	Identification and segmentation of crops and weeds	OA	Downloadable zip file	–	https://projects.asl.ethz.ch/datasets/doku.php?id=weedmapremotesensing2018weedmap	TIF	[120]
Sugar beets 2016	Germany	Sugar beet, weed	Shoot	300	Ground-based field robot	RGB, RGB-depth, MS (NIR), 4 bands ^h	Semantic segmentation	Identification and segmentation of crops and weeds	OA	Download each image separately	–	https://www.ipb.uni-bonn.de/datasets_IJRR2017/	PNG	[121]
Images of soybean leaves	Brazil	Soybean	An intact leaf, shoot	6410	Handheld visible light camera in the field, UAV	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Downloadable zip file	https://doi.org/10.17632/bycbh73438.1	https://data.mendeley.com/datasets/bycbh73438/1	JPG	[122]
Data for: weed detection in soybean crops using ConvNets	Brazil	Soybean, weed	Shoot	400	UAV	RGB	Segmentation masks	Identification and segmentation of crops and weeds	OA	Downloadable zip file	https://doi.org/10.17632/3fmjm7ncc6.2	https://data.mendeley.com/datasets/3fmjm7ncc6/2	JPEG	[123]
Crop vs weed discrimination dataset	UK	Onion, weed	Shoot	20 ^k	Manned mobile platform	RGB, MS (NIR), 2 bands ^h	Semantic segmentation	Identification and segmentation of crops and weeds	OA	Downloadable zip file	–	https://lcas.lincoln.ac.uk/wp/research/datasets-software/crop-vs-weed-discrimination-dataset/	PNG	[124]
		Carrot, weed		20 ^k										
Crop/weed field image dataset (CWFID)	Germany	Carrot, weed	Shoot	60	Ground-based field robot	MS (NIR), 2 bands	Segmentation masks, semantic segmentation	Identification and segmentation of crops and weeds	OA	Downloadable zip file	–	https://github.com/cwfid/dataset	PNG	[125]

Table 1. (continued)

Cassava leaf disease classification	Uganda	<i>Manihot esculenta</i>	Shoot	9436	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Create Kaggle account to download zip file	–	https://www.kaggle.com/c/cassava-disease/data	JPG	[96]
Cassava disease classification		Cassava		21397								https://www.kaggle.com/c/cassava-leaf-disease-classification/data		
Cassava root cross-section images	Uganda, Tanzania	Cassava	Root	10052	Handheld visible light camera in the field	RGB	Semantic segmentation	Quantification of root damage	OA	Downloadable zip file	https://doi.org/10.17632/gvp7vshvnh.3	https://data.mendeley.com/datasets/gvp7vshvnh/3	JPG	[126]
Plant pathology 2020 - FGVC7	NY - USA	Apple	An intact leaf	3651	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Create Kaggle account to download zip file	–	https://www.kaggle.com/c/plant-pathology-2020-fgvc7/data	JPG	[127]
Plant pathology 2021 – FGVC8				23000								https://www.kaggle.com/c/plant-pathology-2021-fgvc8		
LFuji-air dataset	Spain	<i>Malus domestica</i>	Shoot	88	Ground-based field robot	LiDAR	Bounding boxes	Identification and localization of fruits, estimation of yield, canopy geometric characterization	OA	Downloadable zip file	–	https://repositori.udl.cat/handle/10459.1/68782	MAT	[27.128]
MinneApple	MN - USA	Apple	Shoot	1000	Handheld visible light camera in the field	RGB	Polygon annotations	Identification, segmentation, and counting of fruits	OA	Downloadable tar file	http://doi.org/10.13020/8ecp-3r13	https://conservancy.umn.edu/handle/11299/206575	JPG	[129]
Data from: multi-species fruit flower detection using a refined semantic segmentation network	WV - USA	Apple	Shoot	18'	Manned mobile platform	RGB	Segmentation masks, image-level annotations	Classification of species from fruit flowers	OA	Downloadable zip file	http://doi.org/10.15482/USDA.ADC/1423466	https://data.nal.usda.gov/dataset/data-multi-species-fruit-flower-detection-using-refined-semantic-segmentation-network	JPG	[130]
		Peach		24'	Handheld visible light camera in the field									
		Pear		18'										
DiaMOS plant dataset: a dataset for diagnosis and monitoring plant disease	Italy	Pear	Shoot	3505	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Downloadable zip file	https://doi.org/10.5281/zenodo.5557313	https://zenodo.org/record/5557313#_Yv5JrXZByMo	JPG	[131]
PlantaeK: A leaf database of native plants of Jammu and Kashmir	Jammu, Kashmir	Apple	A single detached leaf	351	Handheld visible light camera in a controlled environment	RGB	Image-level annotations	Identification of biotic stress, classification of plant species	OA	Downloadable zip file	http://doi.org/10.17632/t6j2h22jpx.1	https://data.mendeley.com/datasets/t6j2h22jpx/1	JPG	
		Apricot		270										
		Cherry		212										
		Cranberry		212										

