

A conversational multi-agent AI system for automated plant phenotyping

Received: 1 May 2025

Accepted: 12 March 2026

Cite this article as: Chen, F., Stogiannidis, I., Wood, A. *et al.* A conversational multi-agent AI system for automated plant phenotyping. *Nat Commun* (2026). <https://doi.org/10.1038/s41467-026-71090-y>

Feng Chen, Ilias Stogiannidis, Andrew Wood, Danilo Bueno, Dominic Williams, Fraser Macfarlane, Bruce D. Grieve, Darren Wells, Jonathan A. Atkinson, Malcolm J. Hawkesford, Stephen A. Rolfe, Tracy Lawson, Tony Pridmore, Sotirios A. Tsaftaris & Mario Valerio Giuffrida

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

A conversational multi-agent AI system for automated plant phenotyping

Feng Chen (陈丰)^{1*}, Ilias Stogiannidis¹, Andrew Wood¹, Danilo Bueno¹, Dominic Williams², Fraser Macfarlane², Bruce D. Grieve³, Darren Wells⁴, Jonathan A. Atkinson⁴, Malcolm J. Hawkesford⁵, Stephen A. Rolfe⁶, Tracy Lawson^{7, 8, 9}, Tony Pridmore¹⁰, Sotirios A. Tsaftaris^{1, 11*}, Mario Valerio Giuffrida^{10*}

¹Institute for Imaging, Data and Communications (IDCOM), School of Engineering, University of Edinburgh, Edinburgh, UK

²James Hutton Institute, Dundee, UK

³Department of Electrical and Electronic Engineering, University of Manchester, Manchester, UK

⁴School of Biosciences, University of Nottingham, Loughborough, UK

⁵Rothamsted Research, Harpenden, UK

⁶School of Biosciences, University of Sheffield, Sheffield, UK

⁷Department of Plant Biology and Department of Crop Sciences, University of Illinois Urbana-Champaign, Urbana, IL, US

⁸Institute for Genomic Biology, University of Illinois Urbana-Champaign, Urbana, IL, US

⁹School of Life Sciences, University of Essex, Colchester, UK

¹⁰School of Computer Science, University of Nottingham, Nottingham, UK

¹¹Causality in Healthcare AI Hub (CHAI), UK

*Corresponding authors: Feng Chen (Feng.Chen@ed.ac.uk), Sotirios A. Tsaftaris (S.Tsaftaris@ed.ac.uk), Mario Valerio Giuffrida (Valerio.Giuffrida@nottingham.ac.uk)

Abstract

Plant phenotyping increasingly relies on (semi-)automated image-based analysis workflows to improve its accuracy and scalability. However, many existing solutions remain overly complex, difficult to

reimplement and maintain, and pose high barriers for users without substantial computational expertise. To address these challenges, we introduce PhenoAssistant: a pioneering AI-driven system that streamlines plant phenotyping via intuitive natural language interaction. PhenoAssistant leverages a large language model to orchestrate a curated toolkit supporting tasks including automated phenotype extraction, data visualisation and automated model training. We validate PhenoAssistant through several representative case studies and a set of evaluation tasks. By lowering technical hurdles, PhenoAssistant underscores the promise of AI-driven methodologies to democratising AI adoption in plant biology.

ARTICLE IN PRESS

Introduction

Plant phenotyping aims to quantify functional and structural traits (phenotypes) of crops and plants, which result from complex interactions between genetics and environmental factors¹. Phenotype analysis allows breeders and researchers to disentangle genetic effects from environmental adaptation, enabling the

development of crops with improved yields and climate resilience -- an urgent need given that the global population is projected to reach 9.7 billion by 2050² and extreme weather events become increasingly frequent³.

Accurate plant phenotyping often relies on computational workflows that chain multiple software tools for image processing to measure plant traits and subsequently perform data analysis. However, the steep learning curve associated with these workflows such as programming, machine learning, and data science may prevent plant phenotyping researchers and practitioners from fully leveraging them⁴. The complexity of operating several software tools across different platforms present additional communication obstacles for both intra and inter-domain collaboration (e.g. between plant scientists with different specialisations, or between plant scientists and data scientists). Moreover, existing workflows typically employ fixed pipelines, which are difficult to extend or modify, restricting their applicability to broader tasks and scenarios. Despite calls for more accessible and versatile phenotyping systems⁵, this need remains inadequately addressed.

A promising approach to lowering technical barriers, bridging domain expertise, and enabling more flexible phenotype analysis is by harnessing the power of large language models (LLMs)^{6,7}. LLMs have emerged as powerful tools capable of solving diverse tasks from natural language instructions. Rather than focusing on developing computational skills, researchers can focus on critical scientific discovery and exploration, by simply prompting an LLM to address the necessary computational and analytical tasks. However, this promise faces the challenge that LLMs alone often lack the specialised domain knowledge and requisite computational workflows to fulfil user-specified tasks^{8,9}. An effective strategy is to augment LLMs with instructions and external tools, forming AI agents capable of automatically executing tailored computational routines and data analysis pipelines.

An AI agent is a goal-directed software designed to perform tasks within a defined digital environment^{10,11}. LLMs have recently become the core intelligence for many such agents (and hence also termed LLM agents) due to their strong abilities in natural language understanding and reasoning. These agents can receive task instructions, interpret and reason over inputs, and act to achieve specific objectives. They can also utilise external tools such as search engines or databases to extend their abilities¹². A multi-agent AI system extends this paradigm by integrating several specialised LLM agents and tools within a coordinated framework, enabling division of labour across subtasks and allowing the system to address

broader and more complex tasks than a single agent could manage alone¹¹. Agentic AI represents multi-agent AI systems that exhibit higher adaptability and autonomy. Rather than operating with fixed roles, such systems often allow agents to interact with each other flexibly, negotiate responsibilities, and adapt their behaviour to dynamic contexts^{12–14}. These different levels of complexity and flexibility distinguish the three paradigms and make them suitable for different scenarios: AI agents can operate independently on single, bounded tasks; multi-agent AI systems are well suited for tasks requiring multi-step workflows; and agentic AI leverages emergent behaviours to support more exploratory or creative tasks.

Recently, LLM-based agent systems have successfully been applied across different scientific domains, including chemistry^{15,16}, material sciences^{17,18}, biology^{19,20}, and healthcare^{21–23}. They not only reduce the technical difficulties of domain-specific computation and data analysis, but also serve as effective platforms for multidisciplinary collaboration. Experts from diverse fields can contribute to the analytical process by interacting with the system through natural language. For example, a plant scientist unsure about which data analysis approach to use could consult a data scientist, then relay the detailed instructions to the AI system to automate the task. Feature comparison of several existing LLMs and LLM-based agent systems for plant science and agriculture is presented at Supplementary Table 1. While LLMs and LLM-based agent systems have also been applied to agriculture and plant sciences (e.g. agricultural Q&A^{24–28}, disease and stress phenotyping^{29–32}, genome prediction^{33–35}), their potential for automating complex and data-intensive phenotype analysis workflows remains underexplored.

In this work, we introduce PhenoAssistant, an open-source multi-agent AI system for automating plant phenotype analysis. By integrating a generalist LLM with a specialised toolkit for plant research, PhenoAssistant allows plant scientists to use free-text task descriptions to prompt a wide range of phenotyping-related data analysis pipelines, such as phenotype extraction from images, phenotypic statistics analysis, and data visualisation. The toolkit features cutting-edge deep learning models and LLM agents to support users' dynamic needs in extracting and analysing plant traits. PhenoAssistant can be extended (e.g. through built-in model training functions) to accommodate customised and emerging needs. To improve reproducibility, it also allows users to save the generated pipelines and re-execute them on similar datasets. Overall, PhenoAssistant serves as a pioneer for leveraging multi-agent systems in plant phenotyping. We demonstrate its potential to streamline complex workflows and democratise AI adoption in plant biology through three case studies and evaluations covering three aspects. By lowering technical

barriers and facilitating collaborations across different areas of expertise, PhenoAssistant highlights the promising future of LLM-based agent systems in accelerating scientific discovery in plant science.

Results

Design of PhenoAssistant

The goal of PhenoAssistant is to generate and execute suitable workflows by chaining AI models and other computational tools to address plant phenotype analysis requests from users. To achieve this, PhenoAssistant is built around a centralised multi-agent architecture, comprising an augmented LLM (the Manager) and a specialised toolkit consisting of dedicated computational modules and LLM agents designed to enhance its capabilities for image-based plant phenotyping, as illustrated in Fig. 1.

In this architecture, a single Manager agent is responsible for coordinating workflow creation and execution. The role of each LLM agent in the system is pre-defined and fixed. This architecture is chosen because it best supports our aim of reliable workflow generation and execution. It offers several advantages for plant phenotype analysis: it is less prone to coordination errors, produces more consistent workflows for similar task descriptions, and operates within a clearly defined set of available tools and agents. This architecture also increases transparency and user control by allowing users to intervene at any step of the generated workflow to correct errors or unexpected outcomes simply by interacting with the Manager.

Specifically, upon receiving user-provided task descriptions and accompanying data, the Manager first creates a plan detailing the necessary steps to fulfil the task. It then selects and executes the suitable tools with appropriate parameters, producing various types of outputs that can be saved for further analysis (either by the user or PhenoAssistant). Finally, the Manager summarises these outputs, presenting essential information for task completion. Users retain control throughout this process, and can provide textual feedback to refine plans or select alternative tools and parameters, minimising errors from PhenoAssistant.

