

# InfraPhenoGrid: A scientific workflow infrastructure for Plant Phenomics on the Grid

Christophe Pradal<sup>a,b</sup>, Simon Artzet<sup>c</sup>, Jerome Chopard<sup>d</sup>, Dimitri Dupuis<sup>e</sup>,  
Christian Fournier<sup>c,b</sup>, Michael Mielewczik<sup>c,f</sup>, Vincent Negre<sup>c</sup>, Pascal Neveu<sup>d</sup>,  
Didier Parigot<sup>e</sup>, Patrick Valduriez<sup>e</sup>, Sarah Cohen-Boulakia<sup>b,e,g</sup>

<sup>a</sup>*CIRAD, UMR AGAP, Montpellier, France*

<sup>b</sup>*Inria, VirtualPlants, Montpellier, France*

<sup>c</sup>*INRA, UMR459, LEPSE, F-34060 Montpellier, France*

<sup>d</sup>*INRA, UMR729, MISTEA, F-34060 Montpellier, France*

<sup>e</sup>*Inria, Zenith, Montpellier, France*

<sup>f</sup>*ICCH, NHLI, Imperial College London, UK*

<sup>g</sup>*Laboratoire de Recherche en Informatique, Université Paris-Sud, CNRS UMR 8623,  
Université Paris-Saclay, Orsay, France*

---

## Abstract

Plant phenotyping consists in the observation of physical and biochemical traits of plant genotypes in response to environmental conditions. Challenges, in particular in context of climate change and food security, are numerous. High-throughput platforms have been introduced to observe the dynamic growth of a large number of plants in different environmental conditions. Instead of considering a few genotypes at a time (as it is the case when phenomic traits are measured manually), such platforms make it possible to use completely new kinds of approaches. However, the data sets produced by such widely instrumented platforms are huge, constantly augmenting and produced by increasingly complex experiments, reaching a point where distributed computation is mandatory to extract knowledge from data.

In this paper, we introduce InfraPhenoGrid, the infrastructure we designed and deploy to efficiently manage data sets produced by the PhenoArch plant phenomics platform in the context of the French Phenome Project. Our solution consists in deploying scientific workflows on a Grid using a middleware to pilot workflow executions. Our approach is user-friendly in the sense that despite the intrinsic complexity of the infrastructure, running scientific workflows and understanding results obtained (using provenance information) is kept as simple as possible for end-users.

*Keywords:* Phenomics, Scientific Workflows, Provenance, Grid computing

---

## 1. Introduction

Biological research derives its findings from the proper analysis of experiments. However, over the last three decades, both throughput of experiments (from single observations to terabytes of sequences of images produced during a single day) and the breadth of questions studied (from single molecules to entire genomes) have increased tremendously. One of the main challenges remains to efficiently analyze, simulate and model such big data sets while keeping scientist users in the loop.

In this paper we introduce InfraPhenoGrid, the infrastructure we designed and deployed to efficiently manage and analyze data sets produced by the PhenoArch plant phenomics platform. In this context, one difficulty remains to enable users to analyze, simulate and model increasingly huge data sets on a more frequent base. More precisely, the design of InfraPhenoGrid is driven by three needs, described here-after.

First, management of large-scale experiments involving possibly large numbers of interlinked tools has to be supported. Users should be able to analyze and simulate complex structural-functional relationships of plant architectures, integrating multi-disciplinary models developed by different teams. Experiments should be easy to design by users and it is important that over time they can be changed, adapted to new needs (new analysis algorithms are constantly available), and then shared. As a result, the first brick of our infrastructure is a *Scientific Workflow System*.

Second, each experiment can be replayed several times, varying data sets and/or parameter settings. Keeping track of the exact data sets and parameter settings used to produce a given result (provenance) is of paramount importance for scientists to ensure the results reproducibility and allow to properly interpret and understand them. The possibility of comparing results, obtained on several experiments when varying data sets and/or parameter settings are used, is another need directly associated with provenance. Consequently, the second brick of our infrastructure is a *Provenance Layer*.

Last but not least, our infrastructure has to efficiently deal with the analysis of huge data sets, possibly acquired on multiple sites. Analysis may involve combining data produced by platforms with completely different kinds of data, including data obtained from public data sources. Data acquisition

35 is fast compared to the time needed to analyze them. The size of data sets  
36 has reached a turning point at which local infrastructures are no longer suffi-  
37 cient to provide adequate computational power and storage facilities. Hence,  
38 distributed computation has become a major requirement. However, deploy-  
39 ing jobs on a parallel environment might be complex for end users. Therefore  
40 the third brick of our infrastructure introduces a *Middleware* able to pilot  
41 the execution of jobs on parallel (Grid) environments.

42 This paper is organized as follows: Section 2 introduces the precise con-  
43 text of this work, that is, the Phenome Project, PhenoArch platform and  
44 one use case of interest. Section 3 describes in detail the architecture of In-  
45 fraPhenoGrid. Section 4 demonstrates the benefit of using our solution for  
46 managing plant phenomic data sets. Section 5 provides related work while  
47 Section 6 concludes the paper and draws perspectives.

## 48 2. Use Case

### 49 2.1. The Phenome project and the PhenoArch platform

50 Selecting genotypes that maintain and increase crop performance is a par-  
51 ticularly challenging and important topic in the context of societal challenges  
52 such as climate change adaptation, food security and preserving natural re-  
53 sources.

54 A large variety of tasks have to be performed to collect information on  
55 plant traits (called phenotyping), including measuring the size of the leaves,  
56 counting the number of tails... Performing such tasks manually makes it  
57 impossible to consider more than a few plants at a time and it thus cruelly  
58 confines the kind of analyses that can be conducted.

59 In the meantime, massive plant phenotyping in the field, that is, the  
60 evaluation of crop performance (yield) of millions of plants in a large range  
61 of environmental and climatic scenarios, has been very efficient for driving  
62 plant breeding. However plant breeding is now facing a stagnation of genetic  
63 progress in several species. New strategies, such as genomic selection, are  
64 now evolving to directly link the allelic composition of a genome, available  
65 at much higher throughput and lower cost than field phenotyping, to crop  
66 performance. The existence of large marker-environment interactions, *i.e.*  
67 the fact that a given combination of markers has very different genetic val-  
68 ues depending on the climatic scenario, lead concomitantly to a revolution  
69 in phenotyping strategies. Such strategies aim to capture under controlled  
70 conditions, the genetic variability of plant responses to environmental factors

71 for thousands of plants (reference panels), hence identifying more heritable  
72 traits for genomic selection.

73 This first implies the necessity to automate quantification of a large num-  
74 ber of traits, to characterize plant growth, plant development and plant func-  
75 tioning. Second, it requires a tight control or at least accurate measurement  
76 of environmental conditions as sensed by plants. It finally requires fluent and  
77 versatile interactions between data and continuously evolving plant response  
78 models. Such interactions are essential to be considered in the analysis of  
79 a given marker environment interaction and in the integration of processes  
80 to predict genetic values of allelic combinations in different environment sce-  
81 narios.

82 High-throughput phenotyping platforms have thus been designed to allow  
83 growing and observing traits of a large number of plants. These platforms  
84 provide many measurements and imaging functionalities for different plant  
85 species grown in various environmental conditions. They potentially allow to  
86 assess the genetic variability of plant responses to environmental conditions  
87 using novel genetic approaches requiring a large number of genotypes.