		Grapevine		171										
		Peach		331										
		Pear		228										
		Walnut		378										
Data for: identification of plant leaf diseases using a 9-layer deep convolutional neural network	FL - PA - NY - USA	Multiple crops ¹	A single detached leaf	54305	Handheld visible light camera in a controlled environment	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Downloadable zip file	http://doi.org/10.17632/tywbtsjrjv.1	https://data.mendeley.com/datasets/tywbtsjrjv/1	JPG	[13]
ACFR orchard fruit dataset	Australia	Apple	Shoot	1120	Ground-based field robot	RGB	Bounding boxes, circle annotations	Identification and classification of fruits	OA	Downloadable zip file	–	http://data.acfr.usyd.edu.au/ag/treecrops/2016-multifruit/	PNG	[132]
		Almond		620										
		Mango		1964										
Dataset for pest classification in Mango farms	Indonesia	Mango	An intact leaf	510	Handheld visible light camera in the field	RGB	Image-level annotations	Classification of pest	OA	Downloadable zip file	https://doi.org/10.17632/94jf97jzc8.1	https://data.mendeley.com/datasets/94jf97jzc8/1	JPG	[133]
Berries in vineyards-color (BIVcolor)	Germany	<i>Vitis vinifera</i>	Shoot	500	Ground-based field robot	RGB	Image-level annotations	Detection of grape size and color	OA	Downloadable zip file	http://doi.org/10.5073/jki-data.2015.1	https://www.openagrar.de/receive/openagrar_mods_00021925	TIF	[134]
Supporting data for "a novel ground truth multispectral image dataset with weight, anthocyanins and Brix index measures of grape berries tested for its utility in machine learning pipelines"	Spain	<i>V. vinifera</i>	Fruit	1283	Controlled environment stationary platform	MS	Image-level annotations	Prediction of grape variety	OA	Downloadable zip file	http://dx.doi.org/10.524/102220	http://gigadb.org/dataset/102220	TIF	[135]
ESCA-dataset	Italy	Grapevine	Shoot	1770	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Downloadable zip file	http://doi.org/10.17632/89cncx58kj.1	https://data.mendeley.com/datasets/89cncx58kj/1	JPG	[136]
An annotated image dataset of downy mildew symptoms on Merlot grape variety	France	Grapevine	Shoot	99	Ground-based field robot	RGB	Semantic segmentation	Identification of downy mildew	OA	Downloadable zip file	–	https://ars.els-cdn.com/content/image/1-s2.0-S2352340921005345-mm1.zip	JPEG	[137]
The MangoNet semantic dataset	India	<i>Mangifera indica</i>	Shoot	49 ¹	Handheld visible light camera in the field	RGB	Semantic segmentation	Identification and counting of fruits	OA	Downloadable zip file	–	https://github.com/avadesh02/MangoNet-Semantic-Dataset	JPG	[138]

Leaves: India's most famous basil plant leaves quality dataset	India	Basil	Shoot	1131	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Download each image separately	https://dx.doi.org/10.21227/a4f6-4413	https://iee-dataport.org/open-access/leaves-india%E2%80%99s-most-famous-basil-plant-leaves-quality-dataset	JPG	
A database of leaf images: practice towards plant conservation with plant pathology	Jammu, Kashmir	Multiple species"	A single detached leaf	4503	Handheld visible light camera in a controlled environment	RGB	Image-level annotations	Identification of biotic stress	OA	Downloadable zip file	http://doi.org/10.17632/hb74ynkjc4	https://data.mendeley.com/datasets/hb74ynkjc4	JPG	[139]
A citrus fruits and leaves dataset for detection and classification of citrus diseases through machine learning	Pakistan	Citrus	Shoot	150	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Downloadable zip file	http://doi.org/10.17632/3f83gxm57.2	https://data.mendeley.com/datasets/3f83gxm57/2	JPG	[140,
			A single detached leaf	609										141]
Supporting data for "morphometric analysis of Passiflora leaves: the relationship between landmarks of the vasculature and elliptical Fourier descriptors of the blade"	Brazil	40 <i>Passiflora</i> species"	A single detached leaf	5767	Multifunction printer	RGB	Point annotations of leaf edges	Classification of species	OA	Download each image separately	http://doi.org/10.5524/100251	http://gigadb.org/dataset/100251	TIF	[142, 143]
QuinceSet: dataset of annotated Japanese quince images for object detection	Latvia	<i>Chaenomeles japonica</i>	Shoot	1515	Handheld visible light camera in the field	RGB	Image-level annotations, bounding boxes	Identification and localization of fruits	OA	Downloadable zip file	https://doi.org/10.5281/zenodo.6402251	https://zenodo.org/record/6402251	JPG	[144]
Thermal images - diseased & healthy leaves	India	<i>Oryza sativa</i>	An intact leaf	636	Handheld TIR camera in the field	TIR	Image-level annotations	Classification of biotic stress	OA	Create Kaggle account to download zip file	–	https://www.kaggle.com/sujaradha/thermal-images-diseased-healthy-leaves-paddy	JPG	
Date fruit dataset	Saudi Arabia	<i>Phoenix dactylifera</i>	Shoot	8079	Handheld visible light camera in the	RGB	Image-level annotations	Detection and maturity classification of	OA	Create IEEE DataPort account to download zip	http://doi.org/10.21227/x46j-	https://iee-dataport.org/open-access/date-	JPG	[145, 146]