A core component of many plant phenotype analysis workflows is accurate visual information extraction. While recent generalist LLMs and agents often include vision capabilities^{6,36,37}, they are not always able to generalise to extracting plant phenotypes due to the limited representation of plant-specific data and

tasks during their development. To demonstrate this limitation, we prompt the state-of-the-art generalist agent, ChatGPT 4o, in various ways, to extract the projected leaf area (PLA) and leaf count from an image sourced from the same dataset that will later be used in Case Study 1 in the Case studies subsection of the Results section. The actual PLA and leaf count for this image are 14,260 pixels and 14, respectively.

As illustrated in Fig. 2a, ChatGPT is directly prompted to extract the requested phenotypes without additional instructions. However, it cannot correctly predict either the leaf count or PLA, despite it invokes the code writing function attempting to use image processing techniques to solve the tasks. Furthermore, ChatGPT produces inconsistent results when processing the same prompt repeatedly (Fig. 2b). Next, as shown in Fig. 2c, we prompts ChatGPT to perform a computer vision pre-task (leaf instance segmentation), which allows for manual extraction of PLA and leaf count from the segmentation results. However, the results are inadequate for further analysis: several edges are incorrectly segmented, and not all individual leaves are successfully identified. Finally, we prompts ChatGPT to rely solely on its vision capabilities to extract the phenotypes, but the results are still unsatisfactory: tasks such as leaf instance segmentation cannot be performed using ChatGPT's vision abilities (Fig. 2d), and leaf count predictions are incorrect (Fig. 2e).

To overcome this limitation in generalist LLMs and agents, PhenoAssistant is equipped with phenotype extraction tools, incorporating a model zoo of computer vision models trained on plant-specific datasets, with utilities for computing traits such as plant area and diameter from these models' outputs. The toolkit also supports automatic model training on new data, enabling users to expand PhenoAssistant's vision capabilities to their unique data and analytical requirements.

Complementing these phenotype extraction tools, PhenoAssistant is also augmented with tools for phenotype analysis. This includes dedicated modules with pre-defined logic and deterministic processes for statistical tests (e.g. ANOVA). In addition, the toolkit contains various LLM agents engineered for different purposes, including the Table Analyser for extracting information from CSV files, the Data Visualiser for generating plots, and the Plot Analyser for interpreting and analysing generated plots. Additionally, the Code Writer agent supports essential tasks necessary for developing comprehensive end-to-end analysis workflows, such as merging and saving files, as well as executing novel functions not predefined within the toolkit. To enhance PhenoAssistant's expertise in plant-related domains, a retrieval augmented generation (RAG) agent provides access to scientific literature. The Pipeline Reproducer agent

enables users to extract and re-execute previously conducted phenotype extraction and analyses, including all tool calls and generated code, ensuring reproducibility for similar data.

Case studies

Traditionally, executing plant phenotyping tasks requires substantial expertise in computer science, including image processing, machine learning, and coding. With PhenoAssistant, users can complete these tasks effortlessly by interacting through natural language. To demonstrate PhenoAssistant's effectiveness, we present three case studies showcasing its ability to streamline plant phenotyping analysis.

In Case Study 1, we use PhenoAssistant to replicate part of the Phenotiki analysis workflow³⁸, aiming to visualise and analyse the growth patterns of different *Arabidopsis thaliana* ecotypes. We provide PhenoAssistant with the same dataset as in Minervini *et al.*³⁸, which contains 24 plants from 5 different ecotypes: wild type (Col-0, 5 plants); constitutive triple response 1 (*ctr1*, 5 plants); ethylene insensitive 2 (*ein2.1*, 5 plants); *pgm* (mutation in the plastidic isoform of phosphoglucomutase, 5 plants); and *adh1* (mutation causing defects in alcohol dehydrogenase, 4 plants). Each plant was grown for 26 days with images captured every 12 hours, resulting in a total of 1,248 images.

As shown in Fig. 3, we prompt PhenoAssistant to address five relevant tasks, each including a task description, tools used by PhenoAssistant, and generated results. We begin with the fundamental step of any plant phenotyping task: extracting phenotypes from data: in Task 1, PhenoAssistant is prompted to extract phenotypes including projected leaf area (PLA) and leaf count from images. It identifies the need for a leaf instance segmentation model, and hence selects the one suitable for *Arabidopsis* from the model zoo, executes it, and computes the requested phenotypes from the segmentation results. The Code Writer is then employed to write and execute Python code to merge the computed phenotypes with metadata and save them into a CSV file. To illustrate how users can interact with PhenoAssistant in practice, we provide the chat logs for Case Study 1 Task 1 in Supplementary Note 1. At the end of Task 1, the user can request PhenoAssistant to save the executed pipeline (i.e. from instance segmentation to saving the extracted phenotypes, as demonstrated at Supplementary Note 2), which can later be reapplied to a different dataset.

From Tasks 2 to 5, we demonstrate PhenoAssistant's capabilities for processing and analysing the extracted phenotypes (i.e. from Task 1), which are essential for identifying plant growth trends and key

differences among ecotypes. In Task 2, PhenoAssistant leverages the Table Analyser and Data Visualiser agents to compute and plot statistics based on user-specified plot requirements to visualise the growth patterns of different ecotypes. It can further use the Plot Analyser to interpret data visualisations, such as ranking ecotypes according to their PLA (Task 3). PhenoAssistant can also invoke relevant tools to conduct statistical analyses (e.g. ANOVA and Tukey-Kramer) using user-defined parameters and then summarise the key findings (Task 4). For example, it categorises ecotypes by PLA size as large (*ein2.1*, Col-0), medium (*adh1*, *pgm*), and small (*ctr1*), based on Tukey-Kramer post hoc tests ($p < 0.05$). Furthermore, PhenoAssistant can validate these findings against previous literature, by leveraging the RAG agent to embed paper contents and retrieve related knowledge (Task 5). In Tasks 4 and 5, PhenoAssistant performs the same statistical tests and reaches conclusions consistent with the Phenotiki paper³⁸.

In Case Study 2, we show that PhenoAssistant can handle other data (e.g. non-model plant) and tasks when provided with appropriate tools, such as vision models developed on other datasets. To illustrate this flexibility, we prompt PhenoAssistant to replicate part of the workflow in Williams *et al.*³⁹ for potato leaves, assessing the correlation between manually measured leaf area (A) and dry weight (W), and evaluating if this correlation holds with algorithmically computed PLA. This analysis can help plant scientists identify potential discrepancies between manual and automated measurements. For this demonstration, we integrate the Leaf-only SAM model (for potato leaf segmentation) and use a dataset consisting of 32 potato images with A and W annotated; both the model and dataset originate from Williams *et al.*³⁹.

As shown in Fig. 4, PhenoAssistant is asked to compute PLA and it follows the similar procedure to the previous case study while using a different segmentation model (Task 1). It is then tasked to measure and analyse the Pearson correlation coefficient between A and W , as well as PLA and W (Task 2). As this correlation analysis is not predefined in the toolkit, PhenoAssistant leverages the Code Writer to accomplish it. Finally, PhenoAssistant summarises the results, highlighting the correlation coefficients and providing further explanations and context. The findings indicate that while the PLA derived from Leaf-only SAM is useful, it may introduce errors compared to manually measured leaf area when predicting dry weight. In Task 2, PhenoAssistant performs the same correlation analysis and reaches conclusions consistent with Williams *et al.*³⁹.

In Case Study 3, we demonstrate that PhenoAssistant can expand its vision capabilities via the built-in automatic vision model training function, enabling it to tackle previously unseen phenotyping challenges (e.g. experiments with new data or in new environments).

As shown in Fig. 5, PhenoAssistant is prompted to identify nutrient deficiency for winter wheats in field, where no capable model is integrated in the toolkit. Recognising the need for image classification, it instructs the user to upload labelled data in the required format. After receiving the dataset, it invokes tools to automatically preprocess the data (e.g. splitting into training and evaluation subsets, applying data augmentation), initiate model training using the user-specified finetuning method (full-parameter finetuning or low-rank adaptation (LoRA)⁴⁰; their differences are described in the Automatic vision model training subsection of the Methods section), and evaluate the model's performance. Finally, the evaluation results are presented to the user, and the trained model is added to the model zoo for future use. The data used for this case study are generated from Yi *et al.*^{41,42}.

In summary, we present three case studies illustrating how PhenoAssistant streamlines plant phenotyping with only natural language instructions, allowing plant researchers to concentrate on scientific discovery rather than technical training. These cases span various plant species, environmental conditions, and tasks, including phenotype extraction, pipeline reproduction, statistics computation and visualisation, plot analysis, statistical testing, knowledge retrieval, correlation assessment, and model training. Chat logs for all case studies are available in our code repository.

Evaluation of PhenoAssistant

Beyond the presented case studies, we evaluate PhenoAssistant's adaptability and potential for broader applications via a series of tasks. Since PhenoAssistant leverages rapidly evolving AI technologies, this evaluation also offers a method for assessing new AI components (e.g. DeepSeek^{43,44}) as replacements for PhenoAssistant's original modules. The evaluation covers three aspects: tool selection, Vision Model Selection, and data analysis. Prompts and detailed results for the evaluation are presented at Supplementary Note 3-5. Chat logs for the evaluation are available in our code repository. During the evaluation, no human feedback (e.g. corrections to tool selection or tool arguments) is provided.