88 Nine of such platforms, distributed over various regions of France, are  
89 gathered in the Phenome project (Figure 1). More precisely, Phenome con-  
90 sists of two controlled condition platforms (greenhouses with automated ir-  
91 rigation,  $CO_2$  control and temperature control) for 1900 plants, two field  
92 platforms (800 plots) equipped with environment control ( $CO_2$  enrichment,  
93 automated rain shelters) and three larger field platforms (2000 plots) that use  
94 natural gradients of water availability or soil contents. All these platforms  
95 are equipped with environmental sensors and permit automated imaging of  
96 plants in one or multiple wavelengths (thus allowing functional analysis)  
97 using robots to convey plants (for green houses) or to carry instruments  
98 to automatically acquire data in the field (Phenomobile, drones). Finally  
99 two supporting omic platforms enable us to centralize and optimize high  
100 throughput metabolomic and structural measurements associated with the  
101 experiments.

102 The work depicted in this paper is related to the PhenoArch platform<sup>1</sup>,  
103 in the south of France (Montpellier). As depicted in Figure 2, PhenoArch  
104 is composed of a conveyor belt storage structure of 28 lanes carrying 60  
105 carts each (*i.e.* total of 1680 pots), and a conveyor belt system that feeds

---

<sup>1</sup><https://www6.montpellier.inra.fr/lepse/M3P/plateforme-PHENOARCH>

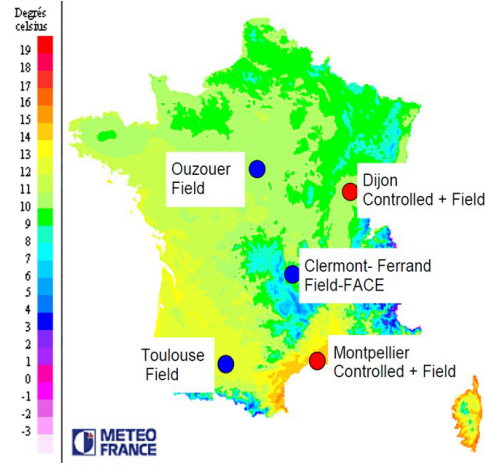


Figure 1: Location of Phenome platforms, superimposed on a map of mean temperature in France. Phenome platforms are representative of the variability of temperature. They also represent different risks of water deficit.

106 either the imaging or the watering units. The imaging unit consists of a  
 107 3-D image acquisition cabin with top and side channel. Five water units  
 108 consist of five weighing terminals and five high-precision watering pump-  
 109 stations, as shown in Figure 2. PhenoArch measures traits associated to the  
 110 plants' adaptation to climate change with a throughput of 1650 plants per  
 111 day. Typical measured variables include the timing of the plant cycle (leaf  
 112 appearance, duration of phenological phases), plant growth rate in terms  
 113 of area and volume, plant organ expansion and plant morphology (angles,  
 114 shape of leaves). The automated irrigation system allows to control various  
 115 water supply scenarios and estimates the responses of these traits to water  
 116 availability. Plants are imaged every day from 12 lateral and one apical view  
 117 (20000 images per day are produced), which allow reconstructing a digital  
 118 'avatar' of each plant of the platform.

119 Three main categories of workflows have to be executed: The first se-  
 120 ries of workflows is related to data acquisition in order to collect, describe,  
 121 and organize data sets while being acquired. The second category consists in  
 122 gathering, standardizing, and making available produced data sets. The third  
 123 category aims at finding answers to research questions: analyzing results ob-  
 124 tained and combining such data with other data sets to extract knowledge.  
 125 Workflows then need to combine highly heterogeneous data, from very dif-



Figure 2: Phenoarch Phenotyping Platform

126 ferent sources such as manual samplings, readings, and human observations  
 127 at different scales (populations, plants, organs, tissues, cells, etc.) and at  
 128 different times and stages. Such data can be either comparative (mutant  
 129 versus wild type), absolute (days to flowering of a cultivar) or relative (rel-  
 130 ative growth per day). Extensive connections to large sets of data types  
 131 are also mandatory (seed stocks, genes, experimental methods, publications,  
 132 etc.) leading to major data integration research questions [1] and calling for  
 133 a new generation of analysis tools.

## 134 2.2. Image analysis workflow

135 As an example, we describe one of the elementary workflows, that mostly  
 136 consists of an image analysis step targeting the estimation of plant leaf area.

137 Most of the raw data (images) are indeed not used directly, but processed  
 138 with an image analysis pipeline to get a trait, and then further analyzed with  
 139 a response model. Analyzing images is part of the important steps of the  
 140 experiments to be performed. It is shared between many workflows and used

141 to produce the traits measurements.

142 A very large variety of algorithms may be used to extract relevant infor-  
143 mation from image analysis. Dedicated plant phenotype commercial packages  
144 use basic functions to estimate total biovolume and leaf area for example.  
145 However current research provides a new landscape of algorithms. They make  
146 the link with ecophysiological models and allow to perform a more precise  
147 analysis of plant traits. In this context, it is particularly important to allow  
148 users to test and compare algorithms on their data sets.

149 Despite its simplicity, this use case already illustrates one important char-  
150 acteristics of phenotyping analysis, that is the intrinsic dependency between  
151 data and models. Consecutively such workflows can then be completed by  
152 other workflows that couple data analysis with a model for analyzing the  
153 response of plant expansion rate to temperature and water availability or  
154 with an integrative model, in a simulation context. Furthermore, such a  
155 kind of *in-silico* experiment can be considered in much wider contexts. For  
156 instance, after an initial segmentation of organs in the image, the global ar-  
157 chitecture of the plant can be reconstructed. This 3D reconstruction can  
158 be interfaced with canopy-level models of light interception to gain access  
159 to physiological parameters like intercepted light and radiation use efficiency  
160 for example. Hence both geometrical parameters attached to a plant (*e.g.*,  
161 leaf surface area) and physiological parameters (*e.g.*, photosynthesis) can be  
162 tracked throughout time and correlated with genotypes.

163 All these steps are particularly challenging and involve multi-disciplinary  
164 teams (biology, statistics, geometry, bioinformatics, computer science...).

### 165 2.3. User Requirements

166 In the introduction we have presented the three main high-level require-  
167 ments we followed to design InfraPhenoGrid: i) the ability for users to design  
168 and exchange experiments where a very large number of tools are interlinked  
169 (handled by a workflow management system), ii) the ability for users to re-  
170 produce experiments and understand the result obtained by such experiments  
171 (handled by a provenance layer), and iii) the ability to deal with large-scale  
172 experiments involving masses of data (handled by parallel computing envi-  
173 ronments). In this subsection we provide precision on the PhenoArch users'  
174 requirements.

175 **A Transparent, Familiar and Flexible Working Environment:** The  
176 classical users of the PhenoArch platform are bioinformaticians, mainly Python  
177 programmers (strongly involved in the design of Jupyter/IPython notebooks

178 [2, 3]), statisticians, image analysts and more generally modelers, all closely  
 179 connected to the Plant community. They are already very familiar with the  
 180 OpenAlea workflow system and in particular they are frequent users of some  
 181 analysis tools and libraries provided by OpenAlea. InfraPhenoGrid should  
 182 thus be designed to be as *transparent* as possible for users, that is, to allow  
 183 them continuing working in the same environment. However, we want our  
 184 infrastructure to be *flexible* to use other workflow systems and/or libraries  
 185 both for our current users to discover them and to welcome next generation  
 186 users.

187 **An Adaptable Operational Infrastructure:** Faced with the amount of  
 188 data to be analyzed, the computing infrastructure of InfraPhenoGrid has  
 189 to be designed in an *operational* distributed infrastructure, already used in  
 190 similar projects. While a National (European) and Open infrastructure has  
 191 to be favored in a first time, InfraPhenoGrid should be *adaptable* to both  
 192 Grid and Cloud solutions.

193 **Reproducibility-Friendly Infrastructure:** InfraPhenoGrid is a workflow  
 194 infrastructure for plant scientists to analyze their datasets and understand  
 195 them. Tracking data used and produced (Provenance) as well as the exact  
 196 description and environments where the tools have been executed is a crucial  
 197 need and should be done following international standards of the domain.  
 198 InfraPhenoGrid should thus be *Reproducibility-Friendly*, welcoming to any  
 199 plugins to export, visualize and analyze Provenance information and more  
 200 generally any tool to enhance reproducibility of experiments.