			Fruit bunch	152	field			fruits		file	sk98	fruit-dataset-automated-harvesting-and-visual-yield-estimation		
			Date fruit	256										
Image set for deep learning: field images of maize annotated with disease symptoms	NY - USA	<i>Zea mays</i>	Shoot	7669	UAV	RGB	Line and spline annotation of lesions	Identification of northern leaf blight infected plants	OA	Downloadable zip file	–	https://osf.io/p67rz	JPG	[147]
				10533	Handheld visible light camera in the field									
Vegetable crops dataset for proximal sensing	France	<i>Z. mays</i>	Shoot	1065	Manned mobile platform	RGB	Image-level annotations, bounding boxes, point annotations	Identification and localization of crop stems	OA	Downloadable zip file	https://doi.org/10.17632/d7kbzjr83k.1	https://data.mendeley.com/datasets/d7kbzjr83k/1	JPG	[148]
		<i>P. vulgaris</i>		779										
		<i>Allium ampeloprasum</i>		601										
Pheno4D	Germany	Maize	Shoot	84	Controlled environment stationary platform	3D laser scanner	Semantic segmentation	Estimation of plant traits, growth analysis	OA	Downloadable zip file	–	https://www.ipb.uni-bonn.de/data/pheno4d	XYZ	[149]
		Tomato		140										
JMuBEN	Kenya	Arabica coffee	A single detached leaf	22591	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Downloadable zip file	https://doi.org/10.17632/t2r6rszp5c.1	https://data.mendeley.com/datasets/t2r6rszp5c/1	JPG	[150]
JMuBEN2				35964										
BRACOL - A Brazilian Arabica coffee leaf images dataset	Brazil	Arabica coffee	A single detached leaf	1747	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Downloadable zip file	https://doi.org/10.17632/yy2k5y8mxg.1	https://data.mendeley.com/datasets/yy2k5y8mxg/1	JPG	[151]
RoCoLe: a robusta coffee leaf images dataset	Ecuador	<i>Robusta</i> coffee	An intact leaf	1560	Handheld visible light camera in the field	RGB	Image-level annotations	Identification and classification of biotic stress	OA	Downloadable zip file	https://doi.org/10.17632/c5yvn32dzg.2	https://data.mendeley.com/datasets/c5yvn32dzg/2	JPG	[152]
KOMATSUNA dataset for instance segmentation, tracking and reconstruction	Japan	<i>B. rapa</i>	Shoot	180	Controlled environment stationary platform	RGB	Semantic segmentation	Leaf segmentation and plant growth measurement	OA	Downloadable zip file	–	https://limu.ait.kyushu-u.ac.jp/~agri/komatsuna	PNG	[153]
				60		RGB-depth								
Single tree point clouds from terrestrial laser scanning	Germany, OR - USA	<i>Quercus petraea</i>	Shoot	22 ^k	Terrestrial laser scanner	3D laser scanner	Image-level annotations	Classification of tree species	OA	Downloadable zip file	https://doi.org/10.25625/FOHUJM	https://data.goe.ttingen-research-online.de/dataset.xhtml?persistentId=doi:10.25625/FOHUJM	XYZ ^o	[154]
		<i>Fraxinus excelsior</i>		39 ^k										
		<i>Picea abies</i>		158										
		<i>Pinus sylvestris</i>		25 ^k										
		<i>Q. rubra</i>		100										

		<i>Fagus sylvatica</i>		163										
		<i>Pseudotsuga menziesii</i>		183										
Early-crop-weed	Greece	<i>S. lycopersicum</i>	Shoot	202	Handheld visible light camera in the field	RGB	Image-level annotations	Identification of crops and weeds	OA	Downloadable zip file	–	https://github.com/AUAGroup/early-crop-weed	JPG	[155]
		<i>Gossypium hirsutum</i>		48 ^k										
		<i>S. nigrum</i>		130										
		<i>Abutilon theophrasti</i>		124										
DeepWeeds	Australia	Multiple species ^p	Shoot	8403	Ground-based field robot	RGB	Image-level annotations	Classification of weed species	OA	Downloadable zip file	–	https://github.com/AlexOlsen/DeepWeeds	JPG	[156]