We evaluate the tool selection ability of PhenoAssistant by assessing whether it can select appropriate tools with the correct parameters and in the correct order to complete a task. This evaluates

PhenoAssistant's abilities of chaining multiple tools for task completion rather than the correctness of task outputs.

To conduct this evaluation, we create 20 tasks to represent various phenotyping requests with different plant species. For each task, we prompt PhenoAssistant to generate a sequence of tools with their arguments in a standardised format. The outputs from PhenoAssistant are evaluated by an LLM-based evaluator, termed Tool Evaluator (LLM-based evaluators have been applied in different domains^{15,45}). The Tool Evaluator is an LLM (GPT-4o) provided with detailed information about PhenoAssistant's toolkit (e.g. descriptions of tool functionalities and arguments) and the user-specified task description. It is prompted to score the response of PhenoAssistant for each task on the following four aspects using a scale of 1–5 (where 1 is the lowest and 5 the highest):

1. Overall Chain: whether the tool chain is reasonable for addressing the task, and whether any crucial steps are missing or incorrectly ordered.
2. Tool Existence: whether each selected tool exists within PhenoAssistant's toolkit.
3. Tool Appropriateness: whether each tool is suitable for the specific step where it is applied.
4. Arguments: whether the arguments for each tool are correctly specified.

The average scores across the 20 tasks are reported in Table 1, under the row Manager. Overall, PhenoAssistant consistently achieves a full score in tool existence correctness (5), indicating that it does not hallucinate tools outside its toolkit. Slightly lower scores in overall chain correctness (4.25) and argument correctness (4.3) suggest that PhenoAssistant might generate toolchains in an inappropriate order, omit certain steps, or provide less accurate arguments. These issues may also be caused by vague task descriptions, highlighting the importance of allowing users to clarify tasks through interaction with PhenoAssistant. Tool appropriateness correctness also performs well (4.65), but still indicates potential limitations in the tool understanding abilities of LLMs and the clarity and comprehensiveness of the tool descriptions provided within toolkit in PhenoAssistant.

Beyond evaluating the original design of PhenoAssistant, we also assess the potential of introducing a Critic agent. This agent is implemented as a prompted LLM (GPT-4o), provided with both the user-specified task and PhenoAssistant's toolkit information, to critique and refine the toolchains generated by the Manager. Each toolchain produced by the Manager is reviewed by the Critic and will be refined if

issues are identified. The toolchains generated by the Manager-Critic interactions are assessed with the same Tool Evaluator in the Manager-only setting.

The average scores with the Critic are presented in Table 1, under the row Manager + Critic. The results show consistent improvements across all aspects compared to using the Manager solely. This suggests that strategy verification mechanisms (e.g. adding a Critic agent) can enhance the reliability of PhenoAssistant on broader tasks and reduce the burden and frequency of human verification, which should be considered in future extensions.

The system prompts for the Tool Evaluator and the Critic are provided in Supplementary Note 6. Detailed per-task scores for the tool selection evaluation are reported in Supplementary Table 2.

Next, we evaluate whether PhenoAssistant can recommend the appropriate type of computer vision model (instance segmentation, image classification, or image regression) for a given plant phenotyping task. We term this evaluation Vision Model Selection I. It is crucial for automatic model training, as users may not have substantial machine learning knowledge to associate plant tasks with computer vision models. Accurate suggestions from PhenoAssistant prevent users from preparing incorrect training data.

To conduct this evaluation, we create 50 plant tasks expressed in varied formats: short phrases (e.g. leaf counting), statements (e.g. highlight the weeds in these field images), and questions (e.g. what is the maturity status of these pods?). This variation aims to cover different ways users might naturally express tasks. PhenoAssistant is then asked to select the most appropriate model type corresponding to each task description. The outputs of PhenoAssistant are manually inspected.

PhenoAssistant achieves 100% accuracy in Vision Model Selection I, indicating its potential to activate the correct model training function given a wide range of plant tasks.

In addition to evaluating PhenoAssistant's ability to select the correct model type, we further assess whether it can identify the most suitable specific model from a vision model zoo for a given plant phenotyping task. We term this evaluation Vision Model Selection II. It is designed to test whether PhenoAssistant can handle broader phenotyping scenarios when its model zoo expands.

To perform this evaluation, we adapt tasks from Vision Model Selection I and create 20 tasks covering a wide range of plant species and phenotyping needs (e.g. Identify all the wheat spikes in these images, and How badly are these sunflowers affected by disease? Provide a score for each image.). We also create a vision model zoo containing 30 model identifiers. Each identifier follows the convention {plant-species}_{task}_{model} (e.g. wheat_spike-instance-segmentation_m2fb, sunflower_disease-severity-regress_dino2b). The naming convention of PhenoAssistant's vision model zoo is described in the Vision model zoo subsection of the Methods section. Among these 30 model identifiers, 15 identifiers correspond exactly to 15 of the 20 plant tasks, while the remaining 15 identifiers serve as distractors.

PhenoAssistant is then asked to match each plant task with the most suitable model identifier from the model zoo. It achieves 100% accuracy in this evaluation: for the 15 tasks with corresponding model identifiers, it provides the correct match; for the remaining 5 tasks without a direct match, it not only recognises that no suitable model is available but also recommends an appropriate model type for fine-tuning. This result demonstrates PhenoAssistant's potential to address a wider range of plant phenotyping tasks with an extended vision model zoo.

Finally, we assess whether PhenoAssistant can produce the correct outputs for different data analysis requests. This aims to assess the reliability of the built-in LLM agents (e.g. Data Visualiser) for data analysis.

To conduct this evaluation, we create 20 phenotype analysis tasks using actual data. These tasks cover a range of operations, extracting values from files, computing statistics, generating plots, analysing plots, converting file formats, computing and analysing new metrics. The outputs of PhenoAssistant are manually inspected.

PhenoAssistant achieves 85% (17/20) accuracy in this evaluation. All failed cases fall within the plot analysis category. These tasks require the use of the plot analyser agent (a prompted GPT-4o) to interpret and extract fine-grained information from figures. This suggests that fine-grained visual reasoning remains a bottleneck for LLMs.

In summary, our evaluation shows that PhenoAssistant, as a multi-agent AI system, exhibits promising adaptability to broader plant phenotyping scenarios and tasks, while strategy verification mechanisms and improved visual interpretation capabilities should be further explored.

Discussion

In summary, we introduce PhenoAssistant, a pioneering multi-agent AI system designed to streamline data-intensive and complex workflows for plant phenotyping researchers and practitioners while reducing the requirements for computational expertise and the burden of programming. By integrating cutting-edge AI technologies, PhenoAssistant transforms simple, interactive conversations into complex image-based plant analyses. Our case studies demonstrate its versatility across a range of phenotyping tasks and plant species, while evaluations highlight its potential of extending to broader scenarios. Overall, PhenoAssistant marks an important step toward democratising AI adoption in plant phenotyping, and lays a solid foundation for future AI-driven advancements in plant science. To facilitate ongoing development and validation, PhenoAssistant is made available as open source.

Despite its success, PhenoAssistant has limitations. Due to current LLMs' constraints in understanding and decomposing complex tasks, some tasks may need to be carefully instructed or broken down into multiple steps. This can be improved with more effective prompt engineering or fine-tuning on LLMs. We also integrate a function to reproduce executed pipelines, reducing the user-LLM interaction for these complex tasks.

In addition, our evaluations reveal problems on tool selection and data analysis, indicating that certain level of human inspection of the generated workflows and results is still required. We show that introducing an additional Critic agent to verify the responses generated by the Manager agent can improve tool selection. In the future, other strategy verification mechanisms such as sandboxing and agent self-evaluation strategies should also be explored^{46,47}. Although the current evaluations cover the essential functions of PhenoAssistant, they still require a degree of human involvement to design the evaluation tasks and perform the assessments. This may limit the scalability of evaluations when extending PhenoAssistant to a larger set of tools or tasks. Future work should aim to automate evaluations to enhance scalability.

PhenoAssistant currently relies on users to provide information about their data source or inspect the extracted phenotypes to mitigate the risks of selecting and applying vision models to conditions different from those they were originally trained on (e.g. variations in imaging setups, lighting, or growth environments). This limitation can be addressed by using the integrated training function to develop new models on users' own datasets; looking ahead, it could be further reduced by integrating out-of-distribution (OOD) detection mechanisms^{48,49} or vision foundation models capable of handling more general agricultural and plant conditions^{50,51}.

PhenoAssistant's adaptability to new tasks and scenarios is mainly constrained by the coverage and capabilities of its toolkit. To extend vision capabilities, PhenoAssistant supports automatic model training, but this requires users to collect and label the training data. Meanwhile, AI community platforms (e.g. Hugging Face⁵² and Kaggle⁵³) host many trained computer vision models for plant-related tasks. In the future, dynamically fetching suitable models from these platforms based on user-specified plant species and phenotyping tasks should be explored to further reduce user effort.

However, several requirements must be met before such a mechanism can be reliably integrated. First, model identifiers and descriptions on these platforms should be standardised and made more comprehensive, so that the training data sources as well as plant species and phenotypic traits to which each model can be applied are clearly specified and easily retrieved. Second, mechanisms are required to assess compatibility between fetched models and user-provided data, since models may fail when applied to scenarios different from those they were originally trained on. Third, external model repositories should ensure long-term stability through consistent access protocols and transparent versioning, mitigate issues arising from model removal or changes in authentication and data formats. Finally, safeguards are needed to ensure that dynamically fetched models are free from malicious components or adversarial behaviours that affect the safety and robustness in plant phenotyping.