### 201 3. InfraPhenoGrid Architecture

202 The InfraPhenoGrid we designed to manage and analyze plant phenotyp-  
 203 ing data sets is an infrastructure based on a number of layers of abstractions.  
 204 First, computational methods for analysis and simulation are expressed by  
 205 means of *scientific workflows* (in the OpenAlea workflow system). Second,  
 206 a *middleware* (SciFloware) maps, manages and optimizes the execution of  
 207 scientific workflows on distributed environments. Third, a *provenance layer*  
 208 captures the workflow execution and reports it to users to further understand  
 209 and explore results of the computation. Last, a large scale infrastructure, the  
 210 *Grid* (France-Grilles) allows to have access to a shared, extensible and very  
 211 large computational power and storage.

212 France-Grilles makes use of two other important components, namely,  
 213 DIRAC and iRODS. DIRAC (Distributed Infrastructure with Remote Agent



Control) [4] is a framework for distributed computing particularly well-suited to deal with large communities of users. iRODS (integrated Rule-Oriented Data System) [5] is a scalable open-source data management software used by research organizations and government agencies worldwide. The focus of iRODS is data. It provides data discovery using a metadata catalog that describes every file, directory, and storage resource in the data grid. iRODS is also in charge of implementing data virtualization. (iRODS will be described in more details in Section 3.3.)

More precisely, the architecture of InfraPhenoGrid is depicted in Figure 3. Circled numbers are related to steps described here-after.

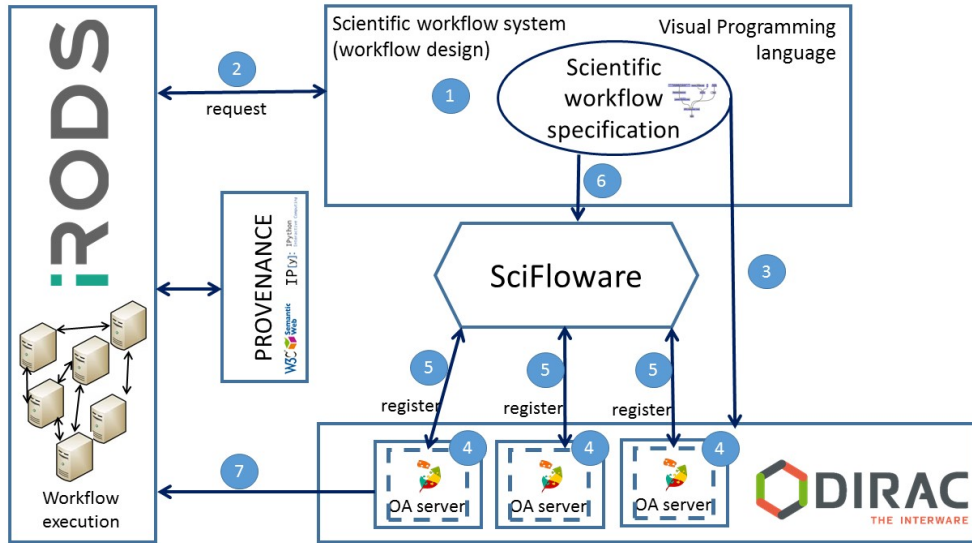


Figure 3: InfraPhenoGrid Architecture

The user interacts with InfraPhenoGrid in a visual programming environment (step 1) by designing a workflow specification from scratch or by selecting an existing workflow in the library of available workflows. To run the workflow, the user has to select the data sets to be taken as input by the workflow at execution time. The distributed infrastructure is thus transparent for the end-user.

Selecting a dataset in InfraPhenoGrid actually corresponds to sending a request to iRODS to concretely get the data (step 2). Resources have then to be allocated; this is performed by DIRAC (step 3). On each allocated worker

233 or job, OpenAlea workflows are deployed by copying an image from iRODS  
234 so that an OpenAlea server is launched (step 4). OpenAlea servers then  
235 register to SciFloware (step 5) which is in charge to distribute computations  
236 of possible subparts of the workflow to the different OpenAlea servers (step  
237 6). The workflow is then concretely executed on the Grid (step 7). At this  
238 stage, provenance information and all information on jobs are stored and  
239 available on iRODS.

240 The next subsections present with more details and discuss the choices we  
241 made on the major components of InfraPhenoGrid, namely, the OpenAlea  
242 workflow system, the SciFloware middleware and the data management and  
243 provenance layer.

### 244 3.1. Scientific Workflow Management system: OpenAlea

245 *OpenAlea - A system targeted to the plant community.* The OpenAlea sci-  
246 entific workflow system is a component-based architecture implemented as  
247 a set of pure Python packages [6]. The visual programming environment  
248 and graphical user interface (GUI) is implemented using the PyQt toolkit, a  
249 Python binding to the Qt application framework. OpenAlea is portable and  
250 available on Linux, Windows, and MacOS/X.

251 OpenAlea has been in constant use since 2004 by users of the French  
252 Plant science community but not only since the system has been downloaded  
253 618000 times and the web site counts international 10000 unique visitors a  
254 month according to the OpenAlea web repository (<https://gforge.inria.fr>).  
255 OpenAlea is distributed under a free software license (L-GPL) and main-  
256 tained and developed by a group of 20 active developers from different re-  
257 search institutes and universities. Development is performed under a collabo-  
258 rative scheme with shared methodologies (*best practices*). Coding sprints are  
259 regularly organized by various sub-groups of developers (pair programming  
260 and test driven development) and scientists (biologists and mathematicians).

261 OpenAlea tools (*e.g.*, models, workflows, components) are published and  
262 shared on the web both through the main OpenAlea web repository and  
263 through web sites of groups which use and contribute to OpenAlea without  
264 being concretely partners of the OpenAlea project<sup>2</sup>.

265 As a consequence, more than 60 researchers have contributed to Ope-  
266 nAlea packages, in France and internationally, published through large meta-

---

<sup>2</sup>See for example the following web sites: <http://www.stse-software.org> and  
<https://www.cpib.ac.uk/research/themes/digital-plant>

267 packages (e.g., *Alinea* to simulate ecophysiological and agronomical processes  
 268 and *VPlants* to analyze, model and simulate plant architecture and its de-  
 269 velopment) to ease the installation for end-users.

270 The strong and long-term experience of PhenoArch users and their inter-  
 271 national collaborators with both using and developing workflows in OpenAlea  
 272 made us choose OpenAlea as the workflow system for InfraPhenoGrid. The  
 273 main technical features of OpenAlea are described here after.

274 *Using OpenAlea (Designing workflows).* From an end-user point-of-view, the  
 275 first feature of OpenAlea exploited is its visual programming environment  
 276 (part A of Figure 4) where users are provided with a set of predefined work-  
 277 flows and libraries of tools (part B of Figure 4) to be combined to form new  
 278 workflows.

279 Users can create new wrapped tools by implementing them in Python  
 280 (in part C of Figure 4). Each tool and workflow is associated with some  
 281 documentation and saved. Ports of actors are typed and widgets can be  
 282 associated with data types to allow users interaction with the data.