1148 ^aAbbreviations: ACFR, Australian center for field robotics; CF, chlorophyll fluorescence; FGVC7/8, the seventh/eight workshop on fine-grained visual categorization; GARNICS, gardening with
1149 a cognitive system; HDF5, hierarchical data format version 5; HS, hyperspectral; JMuBEN, Jepkoech, Mugo and Benson; JPG, joint photographic experts group; LiDAR, light detection and
1150 ranging; MAT, Matlab; MS, multispectral; MSU-PID, Michigan State University-plant imagery database; NIR, near infrared; OA, open access; PID, persistent identifier; RGB, red–green–blue;
1151 TIF, tag image file format; TIR, thermal infrared; UAV, unmanned aerial vehicle.

1152 ^bWe identified 56 publicly available global datasets and their characteristics using Google’s search engine and Google Dataset Search. The search combined terms describing various plant
1153 organ systems, sensors, artificial intelligence (AI) techniques, as well as dataset and database. All pages for each search were systematically collated and screened. Additional datasets are
1154 available in repositories containing large amounts of OA imaging data. Repositories such as the National Ecological Observatory Network ([https://data.neonscience.org/data-](https://data.neonscience.org/data-products/explore)
1155 [products/explore](http://leafsnap.com)), Leafsnap (<http://leafsnap.com>), the Institut National de la Recherche Agronomique (<https://data.inrae.fr>), the United States Geological Survey
1156 (<https://www.usgs.gov/products/data-and-tools/science-datasets>), the National Aeronautics and Space Administration earth science data (<https://earthdata.nasa.gov>), the plant genomics and
1157 phenomics research data repository (<https://edal-pgp.ipk-gatersleben.de>), the computer vision and biosystems signal processing group (<https://vision.eng.au.dk/data-sets>), the Transportation
1158 Energy Resources from Renewable Agriculture Phenotyping Reference Platform (<https://terraref.ncsa.illinois.edu/clowder>), figshare (<https://figshare.com>), Dryad (<http://datadryad.org>), the
1159 International Maize and Wheat Improvement Center (<https://data.cimmyt.org>) and the *Arabidopsis thaliana* phenotyping database (Phenopsis DB, <http://bioweb.supagro.inra.fr/phenopsis>)
1160 provide datasets in downloadable zip files. Similarly, the Oak Ridge National Laboratory Distributed Active Archive (https://daac.ornl.gov/get_data/#themes) provides datasets in
1161 downloadable zip files after registering for an account, as well as the University of Nebraska-Lincoln Plant Vision Initiative (<https://plantvision.unl.edu/dataset>) and the X-Plant ([http://www.x-](http://www.x-plant.org)
1162 [plant.org](http://www.x-plant.org)) after filling a form. Other sources such as the online database for plant image analysis software tools (<https://www.quantitative-plant.org/dataset>) and the registry of research data
1163 repositories (<https://www.re3data.org>) are designed specifically for the discovery of datasets in various repositories.

1164 ^cPlant species were reported whenever they were available in the corresponding referenced paper(s); common names were reported otherwise.

1165 ^dImages of datasets with semantic segmentation annotations are completely annotated images, where a class is assigned to each pixel.

1166 ^eNo datasets were excluded on the basis of access type (i.e., OA, data available on request, or OA with barriers – datasets fulfilling criteria for OA but being inaccessible because of
1167 unpredictable reasons such as broken hyperlinks).

1168 ^fData PID is a long-lasting digital reference to a dataset, such as a digital object identifier (DOI). A dash (–) indicates that no PIDs are available. DOIs for datasets can be issued automatically
1169 by the hosting repositories (e.g., Zenodo, GigaDB, Mendeley Data, and IEEE DataPorts). As datasets should be cited to ensure credit to those who produced and curated them, we
1170 recommend that they should include a PID and the minimum metadata suggested by DataCite (a non-profit membership organization that provides DOIs for research data) and FORCE11 (a
1171 community of scholars, librarians, archivists, publishers and research funders), i.e., author, year, title, and repository. Data producers can be inferred based on the author contributions of the
1172 corresponding referenced paper(s) while data curators can be inferred based on the author(s) that published the dataset to a repository.

1173 ^gFile format defines the structure and encoding of the data stored in it and thus guides researchers on how to programmatically input such data to their AI algorithms.

1174 ^hThe same number of images was taken with each sensor.