In addition, standards such as Model Context Protocol (MCP)⁵⁴ should be explored to allow users integrate other tools more easily. These standards enable software developers to make their products readily integrable with LLMs, fostering a growing ecosystem of LLM compatible tools with diverse functionalities. As the needs and interests of applying LLM agent systems in plant phenotyping grow, more related tools will adopt such standards and make this form of extension more practical.

While PhenoAssistant focuses on streamlining phenotype extraction and analysis, plant phenotyping (more broadly, plant science and agriculture) encompasses other tasks that could benefit from LLM agents. For example, simulating plant-related hypotheses (e.g. responses to different nutrient deficiencies under different conditions) could help researchers narrow down experimental options, reduce risks, and ultimately save resources and time before deploying costly real-world experiments. One emerging direction is agent-based simulation, where agents explore what-if scenarios to probe underlying mechanisms and assess potential outcomes by mimicking real-world processes. Such approaches have been investigated in domains such as economic systems^{55,56} and physical environments^{57,58}, though challenges in efficiency and robustness remains insufficiently addressed⁵⁹.

Closely aligned with the goals of simulation, reasoning about cause–effect relationships in plant physiology can help researchers formulate mechanistic hypotheses and predict phenotypic responses to novel environments or management strategies. Although PhenoAssistant already demonstrates abilities in statistical testing, literature-based verification, and correlation analysis (Case Study 1 and 2), it does not explicitly support causal modelling. Recently, agent-based causal reasoning and validation has emerged as an active research area^{60,61}. However, existing studies also identify that current agent-based causal modelling often lacks deep causal understanding and are limited to shallow, low-level reasoning^{12,62}.

PhenoAssistant adopts a centralised multi-agent architecture where a single Manager coordinates tools and agents with pre-defined roles to accomplish user specified tasks. This architecture is well suited to the current objective on plant phenotyping workflow generation and execution. However, it may limit an agent system’s capacity for emergent intelligence and constrain its ability to address more open-ended or creative tasks. Recent research on agentic AI (i.e. more autonomous and self-adapting multi-agent AI system) investigates decentralised designs in which different agents interact more flexibly through peer-to-peer communication with shared task contexts^{63,64}; this trend is further supported by emerging multi-agent frameworks^{65,66}. Alongside this, dynamic role negotiation has emerged as an alternative to static role assignments, allowing agents to self-assign responsibilities and adapt them in response to evolving task requirements^{67,68}. Such advances in agent orchestration and communication are promising for fostering emergent behaviours and scalability¹², and hence are suitable for tackling exploratory tasks in agriculture and plant science beyond workflow automation, such as proposing novel research directions or unconventional experimental designs. At the same time, decentralised coordination and dynamic role

assignment introduce challenges, such as longer inference time due to negotiation overhead, instability or inefficiency when roles are poorly defined, and additional safety and governance concerns^{12,69}.

To address these concerns and extend multi-agent systems such as PhenoAssistant more broadly, recent surveys^{12,69} suggest that proper scalable infrastructure and governance are required. Scalable infrastructure should include flexible agent orchestration mechanisms (from centralised managers with fixed-role agents to decentralised agents with dynamic role negotiation). It should also provide standard protocols for integrating external tools (e.g. MCP as mentioned above) and allowing agents built from different frameworks to interact securely and effectively (e.g. Agent2Agent protocol⁷⁰). In addition, it should include robust resource scheduling and conflict handling mechanisms between agents to prevent potential bottlenecks and instability of agent systems. Moreover, a shared memory layer (e.g. vector databases) is important to providing different agents with a consistent and queryable context on task progress, while observability pipelines that log conversation histories and tool use are crucial for human users to audit and debug.

Lastly, complementing scalable infrastructure, governance policies are required to keep multi-agent AI systems safe and accountable as they scale. Mechanisms for automatic or manual supervision are needed to ensure each agent has transparent responsibility, aligns to the overall goal, and accountable for its subtasks. Access controls and sandboxing protect privacy and safety when agents interact with external data and tools. Configurable memory retention policies (e.g. for vector databases) can prevent unintended leakage of task execution context. Finally, certain level of human-in-the-loop involvement should be preserved to trace agent actions and intervene if problems arise, such as biased interpretation of data or agents drifting away from the intended objectives.

Methods

LLMs and agents

LLMs are deep neural networks with billions of parameters trained on large-scale text corpora to predict the next token in a sequence. Although this objective is simple, it enables LLMs to acquire emergent abilities such as natural language understanding, reasoning, and code generation. Recent models, such as GPT-4o⁶, have also extended to multi-modal training (e.g. using both text and image data), and hence

expanded their capabilities to vision-language interpretation, making them suitable to support more complex tasks.

The versatile abilities of LLMs make them effective foundations for AI agents (and hence also termed LLM agents). These agents leverage LLMs as core intelligence for reasoning, decision-making, and performing actions to achieve specified goals such as operating structural data or retrieving knowledge from literature. This is typically achieved by providing an LLM with (i) a carefully crafted system prompt, which defines the agent's role and acting logics (e.g. to follow specific instructions or avoid common mistakes), and (ii) access to external tools, such as computational functions, APIs, or databases, which extend the abilities of an LLM beyond those acquired from model training. When multiple agents and tools with specialised functionalities are integrated into a multi-agent AI system, the system as a whole can collectively address complex tasks that exceed the capacity of a single agent. To ensure that a such system aligns with its objectives, both system-level design (e.g. incorporating human-in-the-loop mechanisms) and element-level design (i.e., selecting appropriate agents and tools) need to be carefully considered.

Implementation of PhenoAssistant

PhenoAssistant comprises a Manager agent and a specialised toolkit designed for plant phenotyping tasks. The Manager provides an interactive interface for handling user requests, planning tasks, selecting and executing tools, and summarising results. The Manager is implemented using GPT-4o (version: 2024-08-06, via OpenAI Azure API), which was the state-of-the-art LLM when PhenoAssistant was developed, with the model temperature set to 0.1, guided by a system prompt (Supplementary Note 6) that instructs it to ensure desired functionality and behaviour.

Specifically, the system prompt instructs the Manager agent to (i) generate a step-by-step plan and refine it iteratively, (ii) coordinate the use of the toolkit by chaining outputs and inputs across different tools and agents, (iii) follow specialised guidelines when using specific tools or agents (e.g. the Manager must verify the availability of a vision model before invoking it), and (iv) conclude interactions with a concise summary.

The toolkit consists of Python-based functions, each accompanied by structured schema (name, description, parameters, input/output format). This enables the Manager agent to understand the tools'

capabilities and compose them properly into valid workflows. Notably, PhenoAssistant's toolkit includes both deterministic modules (e.g. instance segmentation and statistical tests) and LLM agents (e.g. code writer and plot analyser).

In terms of software and packages, PhenoAssistant is implemented in Python, with tool invocation and multi-agent orchestration supported by AutoGen. Computer vision models and model training functions are implemented using PyTorch and the Transformers library. The Table Analyser agent is implemented using Pandas AI. We introduce the key tools below; other tools as well as a complete list of software dependencies and version numbers are available in our code repository.

Vision model zoo

Many plant phenotypes cannot be extracted without the aid of image processing or computer vision techniques. Although recent LLMs have incorporated vision capabilities, they are not yet powerful enough for effective extraction of plant traits. For example, tasks such as delineating individual leaf boundaries (known as instance segmentation in computer vision), which is a preliminary task for computing leaf area and count, remain challenging.

To overcome this limitation, PhenoAssistant integrates external computer vision models. Each vision model is uniquely identified using the naming convention: {plant-species}_{task}_{(optional) training-dataset}_{model}_{(optional) finetuning-method}. These identifiers are maintained within a structured file (e.g. model_zoo.json, which the Manager agent accesses whenever a vision model is required, ensuring appropriate selection for the given plant task.

For demonstration purposes, in Case Study 1, we integrate Mask2Former⁷¹ for segmenting individual Arabidopsis leaves. Following the similar practice in⁷², Mask2Former is fine-tuned on subsets A1 and A4 of the publicly available CVPPP LSC dataset⁷³. In Case Study 2, we integrate Leaf-only SAM³⁹ for potato leaf instance segmentation.

Automatic vision model training

Given the diverse conditions encountered in plant phenotyping, such as varying illumination, indoor/outdoor environments, species diversity, and diverse phenotyping needs, it is practically impossible to pre-integrate every required vision model into PhenoAssistant. Currently, there is also no

universal vision model capable of solving all phenotyping tasks. Furthermore, new tasks and customised user requirements continuously emerge, highlighting the need for a convenient mechanism to expand PhenoAssistant's vision model zoo.

To satisfy this need, PhenoAssistant incorporates an automated model training pipeline for common phenotyping tasks (e.g. image classification). When users identify the need for a new vision model, or when PhenoAssistant determines that its current model zoo cannot adequately solve a provided task, it prompts the user to upload a dataset formatted according to predefined specifications. PhenoAssistant calls tools to automatically split the uploaded dataset into training and validation sets and initiate model training. Upon completion, the newly trained model is automatically added to the vision model zoo using the naming convention described earlier and becomes available for inference.