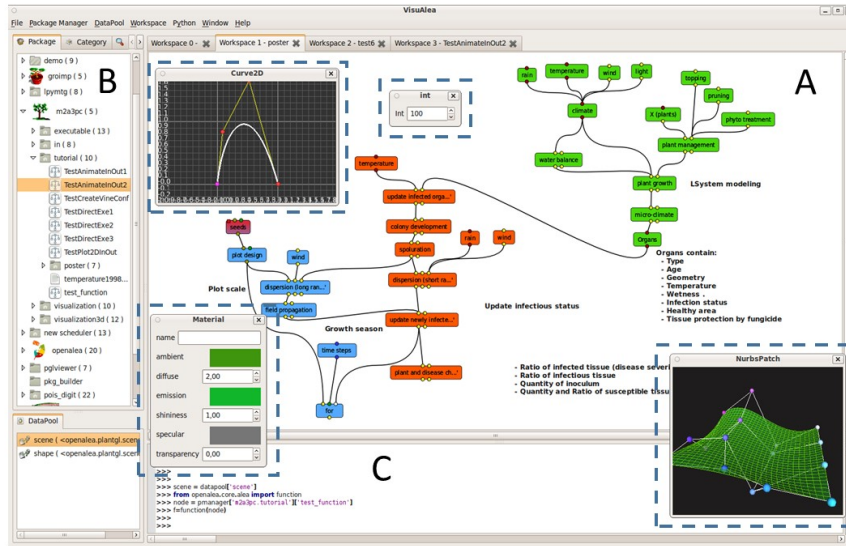


Figure 4: Main graphical user interface of OpenAlea. Users can design and interact with workflows in A. The package manager is in part B and provides users with the structured list of tools available and the list of existing workflows. On part C is the Python interpreter, where OpenAlea actors can be designed. Dotted lines denotes widgets.

283 *Workflow specification.* From a more formal point-of-view, in OpenAlea, a  
 284 workflow is classically represented as a directed multi-graph. Each node is  
 285 called an actor and represents a task to be executed (*a.k.a.* component or  
 286 activity). Each node has a name, a function object (a functor, a program, a  
 287 Web Service or a composite actor), and an explicitly defined set of input and  
 288 output ports. Directed edges are data links which connect output to input  
 289 ports.

290 While OpenAlea can be classically used to perform data analysis as in  
 291 other workflow systems such as Galaxy [7], Taverna [8] or Kepler [9], its orig-  
 292 inality lies in its ability to handle loops expressing retro-action [10]. In other  
 293 words, OpenAlea is able to deal with simulation and modeling. Iteration is  
 294 handled by introducing a specific kind of actor, called *dataflow variable X*. It  
 295 allows to specify that, at a given port, an actor receives an unbound variable  
 296 rather than a value. Connecting an X to an actor transforms a workflow into  
 297 a lambda function. It allows to express higher-order programming providing  
 298 control flow behavior using a set of algebraic operators. An algebraic opera-  
 299 tor is an actor that iterates over first-order function calls, and thus takes one  
 300 or more functions as inputs. Ports that require a function have an associated  
 301 semantic type *Function*.

302 More precisely, the **map** operator is a higher-order function  $map :: (\alpha \rightarrow$   
 303  $\beta) \rightarrow [\alpha] \rightarrow [\beta]$ . Its argument are a function  $f :: \alpha \rightarrow \beta$  (first port) and a  
 304 set of elements of type  $\alpha$  (second input port). The **map** operator applies  $f$  to  
 305 each element of the set and returns the set of resulting elements of type  $\beta$ .  
 306 In Part A of Figure 5, two implementations of the **map** operator are provided.  
 307 The left workflow illustrates the **map** operator running on one single processor  
 308 while the one on the right hand side illustrates the **parallel map** operator  
 309 with the same workflow running on 4 processors.

310 Similarly, the **reduce** operator takes a function  $g$  of two variables and  
 311 a sequence of elements  $[x_i]$  and returns one element. **while** is an iteration  
 312 operator that takes three inputs: an initial element  $t_0$ , a boolean function  
 313 *cond* and function  $h$ . It initializes a variable  $t$  with  $t_0$  and iteratively applies  
 314 the function  $h$  on  $t$  while *cond*( $t$ ) is true.

315 *Workflow execution.* The execution of a given workflow in OpenAlea is launched  
 316 in response to requests for data of one of its actors. Such an actor can sat-  
 317 isfy the request when the upstream sub-workflow has been executed, that is,  
 318 when all the relevant actors connected to its input ports have been executed.  
 319 When such an actor has received its data on its input ports, it executes and

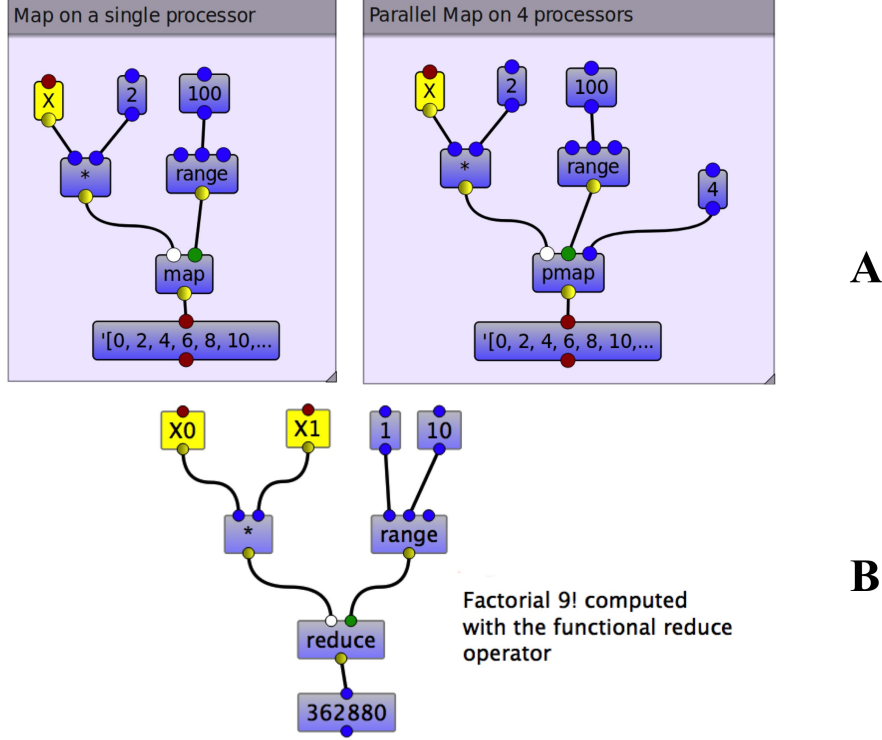


Figure 5: Algebraic operators map (A left), parallel map (A right), and reduce (B)

places data on its output ports.

OpenAlea has introduced  $\lambda$ -dataflow evaluation [10] which differs from the classical evaluation if the workflow contains at least one *dataflow variable*  $X$ . The execution is then decomposed into two stages. First, for each port of type *Function*, a sub-workflow is computed if the upstream sub-workflow contains at least one dataflow variable. This sub-workflow is defined by all the actors needed to produce the data on this port, *i.e.* the upstream sub-workflow and the connected output port. This sub-workflow is dynamically transformed into a function (*i.e.* an actor) of one or several variables corresponding to its dataflow variables. Second, the evaluation of this function by algebraic operators consists in replacing the variables by real data and evaluating the sub-workflow.

Additionally, OpenAlea provides several optimizations in the orchestration of the workflow execution by allowing actors to be *blocked* and *lazy*. If

334 an actor is blocked, the execution is not propagated to the upstream sub-  
335 workflow and if the actor is lazy, the execution is performed only if the actor's  
336 inputs have not changed compared to its previous execution. This type of  
337 orchestration performs only the operations needed to produce the required  
338 result, executing the subset of the graph relevant to the output [11].

### 339 3.2. A middleware for parallel environments: *SciFloware*

340 *The choice of SciFloware.* As stated in the user requirements, InfraPhenoGrid  
341 should be equipped with a system able to (i) hide the complexity of the com-  
342 putation and offer *transparent* access to data and tools for end-users, (ii)  
343 provide a *flexible working environment* by allowing several workflow systems  
344 to be used, (iii) be *adaptable* to pilot tasks both on Grid and Cloud infras-  
345 tructures. Using and tuning such kind of systems remains a very difficult  
346 task.