1175 ⁱDatasets containing more than one object per image (e.g., multiple hypocotyls, fruits, flowers). When segmented, each image could become hundreds of samples to train an AI algorithm.

1176 ^jTwo multispectral cameras were used: a five-band RedEdge-M camera in Germany and a four-band Sequoia camera in Switzerland.

1177 ^kFor datasets with small image number, transfer learning can be applied, giving an AI model a warm start by applying information learned from another previously trained model.

1178 ^lCrops and their corresponding number of images: Apple, 3171; blueberry, 1502; cherry, 1906; corn, 3852; grapevine, 4062; orange, 5507; peach, 2657; pepper, 2475; potato, 2152;
1179 raspberry, 371; soybean, 6925; strawberry, 1565; tomato, 18160.

1180 ^mSpecies and their corresponding number of images: *M. indica*, 435; *Terminalia Arjuna*, 452; *Alstonia Scholaris*, 433; *Psidium guajava*, 419; *Aegle marmelos*, 118; *Syzygium cumini*, 624;
1181 *Jatropha curcas*, 257; *Pongamia Pinnata*, 598; *Ocimum basilicum*, 149; *Punica granatum*, 559; *Platanus orientalis*, 223; *C. limon*, 236.

1182 ⁿ*Passiflora* species and their corresponding number of images: *P. coriacea*, 208; *P. misera*, 215; *P. biflora*, 105; *P. capsularis*, 118; *P. micropetala*, 68; *P. organensis*, 84; *P. pohlii*, 16; *P.*
1183 *rubra*, 87; *P. tricuspis*, 257; *P. caerulea*, 99; *P. cincinnata*, 84; *P. edmundoi*, 111; *P. gibertii*, 192; *P. hatschbachii*, 132; *P. kermesina*, 113; *P. mollissima*, 69; *P. setacea*, 189; *P. suberosa*,
1184 352; *P. tenuifila*, 113; *P. amethystina*, 119; *P. foetida*, 304; *P. gracilis*, 81; *P. morifolia*, 57; *P. actinia*, 95; *P. miersii*, 133; *P. sidifolia*, 145; *P. triloba*, 295; *P. alata*, 235; *P. edulis*, 119; *P.*
1185 *ligularis*, 139; *P. nitida*, 62; *P. racemosa*, 194; *P. villosa*, 58; *P. coccinea*, 169; *P. cristalina*, 220; *P. galbana*, 161; *P. malacophylla*, 168; *P. maliformis*, 156; *P. miniata*, 129; *P. mucronata*, 116.
1186 ^oA point cloud data file in XYZ format contains rows of data, each consisting of x, y, and z coordinates of a point.
1187 ^pSpecies and their corresponding number of images: *Ziziphus mauritiana*, 1125; *Lantana camara*, 1064; *Parkinsonia aculeata*, 1031; *Parthenium hysterophorus*, 1022; *Vachellia nilotica*,
1188 1062; *Cryptostegia grandiflora*, 1009; *Chromolaena odorata*, 1074; *Stachytarpheta spp.*, 1016.

1189 **Glossary**

1190 **Bias:** systematic errors in the ability of AI models to make correct predictions.

1191 **Compute node:** a backend node used for computing in a cluster and reached via a
1192 frontend node.

1193 **Computer cluster:** a group of interconnected computers working together as a single,
1194 integrated computing resource.

1195 **Confusion matrix:** a visual representation that describes the complete performance of
1196 an AI model, summarizing its predictions in four categories: true-positives, true-
1197 negatives, false-positives, and false-negatives.

1198 **Crowdsourced:** the act of collecting data by soliciting contributions from a large group
1199 of people rather than from traditional experiments.

1200 **Explicit knowledge:** the human knowledge that can be readily assembled and passed
1201 on by written or verbal instruction. Metadata is explicit knowledge about data.

1202 **Federated learning (FL):** a collaborative AI training paradigm in which copies of a
1203 model are distributed to devices, where data is stored, for local training, and the
1204 resulting model weights, rather than the data, are sent back to a central server to
1205 update the main model.

1206 **GPU-accelerated:** a backend node used mainly for accelerating computing, connected
1207 in a heterogeneous manner in a computer cluster.

1208 **Human-in-the-loop (HITL):** an approach that aims to achieve what neither a human
1209 nor a machine can do on their own; it leverages a continuous feedback loop between
1210 them to train, evaluate, and deploy AI models that continuously learn and improve their
1211 prediction accuracy.

1212 **Hyperparameters:** a group of variables whose values cannot be estimated from data
1213 and are manually tweaked to determine the optimal configuration to train a specific
1214 model (e.g., learning rate, batch size, number of training epochs).

1215 **Notebook:** a web browser-based interactive computing environment that can be used
1216 to combine software code, computational output, explanatory text, and multimedia
1217 resources in a single document.