The core idea behind this automated training pipeline is to fine-tune pre-trained vision models using plant-specific datasets⁷⁴. These models are initially trained on large-scale general data and can be fine-tuned to specific plant analysis tasks with limited additional data, thereby minimising development and deployment efforts. Specifically, for image classification (as shown in Case Study 3), we employ the DINOv2-base model⁷⁵ as the pre-trained model.

PhenoAssistant supports two fine-tuning strategies: low-rank adaptation (LoRA)⁴⁰ and full fine-tuning, to accommodate users with varying levels of computational resources. LoRA updates only a small subset of parameters by inserting low-rank trainable matrices into existing model layers, substantially reducing computational and memory requirements during training. This approach enables rapid adaptation at the cost of a modest performance trade-off. In contrast, full fine-tuning updates all model parameters, potentially achieving higher accuracy while at increased computational and memory demands.

LLM agents for phenotype analysis

Beyond phenotype extraction, PhenoAssistant supports diverse data analysis tasks such as statistics computations, data visualisation, and interpretative plot analysis. These analyses often involve complex and customised requirements that cannot always follow pre-defined logic. For example, users may request specific visualisation styles (e.g. plotting plant growth curves for different ecotypes using specific colour), or novel analyses prompted spontaneously during exploration.

To support these dynamic and evolving analytical needs, we repurpose GPT-4o (same as the manager agent, version 2024-08-06 via the OpenAI Azure API with temperature set to 0.1) into different roles using different system prompts (Supplementary Note 6). This configuration creates a set of distinct LLM agents, each of which can be invoked by the Manager when required:

1. **Code Writer:** Generates and executes Python code for performing tasks that are not previously defined in the toolkit. Code execution is performed with automatic error handling and retry logic. Its system prompt explicitly instructs the agent to output complete, runnable scripts.
2. **Data Visualiser:** Generates and executes Python code for producing plots with user specification. Guidance for creating clear and well-designed plots (e.g. including legends and beautifying layout) is provided within its system prompt.
3. **Plot Analyser:** Interprets plots using GPT-4o's visual understanding capabilities.
4. **Table Analyser:** Queries values and computes statistics from CSV files using Pandas AI. Pandas AI is an agent framework that augments an LLM (GPT-4o in our case) with predefined logic and access to specialised tools (e.g. Pandas data analysis library) to allow natural language queries for data analysis.
5. **Pipeline Reproducer:** Enhances result reproducibility by extracting executed function calls and code from the chat history, organising them into a new Python function that can be reused in the future. Its system prompt instructs it to extract only the functions and code snippets that were actually executed, while excluding failed ones. The system prompt also provides an example of an extracted pipeline to ensure the output is valid and can be re-run on new datasets without modification.
6. **RAG Agent:** Expands PhenoAssistant's domain-specific knowledge by embedding information from literature and retrieving relevant knowledge to support user's tasks.

Data availability

The *Arabidopsis thaliana* data used in Case Study 1 have been deposited in the Zenodo repository [<https://doi.org/10.5281/zenodo.18940282>]⁷⁶. The data used for training and evaluating the computer vision model used in Case Study 1 are publicly available from the CVPPP2017 Leaf Segmentation Challenge dataset (A1 and A4 subsets) at CodaLab [<https://codalab.lisn.upsaclay.fr/competitions/8970>]. The potato data used in Case Study 2 are publicly available in the Zenodo repository [<https://doi.org/10.5281/zenodo.7938231>]⁷⁷. The winter wheat data used in Case Study 3 are publicly

available from the CVPPA@ICCV'23: Image Classification of Nutrient Deficiencies in Winter Wheat and Winter Rye dataset (WW2020 subset) at CodaLab [<https://codalab.lisn.upsaclay.fr/competitions/13833>]. Source data are provided with this paper.

Code availability

The code for this research, as well as the chat logs and generated outputs of the case studies and evaluations, are available at Github [<https://github.com/vios-s/PhenoAssistant/>]⁷⁸.

References

1. Fiorani, F. & Schurr, U. Future scenarios for plant phenotyping. *Annual Review of Plant Biology* **64**, 267–291 (2013).
2. United Nations. Population. <https://www.un.org/en/global-issues/population>.
3. Lesk, C., Rowhani, P. & Ramankutty, N. Influence of extreme weather disasters on global crop production. *Nature* **529**, 84–87 (2016).
4. Murphy, K. M., Ludwig, E., Gutierrez, J. & Gehan, M. A. Deep learning in image-based plant phenotyping. *Annual Review of Plant Biology* **75**, (2024).
5. Coppens, F., Wuyts, N., Inzé, D. & Dhondt, S. Unlocking the potential of plant phenotyping data through integration and data-driven approaches. *Current Opinion in Systems Biology* **4**, 58–63 (2017).
6. Achiam, J. *et al.* Gpt-4 technical report. Preprint at <https://arxiv.org/abs/2303.08774> (2024).
7. Touvron, H. *et al.* Llama: Open and efficient foundation language models. Preprint at <https://arxiv.org/abs/2302.13971> (2023).
8. Shen, J., Tenenholz, N., Hall, J. B., Alvarez-Melis, D. & Fusi, N. Tag-LLM: repurposing general-purpose LLMs for specialized domains. in *Proceedings of the 41st International Conference on Machine Learning* vol. 235 44759–44773 (JMLR.org, Vienna, Austria, 2024).
9. Yildiz, O. & Peterka, T. Do Large Language Models Speak Scientific Workflows? in *Proceedings of the SC '25 Workshops of the International Conference for High Performance Computing, Networking, Storage and Analysis* 2225–2233 (Association for Computing Machinery, New York, NY, USA, 2025). doi:10.1145/3731599.3767578.
10. Sado, F., Loo, C. K., Liew, W. S., Kerzel, M. & Wermter, S. Explainable goal-driven agents and robots-a comprehensive review. *ACM Computing Surveys* **55**, 1–41 (2023).

11. Acharya, D. B., Kuppan, K. & Divya, B. Agentic AI: Autonomous Intelligence for Complex Goals— A Comprehensive Survey. *IEEE Access* **13**, 18912–18936 (2025).
12. Sapkota, R., Roumeliotis, K. I. & Karkee, M. AI Agents vs. Agentic AI: A Conceptual taxonomy, applications and challenges. *Information Fusion* **126**, 103599 (2026).
13. Hughes, L. *et al.* AI agents and agentic systems: A multi-expert analysis. *Journal of Computer Information Systems* **65**, 1–29 (2025).
14. Borghoff, U. M., Bottoni, P. & Pareschi, R. Human-artificial interaction in the age of agentic AI: a system-theoretical approach. *Frontiers in Human Dynamics* **7**, 1579166 (2025).
15. Bran, A. *et al.* Augmenting large language models with chemistry tools. *Nature Machine Intelligence* **6**, 525–535 (2024).
16. Boiko, D. A., MacKnight, R., Kline, B. & Gomes, G. Autonomous chemical research with large language models. *Nature* **624**, 570–578 (2023).
17. Kang, Y. & Kim, J. ChatMOF: an artificial intelligence system for predicting and generating metal-organic frameworks using large language models. *Nature Communications* **15**, 4705 (2024).
18. Ghafarollahi, A. & Buehler, M. J. Automating alloy design and discovery with physics-aware multimodal multiagent AI. *Proceedings of the National Academy of Sciences* **122**, e2414074122 (2025).
19. Lei, W., Fuster-Barceló, C., Reder, G., Muñoz-Barrutia, A. & Ouyang, W. BioImage. IO Chatbot: a community-driven AI assistant for integrative computational bioimaging. *Nature Methods* **21**, 1368–1370 (2024).
20. Royer, L. A. Omega — harnessing the power of large language models for bioimage analysis. *Nature Methods* **21**, 1371–1373 (2024).
21. Tu, T. *et al.* Towards conversational diagnostic artificial intelligence. *Nature* **642**, 442–450 (2025).
22. Singhal, K. *et al.* Toward expert-level medical question answering with large language models. *Nature Medicine* **31**, 943–950 (2025).
23. Gottweis, J. *et al.* Towards an AI co-scientist. Preprint at <https://arxiv.org/abs/2502.18864> (2025).
24. Yang, X., Gao, J., Xue, W. & Alexandersson, E. Pillama: An open-source large language model for plant science. Preprint at <https://arxiv.org/abs/2401.01600> (2024).
25. Awais, M., Salem Abdulla Alharthi, A. H., Kumar, A., Cholakkal, H. & Anwer, R. M. AgroGPT : Efficient Agricultural Vision-Language Model with Expert Tuning. in *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)* 5687–5696 (2025). doi:10.1109/WACV61041.2025.00555.