347 Our research groups have recently developed *SciFloware*, a generic lightweight  
348 middleware able to coordinate and pilot the tasks to be executed in a trans-  
349 parent way for the user. *SciFloware* is based on the Shared-data Overlay  
350 Network [12] (SON) which follows a Software as a Service (SaaS) model,  
351 eliminating the need to install and maintain the software and allows users to  
352 run HPC programs through graphical interfaces (*e.g.*, graphical interfaces of  
353 scientific workflow management systems).

354 *SciFloware* has been chosen to be the InfraPhenoGrid middleware, be-  
355 cause it was a system we are familiar with and it matched our demands on  
356 requirements. The technical features of *SciFloware* are described here after.

357 *Internal representation of workflows in SciFloware.* SciFloware has been de-  
358 signed to allow interoperability of different workflow systems. In absence of  
359 standards to represent scientific workflows<sup>3</sup>, SciFloware defines its own XML  
360 workflow specification to describe a *master workflow*. A master workflow is  
361 a meta-level workflow used to orchestrate and compose concrete workflows.  
362 Each workflow is run independently by different scientific workflow systems.  
363 As an example, in Figure 6, the SciFloware master workflow consists of four  
364 steps, including steps *a* and *b* (sub-workflows of the master workflow) which  
365 are respectively executed by the *X workflow system* (for the *a sub-workflow*)  
366 and the *OpenAlea* workflow system (for the *b sub-workflow*). More precisely,

---

<sup>3</sup>The Common Workflow Language (CWL) Initiative may become a solution in the future but it has not reached the right level of maturity yet.

367 SciFloware is responsible for (i) sending to workflow systems the execution  
 368 of such two master sub-workflows with the right input data produced by the  
 369 previous executed master workflow step, (ii) collecting output data gener-  
 370 ated at the end of each execution of the sub-workflows, and (iii) launching  
 371 the execution of the last step of the master workflow with such collected  
 372 data.

373 While OpenAlea has been the first workflow system orchestrated by Sci-  
 374 Floware (based on our user requirements), Galaxy [7] and Taverna [8] are  
 375 currently under consideration (to play the role as "X workflow system" in  
 376 Figure 6).

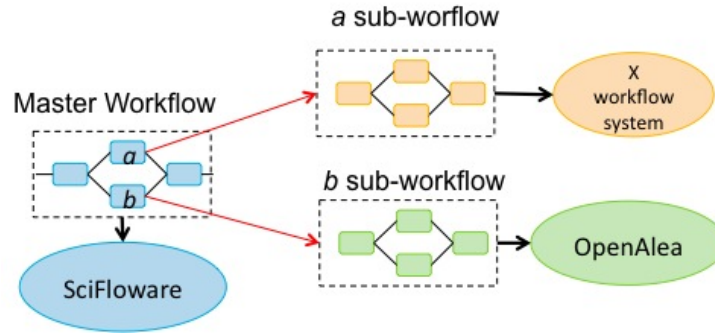


Figure 6: SciFloware distributed middleware.

377 *Algebraic expressions in SciFloware.* SciFloware uses a relational data model  
 378 and an algebraic language to represent data-intensive scientific workflows [13].  
 379 Data flows are represented as relations and workflow nodes and activities  
 380 as algebraic expressions. A relation is a set of tuples composed of basic  
 381 attributes (*e.g.*, int, float, string, file references).

382 An algebraic expression consists of algebraic activities, additional operands,  
 383 operators, input relations and output relations. It is comprised of a workflow,  
 384 a program or an SQL expression, with input and output relation schemes.

385 More precisely, SciFloware provides six algebraic operators: Map, SplitMap,  
 386 Reduce, Filter, SRQuery and MRQuery. The semantics of these operators  
 387 has been defined in [13].

388 *Software as a Service, communication with workflow systems.* As previously  
 389 stated, SciFloware is particularly modular, following the SAAS ("Software

390 as a Service”) principle. SciFloware is a service-and-component-based dis-  
391 tributed architecture and built from software components. Each compo-  
392 nent provides services to other components of the workflow. A service is a  
393 self-contained unit of functionality which exposes components through well-  
394 defined interfaces. SciFloware provides services that can be combined to build  
395 large distributed applications. In InfraPhenoGrid, SciFloware runs on a ded-  
396 icated server and provides a registration service, where distributed workers  
397 can register themselves. Using a communication protocol, SciFloware dis-  
398 tributes the computation among workers, with each worker running behind  
399 a dedicated server. The tech stack (*e.g.*, authentication) is implemented  
400 using services and can easily be evolved further.

401 More precisely, three main kinds of services are considered in SciFloware:  
402 algebraic operators, scientific workflow systems and the communication pro-  
403 tocol between algebraic operators and workflow systems.

404 SciFloware schedules the computation using algebraic operators. Services  
405 associated with these are managed by the SciFloware server. Each compo-  
406 nent has its own decentralized execution strategy, which allows to simply  
407 distribute the execution on several sites. A component type is associated  
408 with each algebraic operator. For example, the *Map* operator has an associ-  
409 ated *map* type of component, with an input and output service for relations  
410 and a service to schedule activity among workers.

411 As for workflow systems, each worker runs a service, running a workflow  
412 within a given workflow system. It receives a workflow description and input  
413 data IDs and returns the produced output data. On a Grid infrastructure,  
414 workers will run on different nodes, each one executing its own server to  
415 communicate with SciFloware.

416 Eventually, a set of services have been defined to manage communication  
417 between SciFloware and different workflow systems. Using these, SciFloware  
418 can request the execution of a workflow on any workflow system server, or be  
419 notified at the end of the execution of an OpenAlea workflow. Other services  
420 manage database relationships and communication protocols.

421 Figure 7 illustrates a case where a component manages local task execu-  
422 tions depending on a given algebraic operator. More precisely, in Figure 7,  
423 the Map component distributes the computation among workers (OA servers,  
424 that is, OpenAlea servers) depending on its own scheduling policy.

425 SciFloware allows new components to be added and instantiated dynam-  
426 ically to extend the middleware.



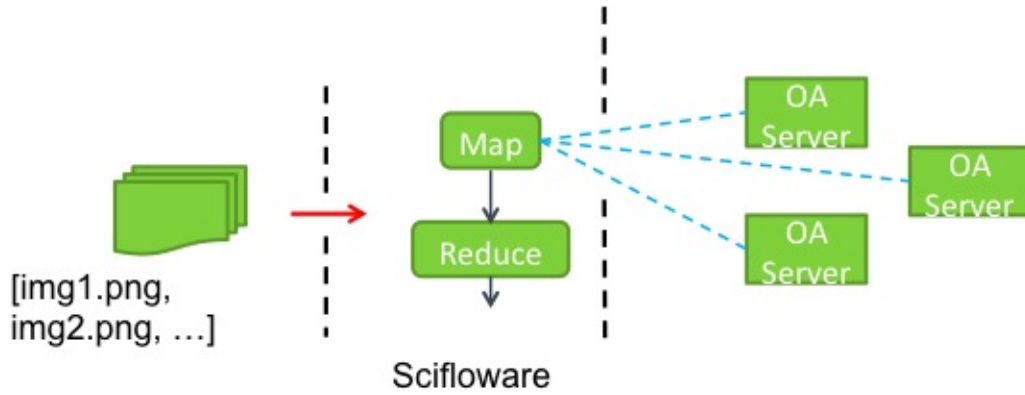


Figure 7: Distributed execution of a scientific workflow on the Grid with SciFloware.