1218 **Parallel computing:** a form of computation in which multiple compute nodes operating
1219 simultaneously are used to solve a large problem broken into independent smaller parts
1220 that can be processed concurrently.

1221 **Quantum bit (qubit):** the quantum analogue of a classical bit; it may adopt the states 0,
1222 1, or any possible combination of both states.

1223 **Quantum processing unit (QPU):** a computational unit that leverages quantum
1224 mechanical phenomena to manipulate information, relying on qubits.

1225 **Quantum volume (QV):** a metric that measures the performance of a quantum
1226 computer taking into account its number of qubits and error rates.

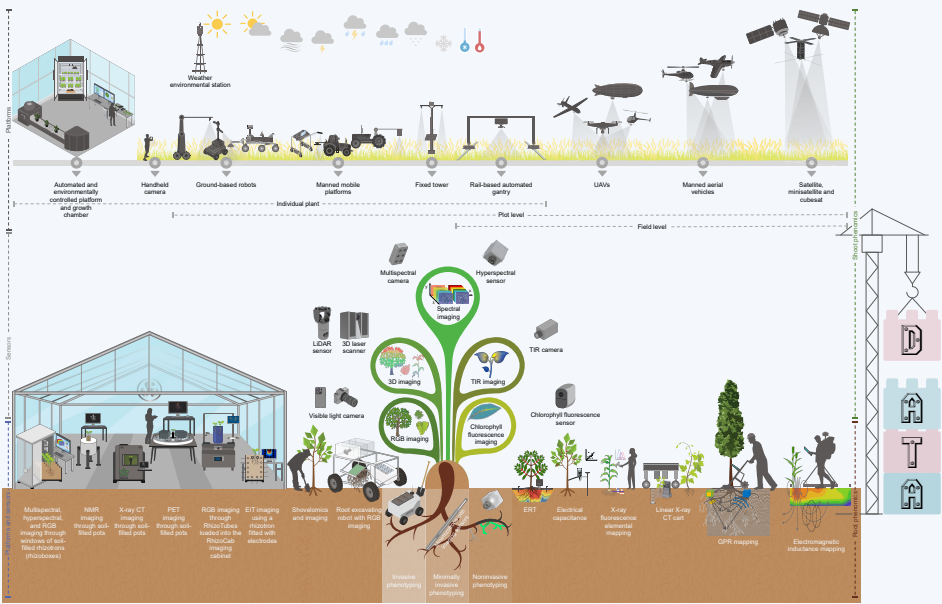
1227 **Synthetic data:** data generated artificially using AI algorithms when real data cannot be
1228 collected in sufficient amounts.

1229 **Tacit knowledge:** the know-how, skills, and intuition that live in the individual's
1230 experiences and are hard to impart or transfer to others. It can be shared through
1231 advances in information and communications technology, and thus becomes explicit.

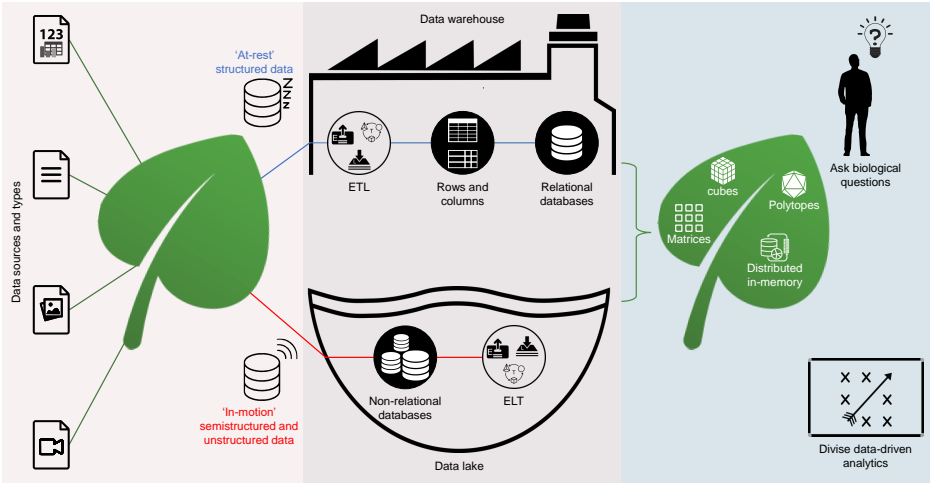
1232 **Transfer learning (TL):** a technique in which an AI algorithm reuses parts of a
1233 previously trained model on a new model to perform a different but similar task.

1234 **Trustworthiness:** a quality of an AI model working reliably in ways that anyone can
1235 trust; it should be (i) lawful, ensuring compliance with all applicable laws and
1236 regulations; (ii) ethical, demonstrating adherence to ethical principles and values, (iii)
1237 robust, able to deal with bias during all of its lifecycle; and (iv) explainable.