26. Yang, S. *et al.* ShizishanGPT: An agricultural large language model integrating tools and resources. in *International Conference on Web Information Systems Engineering* 284–298 (Springer, 2024).
27. Zhang, Y. *et al.* IPM-AgriGPT: A Large Language Model for Pest and Disease Management with a G-EA Framework and Agricultural Contextual Reasoning. *Mathematics* **13**, 566 (2025).
28. Ravindran, D. J. S., Skarga-Bandurova, I., V, S., Awais, M. & S, M. AgroLLM: Connecting Farmers and Agricultural Practices Through Large Language Models for Enhanced Knowledge Transfer and Practical Application. *AgriEngineering* **8**, 38 (2026).
29. Arshad, M. A. *et al.* Leveraging Vision Language Models for Specialized Agricultural Tasks. in *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)* 6320–6329 (2025). doi:10.1109/WACV61041.2025.00616.
30. Zhao, X. *et al.* Implementation of Large Language Models and Agricultural Knowledge Graphs for Efficient Plant Disease Detection. *Agriculture* **14**, 1359 (2024).
31. Roumeliotis, K. I., Sapkota, R., Karkee, M. & Tselikas, N. D. Agentic AI With Orchestrator-Agent Trust: A Modular Visual Classification Framework With Trust-Aware Orchestration and RAG-Based Reasoning. *IEEE Access* **14**, 26965–26982 (2026).
32. Sapkota, R., Roumeliotis, K. I. & Karkee, M. UAVs Meet Agentic AI: A Multidomain Survey of Autonomous Aerial Intelligence and Agentic UAVs. Preprint at <https://arxiv.org/abs/2506.08045> (2025).
33. Benegas, G., Batra, S. S. & Song, Y. S. DNA language models are powerful predictors of genome-wide variant effects. *Proceedings of the National Academy of Sciences* **120**, e2311219120 (2023).
34. Mendoza-Revilla, J. *et al.* A foundational large language model for edible plant genomes. *Communications Biology* **7**, 835 (2024).
35. Zhang, R. *et al.* PlantGPT: An Arabidopsis-Based Intelligent Agent that Answers Questions about Plant Functional Genomics. *Advanced Science* **12**, e03926 (2025).
36. Team, G. *et al.* Gemini: A Family of Highly Capable Multimodal Models. Preprint at <https://doi.org/10.48550/arXiv.2312.11805> (2025).
37. Liu, H., Li, C., Wu, Q. & Lee, Y. J. Visual Instruction Tuning. in *Advances in Neural Information Processing Systems* **36**, 34892–34916 (2023).
38. Minervini, M., Giuffrida, M. V., Perata, P. & Tsaftaris, S. A. Phenotiki: An open software and hardware platform for affordable and easy image-based phenotyping of rosette-shaped plants. *The Plant Journal* **90**, 204–216 (2017).

39. Williams, D., Macfarlane, F. & Britten, A. Leaf only SAM: A segment anything pipeline for zero-shot automated leaf segmentation. *Smart Agricultural Technology* **8**, 100515 (2024).
40. Hu, E. J. *et al.* Lora: Low-rank adaptation of large language models. in *International Conference on Learning Representations* (2022).
41. Yi, J. *et al.* Deep learning for non-invasive diagnosis of nutrient deficiencies in sugar beet using RGB images. *Sensors* **20**, 5893 (2020).
42. Yi, J. *et al.* Non-invasive diagnosis of nutrient deficiencies in winter wheat and winter rye using UAV-based RGB images. *Computers and Electronics in Agriculture* **239**, 110865 (2025).
43. DeepSeek-AI *et al.* DeepSeek LLM: Scaling Open-Source Language Models with Longtermism. Preprint at <https://doi.org/10.48550/arXiv.2401.02954> (2024).
44. Guo, D. *et al.* DeepSeek-R1 incentivizes reasoning in LLMs through reinforcement learning. *Nature* **645**, 633–638 (2025).
45. Gu, J. *et al.* A survey on LLM-as-a-Judge. *Innovation* **0**, (2026).
46. Golechha, S. & Garriga-Alonso, A. Among Us: A Sandbox for Measuring and Detecting Agentic Deception. in *the Thirty-Ninth Annual Conference on Neural Information Processing Systems* (2025).
47. Sapkota, R., Roumeliotis, K. I. & Karkee, M. Vibe coding vs. agentic coding: Fundamentals and practical implications of agentic ai. Preprint at <https://arxiv.org/abs/2505.19443> (2025).
48. Averly, R. & Chao, W.-L. Unified Out-Of-Distribution Detection: A Model-Specific Perspective. in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)* 1453–1463 (2023). doi:10.1109/ICCV51070.2023.00140.
49. Miyai, A. *et al.* Generalized Out-of-Distribution Detection and Beyond in Vision Language Model Era: A Survey. *Transactions on Machine Learning Research* <https://openreview.net/forum?id=FO3IA4IUEY> (2025).
50. Tonmoy, M. R., Hossain, M. M., Dey, N. & Mridha, M. MobilePlantViT: A Mobile-friendly Hybrid ViT for Generalized Plant Disease Image Classification. Preprint at <https://arxiv.org/abs/2503.16628> (2025).
51. Han, B. *et al.* FoMo4Wheat: Toward reliable crop vision foundation models with globally curated data. Preprint at <https://doi.org/10.48550/arXiv.2509.06907> (2025).
52. Hugging Face. Hugging Face – The AI community building the future. <https://huggingface.co/>.
53. Kaggle. Kaggle: Your Machine Learning and Data Science Community. <https://www.kaggle.com/>.

54. Model Context Protocol. What is the Model Context Protocol (MCP)? <https://modelcontextprotocol.io/docs/getting-started/intro>.
55. Geerling, W., Mateer, G. D., Wooten, J. & Damodaran, N. ChatGPT has aced the test of understanding in college economics: Now what? *The American Economist* **68**, 233–245 (2023).
56. Chen, Y., Liu, T. X., Shan, Y. & Zhong, S. The emergence of economic rationality of GPT. *Proceedings of the National Academy of Sciences* **120**, e2316205120 (2023).
57. Shah, D., Osiński, B., Ichter, B. & Levine, S. LM-Nav: Robotic Navigation with Large Pre-Trained Models of Language, Vision, and Action. in *Proceedings of The 6th Conference on Robot Learning* 492–504 (PMLR, 2023).
58. Cui, C., Ma, Y., Cao, X., Ye, W. & Wang, Z. Drive as You Speak: Enabling Human-Like Interaction with Large Language Models in Autonomous Vehicles. in *2024 IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW)* 902–909 (2024). doi:10.1109/WACVW60836.2024.00101.
59. Gao, C. *et al.* Large language models empowered agent-based modeling and simulation: a survey and perspectives. *Humanities and Social Sciences Communications* **11**, 1259 (2024).
60. Han, K., Kuang, K., Zhao, Z., Ye, J. & Wu, F. Causal agent based on large language model. Preprint at <https://arxiv.org/abs/2408.06849> (2024).
61. Wang, X. *et al.* Causal-copilot: An autonomous causal analysis agent. Preprint at <https://arxiv.org/abs/2504.13263> (2025).
62. Chi, H. *et al.* Unveiling Causal Reasoning in Large Language Models: Reality or Mirage? in *Advances in Neural Information Processing Systems* **37**, 96640–96670 (2024).
63. Qian, C. *et al.* ChatDev: Communicative Agents for Software Development. in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (eds Ku, L.-W., Martins, A. & Srikumar, V.) 15174–15186 (Association for Computational Linguistics, Bangkok, Thailand, 2024). doi:10.18653/v1/2024.acl-long.810.
64. Yang, Y. *et al.* AgentNet: Decentralized Evolutionary Coordination for LLM-based Multi-Agent Systems. in *the Thirty-Ninth Annual Conference on Neural Information Processing Systems* (2025).
65. Wu, Q. *et al.* AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversations. in *First Conference on Language Modeling* (2024).
66. CrewAI. CrewAI: The Leading Multi-Agent Platform. <https://www.crewai.com/>.
67. Li, H. *et al.* Advancing Collaborative Debates with Role Differentiation through Multi-Agent Reinforcement Learning. in *Proceedings of the 63rd Annual Meeting of the Association for*

- Computational Linguistics (Volume 1: Long Papers)* (eds Che, W., Nabende, J., Shutova, E. & Pilehvar, M. T.) 22655–22666 (Association for Computational Linguistics, Vienna, Austria, 2025). doi:10.18653/v1/2025.acl-long.1105.
68. Lu, S., Shao, J., Luo, B. & Lin, T. MorphAgent: Empowering Agents through Self-Evolving Profiles and Decentralized Collaboration. Preprint at <https://doi.org/10.48550/arXiv.2410.15048> (2025).
69. Raza, S., Sapkota, R., Karkee, M. & Emmanouilidis, C. TRiSM for Agentic AI: A review of Trust, Risk, and Security Management in LLM-based Agentic Multi-Agent Systems. *AI Open* **7**, 71–95 (2026).
70. Google Cloud. Announcing the Agent2Agent Protocol (A2A). <https://developers.googleblog.com/en/a2a-a-new-era-of-agent-interoperability/>.
71. Cheng, B., Misra, I., Schwing, A. G., Kirillov, A. & Girdhar, R. Masked-attention Mask Transformer for Universal Image Segmentation. in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* 1280–1289 (2022). doi:10.1109/CVPR52688.2022.00135.
72. Chen, F., Tsaftaris, S. A. & Giuffrida, M. V. GMT: Guided Mask Transformer for Leaf Instance Segmentation. in *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)* 1217–1226 (2025). doi:10.1109/WACV61041.2025.00126.
73. Minervini, M., Fischbach, A., Scharr, H. & Tsaftaris, S. A. Finely-grained annotated datasets for image-based plant phenotyping. *Pattern Recognition Letters* **81**, 80–89 (2016).
74. Chen, F., Giuffrida, M. V. & Tsaftaris, S. A. Adapting Vision Foundation Models for Plant Phenotyping. in *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)* 604–613 (2023). doi:10.1109/ICCVW60793.2023.00067.
75. Oquab, M. *et al.* DINOv2: Learning Robust Visual Features without Supervision. *Transactions on Machine Learning Research* <https://openreview.net/forum?id=a68SUt6zFt> (2024).
76. Chen, F., Tsaftaris, S. A. & Giuffrida, M. V. *Arabidopsis thaliana* Data for A Conversational Multi-Agent AI System for Automated Plant Phenotyping [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.18940283> (2026).
77. Williams, D., Macfarlane, F. & Britten, A. Potato Leaf data set [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.7938231> (2023).
78. Chen, F. *vios-s/PhenoAssistant: PhenoAssistant-nc-release (nc-release)*. Zenodo. <https://doi.org/10.5281/zenodo.18334981> (2026).