427 *SciFloware Architecture.* The SciFloware architecture is composed of the fol-  
 428 lowing components (see Figure 8):

- 429 • **Algebraic Operators.** For each algebraic operator, a specific compo-  
 430 nent type is associated, which can be instantiated several times during  
 431 the execution of a given workflow.
- 432 • **Execution Manager.** The Execution Manager takes in a workflow  
 433 specification and instantiates the needed components (*i.e.* workflows  
 434 or part-of workflows), such as the algebraic components. A message is  
 435 sent to trigger execution of the first workflow component. The Execu-  
 436 tion Manager also manages and runs the available OpenAlea servers.  
 437 At each registration of an instance of an OpenAlea server, a worker  
 438 component is added to the list of available workers.
- 439 • **Data Manager.** The Data Manager provides workers with access to  
 440 the storage database.

441 *Implementation of SciFloware.* The implementation of SciFloware is based  
 442 on the Shared-data Overlay Network (SON) [12], written in JAVA, fully  
 443 integrated into the Eclipse environment and implemented on top of OSGi<sup>4</sup>.

---

<sup>4</sup>The OSGi specification describes a modular system and a service platform for the Java programming language that implements a complete and dynamic component model.

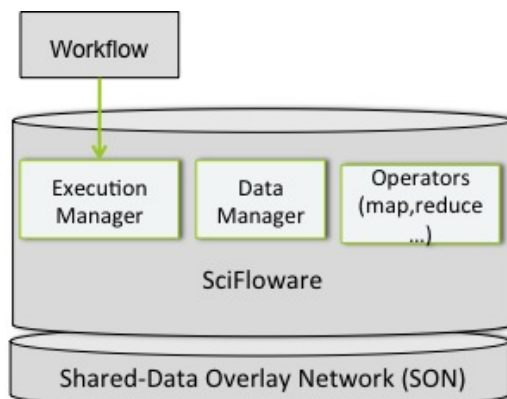


Figure 8: SciFloware architecture.

444 The communication protocol between SciFloware and scientific workflow  
 445 systems being run on different workers use a REST (Representational State  
 446 Transfer) interface. The RESTfull communication protocol is used to regis-  
 447 ter new workers, exchange the workflow specification and the data identifiers  
 448 stored on iRODS and start and stop the execution. Each service in Sci-  
 449 Floware is described using description files similar to the Web Service Def-  
 450 inition Language (WSDL). Each SciFloware component has an associated  
 451 description file defining the required and provided services.

452 SON is used for the internal composition of SciFloware components. SON  
 453 allows to build an application following the SAAS model (Software as a Ser-  
 454 vice) using a set of components that can be executed in a distributed way  
 455 (*e.g.*, Peer-to-Peer) on a Grid or on a Cloud infrastructure. The SON middle-  
 456 ware allows to define new component types by specifying the services offered  
 457 by each type. During execution it allows to dynamically combine different  
 458 component instances and to manage the life cycle of these components (cre-  
 459 ate, init, stop,...).

460 Concerning the authorization framework, SciFloware uses OAuth2 [14] to  
 461 enable safe registration of the OpenAlea servers to SciFloware. The OAuth2  
 462 flow used is the Resource Owner Password Credentials Grant flow<sup>5</sup>. During  
 463 the deployment of the OpenAlea server, the user provides its SciFloware's  
 464 username and password. At startup, each OpenAlea server requests an access

---

<sup>5</sup>RFC 6749 - <http://tools.ietf.org/html/rfc6749>

465 token from SciFloware using login authentication. After registration, all  
466 the communication between SciFloware and OpenAlea is mediated through  
467 iRODS, whose access requires authentication based on certificates delivered  
468 by the French Grid Certification Authority.

469 As for the Data Manager component, it has been extended to support  
470 iRODS. The connection with iRODS is established using the REST interface.  
471 iRODS is the focus of the next subsection.

472 *Packaging OpenAlea into a virtual environment.* OpenAlea and all its de-  
473 pendencies have been built and packaged into a virtual environment [15] on  
474 a virtual machine. The operating system (*i.e.*, a scientific Linux version 6)  
475 of the virtual machine is the same as the one deployed on each worker of  
476 the Grid. The virtual environment is packed and stored on iRODS. When  
477 a worker is reserved by DIRAC, the bundle is uploaded locally and uncom-  
478 pressed. A shell script updates the environment variables and an OpenAlea  
479 server is launched. This method has been preferred to virtualization (*e.g.*,  
480 vagrant or docker [16]) for performance reasons and because all the workers  
481 have the same operating system. The size of the compressed bundle con-  
482 taining all the packages and their dependencies is 210 MB. The latency (or  
483 delay) due to the copy of the bundle from iRODS and its installation is less  
484 than 1 min, while each worker is deployed once for a maximum of 24 hours.

### 485 3.3. Data management and Provenance

486 *Data management with iRODS.* As previously stated, iRODS (v3.3) is an  
487 open-source data system chosen by France-Grilles for its ability to provide a  
488 technology enabling data and policy virtualization for multiple and geograph-  
489 ically separated users [17]. iRODS federates distributed and heterogeneous  
490 data resources into a single logical file system and provides a modular in-  
491 terface to integrate new client-side applications. iRODS does not only allow  
492 the worker nodes of the grid infrastructure to access data sets but it provides  
493 end-users with access to data, through GUI, while enabling user to annotate  
494 such datasets with rich metadata.

495 Both input data and the results of a scientific workflow computation as-  
496 sociated with their provenance are stored on iRODS. This drastically reduces  
497 the volume of data to be transferred through SciFloware to launch a compu-  
498 tation and retrieve its result (only the address of data in the central catalog  
499 has to be exchanged). We also implemented a communication protocol by file  
500 to bypass the limitations enforced by the Grid on network communications.

501 *Provenance Layer.* One of the major aims of our infrastructure is to be  
502 *reproducibility-friendly* [18]. The starting point to make a scientific result  
503 reproducible is to keep track of the exact data sets and exact tools (including  
504 parameter settings) used to obtain a given data item. To answer such  
505 needs a *Provenance Layer* has been designed and implemented in our infrastructure.  
506 A layer currently has two main components: a provenance model,  
507 based on the W3C standard PROV [19] and a notebook generator, able  
508 to automatically generate notebooks from some workflow executions. The  
509 Provenance layer is flexible in the sense that new modules taking in PROV  
510 data and making it possible to visualize and analyze provenance information  
511 can be integrated.

512 While iRODS is in charge of concretely managing data, the provenance  
513 module reconstructs (by querying iRODS) the history of each data item. In  
514 other words, the provenance module is able to provide for each produced  
515 data item the exact series of workflow node executions, including the data  
516 sets used as input of such nodes. Such provenance information is represented  
517 according to the W3C PROV standard [19]. Both prospective (the workflow  
518 specification) and retrospective (execution and data sets) provenance are  
519 stored.

520 *Notebooks generator.* The second component of the provenance layer aims  
521 at helping the PhenoArch users understand the results they obtained by  
522 allowing them to visualize and interact with the main steps of the process  
523 used to produce such results. In other words, we want PhenoArch users to  
524 be able to interact and visualize (part of) the execution of some workflows  
525 to follow how some final results have been obtained.

526 The current solution used by an increasing number of scientists to answer  
527 such kinds of need is to (manually) design notebooks using the Jupyter/IPython  
528 Notebook web application [2]. From a developer point-of-view, a notebook  
529 is a JSON document (convertible into a number of open standard output  
530 formats including HTML, LaTeX, PDF ... and which can be designed using  
531 a web-base user interface) containing an ordered list of input/output cells  
532 with iPython code able to generate text, mathematics, plots and rich media  
533 (images, video...). From an end-user point-of-view, a notebook is a rich web  
534 page, where code, text, mathematics, plots and multimedia (images, video...) can  
535 be displayed and interacted with. In particular any user can modify the  
536 input of a notebook cell to observe the impact of this change on the objects  
537 displayed (*e.g.*, plots).