1238 **Unbalanced dataset:** a dataset having certain classes contain substantially more
1239 training examples than other classes, misleading the classifier algorithm to overlearn
1240 the majority classes and to perform poorly in the prediction of the minority classes.



Data architecture



Data sources and types

'At-rest' structured data

'In-motion' semistructured and unstructured data

Data warehouse

ETL

Rows and columns

Relational databases

Data lake

Non-relational databases

ELT

cubes

Matrices

Polytopes











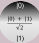






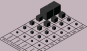











Distributed in-memory

Ask biological questions

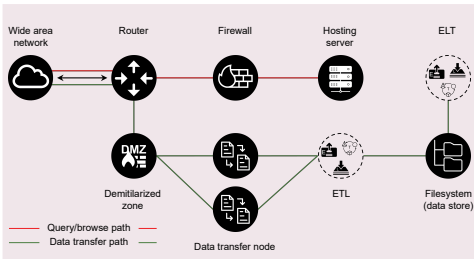
Divise data-driven analytics

Technology infrastructure

(A) Hardware

	Pre-exascale supercomputer	Exascale supercomputer	Quantum computer
Features	 Classical bits (0 or 1)  CPUs (e.g., Graham: 2,532; Sierra: 8,640)  GPUs (e.g., Graham: 520; Sierra: 17,280)  Storage capacity (e.g., Graham: 50 PB; Sierra: 154 PB)  Floor space (e.g., Graham: = 160 m ² ; Sierra: = 650 m ²)	 Classical bits (0 or 1)  CPUs (e.g., Summit: 9,216; Frontier: 9,408)  GPUs (e.g., Summit: 27,648; Frontier: 37,632)  Storage capacity (e.g., Summit: 250 PB; Frontier: 700 PB)  Floor space (e.g., Summit: = 520 m ² ; Frontier: = 372 m ²)	 Qubit (0, 1, or a superposition of 0 and 1)  Quantum entanglement  Floor space (e.g., System One and H1-2: room-sized computers)
Performance measures	 A million billion (10 ¹⁹) FLOPS  Double precision LINPACK benchmark (e.g., Graham: 2.6 petaFLOPS; Sierra: 125 petaFLOPS)	 A billion billion (10 ²¹) FLOPS  Mixed precision HPL-AI benchmark (e.g., Summit: 1.4 exaFLOPS; Frontier: 6.86 exaFLOPS)	 QV (e.g., System One: 32; H1-2: 4,096)
Benefits	 Pattern searching in big phenomics data  Crop design  Leveraging AI abilities for digital phenomics	 Faster pattern searching in big phenomics data  Speeding up crop design  Boosting AI abilities for digital phenomics	 Searching big phenomics data to uncover patterns in seconds; Accessing all items in a database at the same time in seconds  Designing faster, better crops  Supercharging AI abilities for digital phenomics  Debugging millions of lines of software code  Providing sustainability-friendly solutions

(B) Network flow



(C) Software frameworks

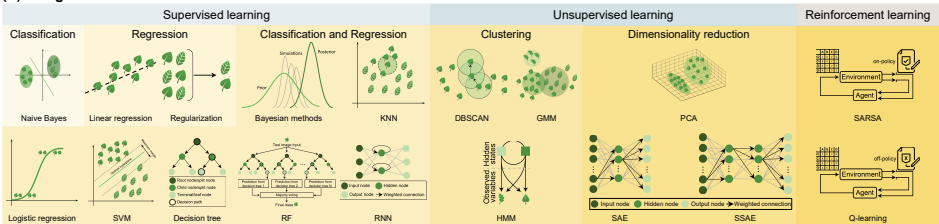
Tensorflow Keras
PyTorch
scikit-learn
XGBoost
Apache MXNet

(D) Programming languages

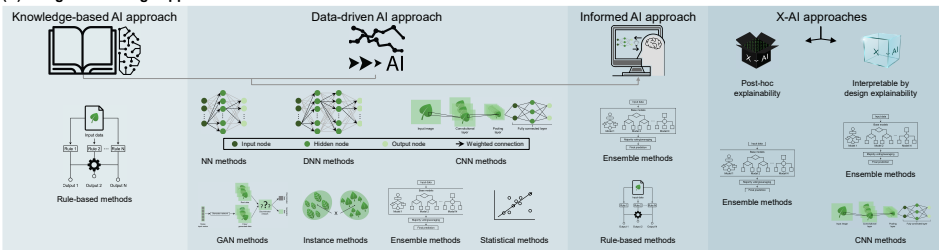
Python, JavaScript, C++, Java
Python, C++, Java
Python
Python, C, C++, Java, Julia, R, Ruby, Swift
Python, C++, Java, Scala, Julia, Clojure, R, and Perl