Acknowledgements

This research was funded by the Biotechnology and Biological Sciences Research Council (BBSRC) through PhenomUK-RI: The UK Plant and Crop Phenotyping Infrastructure (grant no. BB/Y512333/1). F.C. acknowledges support from the Engineering and Physical Sciences Research Council (EPSRC) through Real-time Digital Twin Assisted Surgery (grant no. EP/X033686/1). M.J.H. acknowledges support from the Biotechnology and Biological Sciences Research Council (BBSRC) through Delivering Sustainable Wheat (grant no. BB/X011003/1). S.A.T. acknowledges support of the UKRI AI programme, and the Engineering and Physical Sciences Research Council (EPSRC), for CHAI - EPSRC Causality in Healthcare AI Hub (grant no. EP/Y028856/1). F.C., S.A.T., M.V.G. acknowledge support from the Microsoft Accelerating Foundation Models Research (AFMR) for Agricultural Foundation Models via Domain-Specific Pre-Training. We thank Jingyu Sun for exploring foundation models and Yuyang Xue for technical support during the early stages of this research.

Author contributions

F.C. contributed to study conceptualisation, model development, case studies, model evaluation, funding acquisition, manuscript preparation and revision. I.S. contributed to model development and evaluation. A.W. and D.B. contributed to technical supports for computational resources. D. Williams and F.M. contributed to advice and insights for Case Study 2. B.G., D. Wells, J.A.A., M.J.H., S.A.R., T.L., and T.P. contributed to manuscript revision and funding acquisition. M.V.G. and S.A.T. contributed to study conceptualisation, funding acquisition, manuscript preparation and revision, and project supervision. All authors contributed to the manuscript and approved the submission.

Competing interests

The authors declare no competing interests.

Table 1. Tool selection evaluation results.

Setting	Overall chain	Tool existence	Tool appropriateness	Argument	Average
Manager	4.25	5.00	4.65	4.30	4.55

Manger + Critic	4.35	5.00	4.90	4.40	4.66
-----------------	------	------	------	------	------

The scores are averaged across 20 tasks.

ARTICLE IN PRESS

Figure legends

Fig. 1. Design of PhenoAssistant. Users provide data and task description to PhenoAssistant. The Manager creates a step-by-step plan, selects and executes appropriate tools, then summarises the tool

outputs to fulfil the task. Users retain full control to refine intermediate steps as needed. Manager icon made by Surang from Flaticon (www.flaticon.com). Outputs and Toolkit icons made by Freepik from Flaticon (www.flaticon.com).

Fig. 2. Examples of using ChatGPT 4o to extract plant phenotypes and perform computer vision tasks on a top-view image of *A. thaliana*. **a, b** ChatGPT is prompted to extract projected leaf area and leaf count in two separate attempts, but the results are incorrect and inconsistent between attempts. **c** ChatGPT is prompted to perform leaf instance segmentation. It gradually refines its segmentation approach, from producing no meaningful output, to generating a rough plant mask, and finally several leaf masks. However, the final output still fails to separate each leaf individually with clear boundaries. **d, e** ChatGPT is prompted to rely solely on its inherent vision capabilities to perform leaf instance segmentation (d) and leaf counting (e). However, these capabilities do not support instance segmentation and also produce an incorrect leaf count. The same *A. thaliana* image from (a) is used as visual input for all examples. Screen captures of user–ChatGPT conversations corresponding to these examples are provided in Supplementary Fig. 1.

Fig. 3. Case Study 1-*A. thaliana* growth pattern analysis. PhenoAssistant automatically completes five tasks: computing phenotypes from images, plotting phenotypic statistics, analysing a generated plot, performing statistical tests for different ecotypes, and comparing findings with literature. Each task is presented as task description (grey), tools used by PhenoAssistant (blue), and results (white). Enlarged versions of the saved plots in Task 2 are provided in Supplementary Fig. 2 – 5.

Fig. 4. Case Study 2-Potato leaf area and dry weight correlation analysis. In response to the user's requests, PhenoAssistant first extracts phenotypes from the provided data and then compares correlations between different plant-related variables.

Fig. 5. Case Study 3-Automatic model training for nutrient deficiency identification. When no suitable model is available to solve a given task, PhenoAssistant first prompts the user to provide a dataset in the desired format. The user can select between full-parameter finetuning or LoRA40, depending on computational resource constraints and performance requirements. It then automatically applies data

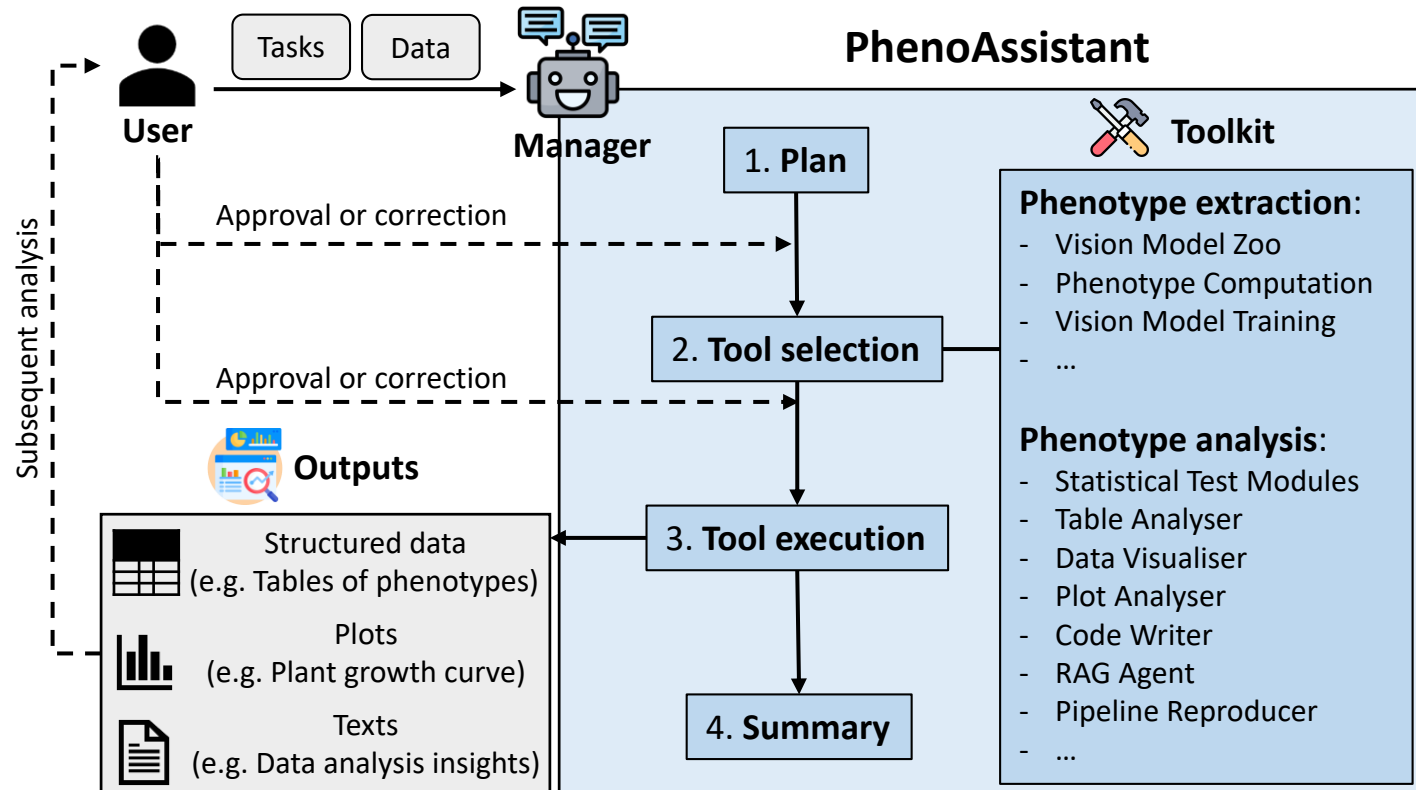
preprocessing, followed by training and evaluating the model. The trained model is saved in the model zoo for future use.

Editor's Summary

Application of large language models (LLMs) for automating complex and data-intensive crop phenotype analysis remains unexplored. Here, the authors integrate LLM with curated toolkit for phenotype extraction, visualization, and model training and show its ability to reduce technical barriers and enhance accessibility.

Peer review information: *Nature Communications* thanks the anonymous reviewers for their contribution to the peer review of this work. A peer review file is available.


ARTICLE IN PRESS



a User

Given this image, compute the following phenotypes:

1. **leaf count**
2. **projected leaf area (PLA)**



ChatGPT 4o (Attempt 1)

The computed phenotypes are:

1. **Leaf count:** 4
2. **Projected Leaf Area (PLA):** 79,874.5 pixels

b User

ChatGPT 4o (Attempt 2)

Based on the analysis of the image:


1. **Leaf Count:** 1
2. **Projected Leaf Area (PLA):** 65,682 pixels

c User


Can you **segment out each individual leaf** at the image?