538 In InfraPhenoGrid, we implement a *notebook generator* where a set of  
 539 OpenAlea workflow executions are automatically converted into Jupyter/  
 540 IPython notebooks. More precisely, the generator (i) converts each Ope-  
 541 nAlea workflow actor (natively coded in Python) into an IPython cell and  
 542 (ii) queries iRODS to extract automatically information on input and output  
 543 datasets respectively used and produced by the execution of each OpenAlea  
 544 actor. Users can then visualize and interact with data used and produced  
 545 during an execution. A concrete example of generated notebooks is provided  
 546 in the next Section (Results).

## 547 4. Results

548 In this section, we present the benefits of using InfraPhenoGrid in Plant  
 549 Phenotyping by showing how the various components of our infrastructure  
 550 make it possible to perform complex experiments on huge data sets.

551 *Designing workflow.* Figure 9 provides an example of a workflow designed  
 552 and executed by our end-users to estimate the growth of a plant. Such a  
 553 workflow starts with querying iRODS to get input data ("import image"  
 554 node). In this concrete example, the 1407th individual plant of genotype  
 555 A310 has been considered. As a consequence, "1407" and "A310" appears  
 556 in the name of the actors of the workflow. The *import image* node returns  
 557 a time series of images. Each acquisition, taken every two days, records  
 558 pictures of the plant at various angles by rotating the plant. The node *keys*  
 559 returns the order sets of the dates, while *values* returns a list of images  
 560 corresponding to the different side views of the growing plant, along time.  
 561 The *map* node applies the sub-workflow illustrated on Part B of Figure 9 on  
 562 each set of images of the time series. The result is a list of estimated areas  
 563 of the individual plant over time.

564 The *X* node represents a dataflow variable and abstracts the sub-workflow  
 565 (part B) as a function. The sub-workflow B implements a mean-shift algo-  
 566 rithm. It receives the multiple views of the plant at a given time. All these  
 567 views are combined (node *cv.mean*) using an OpenCV algorithm to separate  
 568 the plant from the cabin background in the image. The *macro\_side\_binarization*  
 569 node subtracts the cabin background of each image and the "green" pixels are  
 570 counted (*countNonZero*). Again, the *map* algebraic operator is used to run  
 571 the same treatment on a set of images. Note that the mean shift algorithm  
 572 is only computed once due to lazy evaluation. The binary image (Part D)

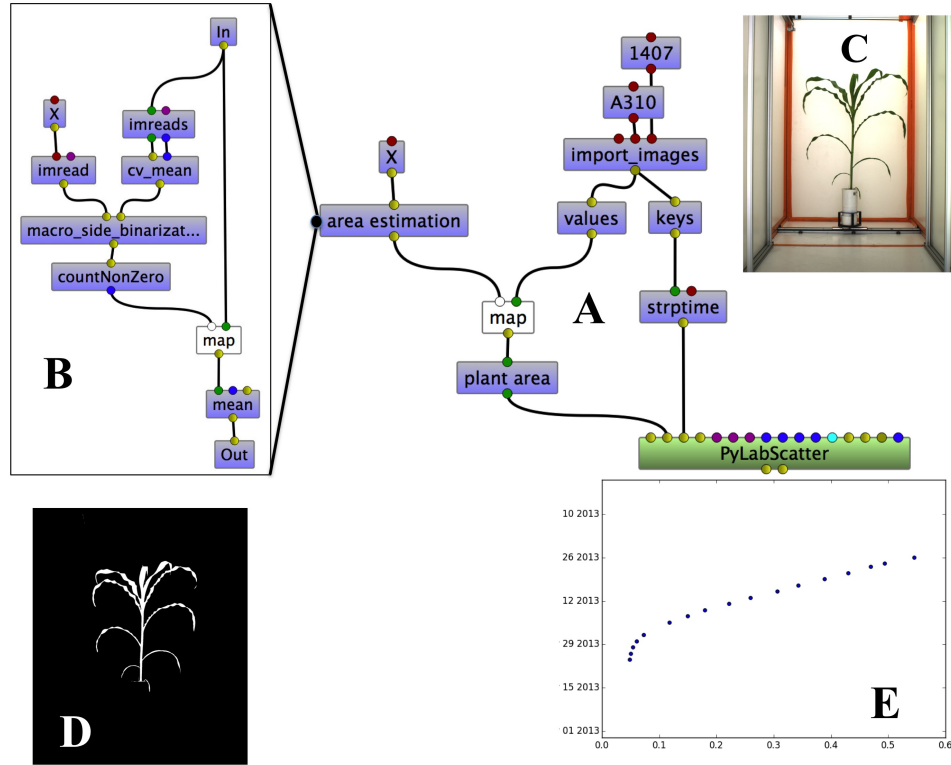


Figure 9: Workflow. Part A represents the main workflow while Part B is a sub-workflow corresponding to the node "area estimation" of workflow A. Part C depicts one of the set of real images of a plant of a given genotype obtained from the imaging system of the phenoarch platform. D represents the binary image. E plots the growth of the plant.

573 is produced by one execution of the *macro\_side\_binarization* for each lateral  
 574 perspective image using an optimized HSV segmentation algorithm.

575 Finally, the *plant area* node receives as input a number of "green pixels"  
 576 for each plant and estimates, using a linear model, the plant area along its  
 577 growth and development. The *PyLabScatter* node plots the graph of the  
 578 growth of the plant (Part E) using the Matplotlib library [20] wrapped  
 579 within OpenAlea.

580 *Exploring alternative methods.* OpenAlea is modular and allows users to eas-  
 581 ily test various alternative methods. In particular, several variations of sub-  
 582 workflow B can be designed, resulting in the use of several Binarization al-  
 583 gorithms (as mentioned in Table 1). The results obtained by such variations

Binarization algorithms	Time (seconds)
Adaptive Threshold [21]	62.1
HSV [22]	85.7
Mean Shift [23]	73.9

Table 1: Execution time of the workflow illustrated in Figure 9 with variations of sub-workflow B on one plant measured during one seasonal growth (5 weeks). This experiment gathered 124 images.

584 can then be compared (mainly qualitatively) by the user.

585 *Exploiting the Grid.* Execution times of Binarization algorithms reported in  
586 Table 1 are related to one single plant of one single genotype for which 13  
587 images have been collected during one month. The challenge then lies in  
588 considering 300 genotypes, with 1900 plants in each genotype. Per day, a  
589 PhenoArch platform produces 20000 images of 50M, equivalent to 1To/day,  
590 5To/week and 250 To/year. Without any distributed infrastructure, process-  
591 ing these huge amounts of data would take between 409 days and 565 days  
592 (depending on the strategy followed for the sub-workflow B).

593 For this project, France-Grilles provides us with 32 000 logical process  
594 units. With only a subset of this resource (2%), the whole computation,  
595 scheduled by SciFloware, can be performed in less than 12 hours (night time).

596 *Exploiting Provenance.* During any execution, provenance data, that is, all  
597 the data items processed and produced (including intermediate and final  
598 images) are stored on iRODS. Based on this provenance information, IPython  
599 notebooks are automatically generated for each individual plant either at a  
600 given time, or for the entire growth period, or for a given genotype. Each  
601 cell of the notebook contains the script of the workflow node executed. The  
602 input/output data are direct references to the data produced and stored in  
603 iRODS.

604 Very interestingly, users can upload a given notebook (see Figure 10)  
605 with the corresponding data, and modify the execution parameters directly  
606 on their computer to visualize the impact of such modifications and better  
607 understand their results. Thus, biologists can explore the obtained results a  
608 posteriori to discover for instance why some outliers on the "growth curve"  
609 (see Figure 9E) have appeared. Reasons for this situation may actually be  
610 numerous, including problems occurring during acquisition (*e.g.*, the plant

may have fallen...), limitation of the method used (bugs in the implementation...) or wrong set of parameters.