AI architecture design

(A) AI algorithms

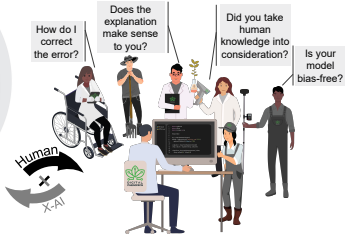
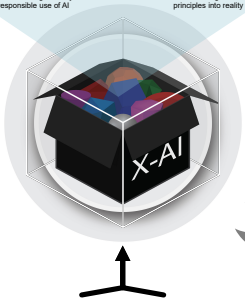
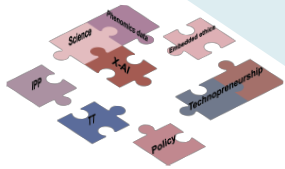


(B) AI algorithm design approaches and methods



What it means to look inside the X-AI box – The seven outcomes of X-AI

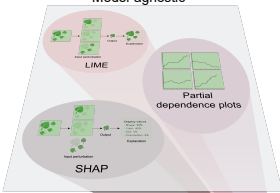
Outcomes



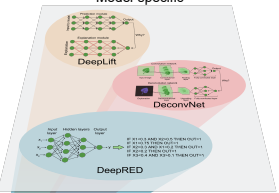
Monitoring inputs, outputs, and outcomes in plant digital phenomics

Post-hoc explainable models

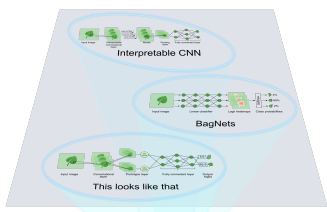
Model-agnostic



Model-specific



Interpretable by design models



Opening the black box



Designing a transparent glass box

Inputs



Data preprocessing

Planning



Identification and classification of cassava healthy and diseased leaves

1. Problem formulation

- What is the phenomics question to be answered?
- What are the target classes or variables?
- Is the target a mutually exclusive class?

2. Data preparation

- Are images taken in a controlled environment or in the field?
- Are there sufficient data?
- Are there enough features to distinguish between classes?
- Are the labels accurate?
- What is a good train-validate-test data split ratio?
- Is the data balanced? Is there a need for oversampling or undersampling?

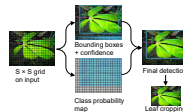
- #### 3. Choosing and training an AI algorithm
- Does the algorithm need to perform a classification or a regression task?
 - How much training data does the algorithm require?
 - Can it handle missing values?
 - What programming language and software framework to choose?
 - Are there sufficient compute resources to train the algorithm?

(i) Leaf dataset

Healthy: 2,077 images
Bacterial blight: 1,087 images
Brown streak: 2,180 images
Green mottle: 2,380 images
Mosaic: 13,100 images
Total: 21,307 images

Healthy: 1,360 images
Bacterial blight: 963 images
Brown streak: 1,820 images
Green mottle: 2,010 images
Mosaic: 11,091 images
Total: 17,190 images

(ii) Data cleaning



(iii) Data splitting

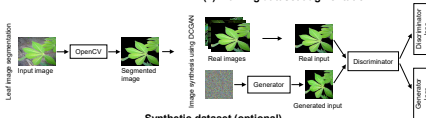


(iv) Image cropping using YOLO

Healthy: 20,110 images
Bacterial blight: 10,241 images
Brown streak: 10,674 images
Green mottle: 20,602 images
Mosaic: 19,962 images
Total: 99,619 images

Augmented training set

(v) Training dataset augmentation



4. Results and evaluation

- Are the results meaningful?
- Did the model overfit the training data? Did it underfit?
- Should the prediction be trusted? What were the most important features?
- Do the explanations make sense? Do they provide enough details?



Desktop server, training on two CPUs



Researcher's local computer, training on two GPUs



Reserved and on-demand cloud services, training on one or more GPUs

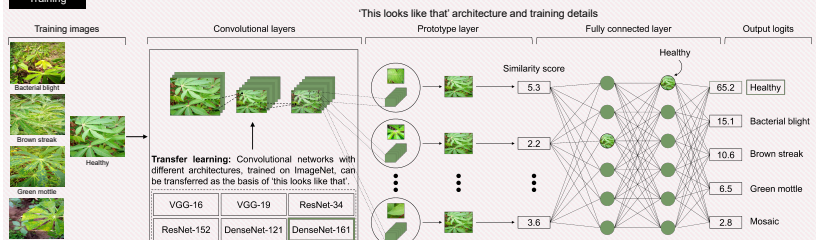


Lab cluster, training on eight GPUs



Exascale supercomputer, training on 26 GPUs or up to 27,648 GPUs

Training



Prototype pruning



Application