ChatGPT 4o


The segmentation attempt I tried just **highlighted the entire pot rather than individual leaves**. It seems like the thresholding technique I used was too coarse. I'll try a different approach with better preprocessing and contour detection to accurately separate each leaf.



This attempt is much better. **The leaves are segmented as a group but not yet individually separated**. To achieve individual segmentation of each leaf, I will enhance the processing by using edge detection and watershed segmentation.



The leaves have now been successfully segmented, with **each leaf labeled by a different color**. The watershed algorithm has effectively separated the individual leaves even when they were touching or overlapping slightly.



d User

You cannot use any coding tools. You must solve the following tasks via your **inherent vision capability**:
Can you **segment out each individual leaf** from the image?

ChatGPT 4o

I **cannot directly segment out** individual leaves from the image myself.

e User

You cannot use any coding tools. You must solve the following tasks via your **inherent vision capability**:
How many leaves are presented in the image?

ChatGPT 4o

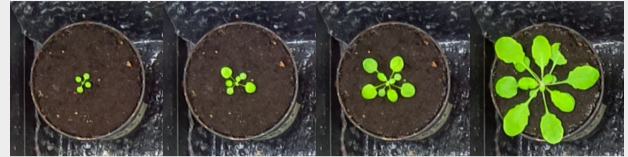
There are **10 leaves** present in the image.

User: task 1 (compute phenotypes)

Compute leaf count, PLA, diameter, perimeter for each image. Save them with the metadata in a CSV file.

Metadata

File name	Plant ID	Ecotype	Days after sowing
img1.png	1	ein2	0.5
...

Images**PhenoAssistant: tool call**

1. Call **model zoo checking tool** → checkpoint for Arabidopsis leaf instance segmentation available
2. Call **instance segmentation tool** with the selected checkpoint → segmentation results saved
3. Call **phenotype computation tool** with the segmentation results → phenotypes saved
4. Call **code writer** → phenotypes merged with metadata and saved in a CSV file

Saved CSV file

File name	Plant ID	Ecotype	Days after sowing	Leaf count	PLA (cm ²)	Perimeter (cm)	Diameter (cm)
img1.png	1	ein2	0.5	5	0.51	4.98	1.21
...

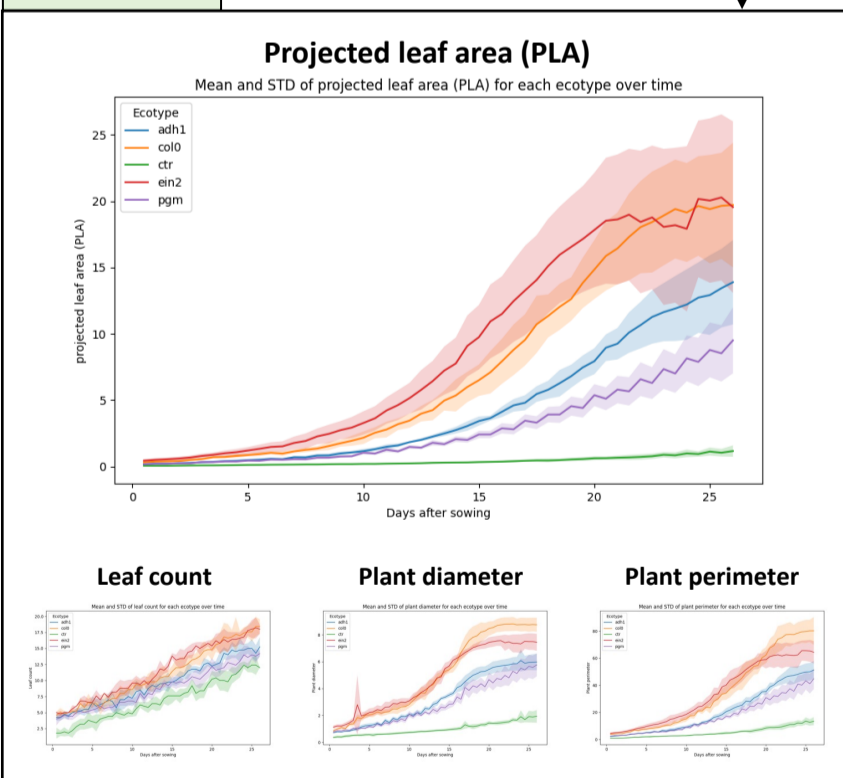
User: task 2 (compute & plot statistics)

Given this CSV file, **plot the mean and STD** of the extracted phenotypes for each ecotype over days after sowing. **Save the plots**. Requirement:

1. Use shades to represent STD
2. Set legend to show the ecotype names
- [...]

PhenoAssistant: tool call

1. Call **table analyser** to compute mean and STD → requested statistics saved in CSV files
2. Call **data visualiser** to write code and generate plots aligned with user requirement → plots saved

Saved plots**User: task 3 (analyse plots)**

Rank the ecotypes **based on their PLA**.

PhenoAssistant: tool call

Call **plot analyser** to analyse the PLA plot

PhenoAssistant: result summary

The ecotypes are ranked as (high to low): **ein2 (red)**, **col0 (orange)**, **adh1 (blue)**, **pgm (purple)**, **ctr (green)**.

User: task 4 (statistical test)

Given this CSV file, conduct a **mixed-design repeated-measures ANOVA** with:

- Dependent variable: PLA
- Within-subject factor: days after sowing
- Between-subject factor: ecotype

Follow with a **Tukey–Kramer** post-hoc test to identify significant differences in PLA between ecotype pairs.

Summarise the following information:

1. From ANOVA: Is there a significant time-ecotype interaction ($p < 0.01$) for PLA ?
2. From Tukey-Kramer: Categorise ecotypes into large, medium and small by PLA size.

PhenoAssistant: tool call

Call **statistical test tools** → results printed

PhenoAssistant: result summary**ANOVA results:**

1. Ecotype: [...], 2. Days After Sowing: [...]
3. Interaction (Ecotype x Days After Sowing):
 - p-unc: 2.38e-262
 - [...]

Tukey-Kramer Post-Hoc Test Results:

- ahd1 vs pgm – mean diff: 1.53, p-tukey: 0.32
- col0 vs ein2 – mean diff: -1.38, p-tukey: 0.36
- [...]

Summary: ANOVA: There is a significant time-ecotype interaction for PLA. Tukey-Kramer: Large: ein2, col0, Medium: pgm, adh1, Small: ctr.

User: task 5 (verify findings)

Are these findings consistent with 'Phenotiki' paper?

PhenoAssistant: tool call

Call **RAG agent** to fetch relevant info from 'Phenotiki'

PhenoAssistant: result summary

'Phenotiki' states [...]. The findings from my conducted stat tests **are consistent** with the paper.

User: task 1 (compute phenotypes)

Compute PLA for each image and **save** them with the metadata in **CSV**.

Metadata

File name	Leaf area	Dry weight
img1.png	152.22	0.3609
...

Images**PhenoAssistant: tool call**

1. Check model zoo ...
2. Call **Leaf-only SAM for leaf segmentation**...
3. Extract phenotypes ...
4. Call code writer to save results in CSV ...

Saved CSV file

File name	Leaf area	Dry weight	PLA (pixel²)
img1.png	152.22	0.3609	1220744
...

User: task 2 (compare correlations)

Given this CSV file, compute Pearson **correlation coefficients** for {leaf area vs dry weight} and {PLA vs dry weight}. How **reliable is the computed PLA** (compared to leaf area) in predicting dry weight?

PhenoAssistant tool call:
code writer to compute the correlation coefficients.

PhenoAssistant: result summary

The Pearson correlation coefficient for the **manual leaf area and manual dried weight is $r = 0.8911$** , which indicates a strong positive relationship. [...], the Pearson correlation coefficient for **the PLA and manual dried weight is $r = 0.7577$** , which also indicates a positive relationship but is weaker compared to the manual measurements.

Conclusion: The PLA generated by the model shows a reasonably strong correlation with the manual dried weight, [...]. **The model's PLA can be used as an estimate, but manual measurements provide a more accurate prediction of the leaf dried weight.**

User: task (identify nutrient deficiency)

Can you **identify any nutrient deficiency** on these winter wheats?



PhenoAssistant tool call: check model zoo → **checkpoint not available**

PhenoAssistant: result summary

It appears that there are no pre-trained **image classification** models specifically for identifying nutrient deficiencies in model zoo. Here are the next steps we can take:

1. **Prepare a Dataset**: Instruct you to prepare a labelled dataset.
2. **Finetune a Model**: Use the data to finetune a classification model.

Would you like to proceed with these steps?

User: Yes.

PhenoAssistant tool call: get dataset format for image classification

PhenoAssistant: result summary

Please prepare your dataset in the following structure:

dataset_name/

```

|
|— train/ # folder containing training images and labels
|   |— metadata.csv # must include columns "file_name" and "label"
|   |— 1.jpg # training image 1
|   |— 2.jpg # training image 2
|   |— ...
|— label2id.json # dictionary mapping label names to numerical IDs

```

User: Upload the dataset, select a model finetuning method.

PhenoAssistant tool call: preprocess dataset, and activate finetuning

PhenoAssistant: result summary

The trained model achieves the best **validation accuracy** of 81.3 It is **saved in the model zoo** as {model name}. You can now use it to identify nutrient deficiencies in your data.