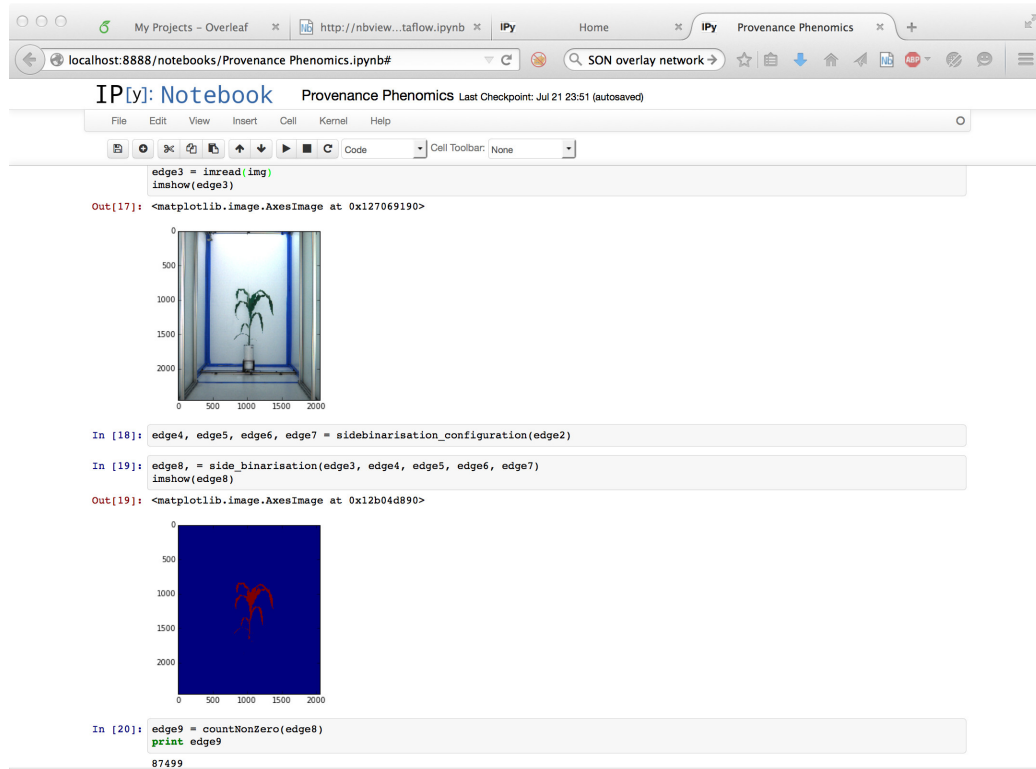


Figure 10: Notebook generated from the provenance of the execution of the workflow illustrated in Figure 9.B. Each node corresponds to a cell containing the equivalent Python function. Values flowing through edges have a name 'edge' with a value captured by the provenance module and stored into iRODS. The two images displayed correspond to the plant before and after its binarization. The total number of pixels of the plant is 87 489, which is the last cell value.

## 5. Discussion

In the last ten years, several approaches have been designed to distribute scientific workflows on parallel environments. Survey papers and books include [24, 25, 26, 27]. While cloud computing is increasingly used to manage Life Science data, our project involves large numbers of production sites and collaborating users, and such numbers are growing over time, making Grid technologies particularly well-suited [28].



620 More precisely, we wrap-up and discuss the four key aspects of our ap-  
621 proach.

622 First of all, the infrastructure we introduce in this paper works on a  
623 production environment with real and huge data sets produced on a daily  
624 basis. In Plant phenotyping, which is a field of growing importance in the  
625 context of climate change and food security [29, 30, 31], complex data sets  
626 produced by high-throughput image based phenotyping platforms need to  
627 be combined with highly heterogeneous data. In particular, compared to  
628 classical bioinformatics (largely driven by molecular biology), the need in  
629 plant phenomics is on modeling, simulation, designing statistical approaches,  
630 and performing complex image analysis. This poses new challenges in data  
631 integration and calls for new kinds of analytical tools [32, 33].

632 Second, the workflow system we use is well-established in the Plant com-  
633 munity. It is well-known by our first generation users, and has very specific  
634 features, making it possible to perform both data analysis and simulation  
635 (retro-action loops, modeling) while allowing visualization of complex data.  
636 As shown previously in [10, 6], OpenAlea belongs to the family of functional  
637 workflow systems (such as [9, 34, 35, 36]). It provides a unique solution,  
638 which is able to extend the dataflow model of computation by introducing  
639 higher-order language constructs in a visual programming environment, thus  
640 allowing to design highly expressive workflows in a fully uniform way.

641 Third, while the context of this work is data intensive and involves very  
642 complex experiments on huge data sets, we use a middleware approach to  
643 make the distribution and coordination of jobs transparent to the user. Sci-  
644 Floware is data-driven [37] in the sense that its internal workflow language  
645 clearly separates the definition of data to be processed from the graph of  
646 activities to be applied to the data. This separation is particularly suited to  
647 scientific workflows where the same experiment has to be reused for analyzing  
648 different data sets without any change. Optimization in SciFloware focuses  
649 on two main aspects: It uses asynchronous messages to execute workflows  
650 on a distributed infrastructure such as Grid or Cloud [37, 38]. Furthermore,  
651 while it uses generally coarse grain parallelism, it is able to exploit the fact  
652 that some workflow actors may be algebraic operators (*i.e.*, some actors are  
653 not black-boxes) to optimize the workflow execution. Additionally, basing  
654 the SciFloware implementation on the middleware SON [12] which follows the  
655 SAAS concept (Software as a Service) makes it very modular and flexible.

656 Last but not least, we have developed a large series of reproducibility-  
657 friendly features. Our approach allows users to share their experiments in

658 the spirit of [39], understand and compare their results [40] and possibly  
659 refine their analysis process to augment quality of their data sets. To do so,  
660 we have followed the recommendations and current standards on provenance  
661 [41] [19] and introduced a generator of IPython/Jupyter notebooks [2].

## 662 6. Conclusion

663 High-throughput phenotyping platforms provide a unique and particu-  
664 larly novel kind of solution to study the behavior of plants in context of  
665 climate change and food security. At the same time, size and complexity of  
666 data sets produced by such platforms are huge, constantly augmenting and  
667 the experiments to be performed are becoming increasingly complex, posing  
668 particularly novel challenges.

669 This paper introduces the InfraPhenoGrid infrastructure we designed and  
670 deployed to efficiently manage the data sets produced by the PhenoArch plant  
671 phenomics platform in the context of the Phenome Project.

672 Our solution consists in deploying scientific workflows on a Grid (France-  
673 Grilles) using a middleware to pilot workflow executions. InfraPhenoGrid  
674 is user-friendly in the sense that despite the intrinsic complexity of the in-  
675 frastructure, running scientific workflows and understanding results obtained  
676 (using provenance information) is kept as simple as possible for end-users.

677 Future work includes considering automatic transformation of scripts de-  
678 signed by scientists into OpenAlea scientific workflows to augment the size  
679 of the workflow library. We are also working on techniques to augment the  
680 reuse of workflows (and scripts) by guiding the design of such experiments  
681 ([42, 43]).

682 Another very important point we are actively working on is the repro-  
683 ducibility of scientific results. From the data point-of-view, one of the main  
684 challenge lies in finding the right level of granularity at which visualizing and,  
685 to some extent, recording the data [44]. Currently, the complete data sets  
686 are kept while we investigate several techniques (inspired from [42]) to re-  
687 duce the amount of data stored while ensuring a good level of reproducibility.  
688 From the environment point-of-view, we are currently able to consider virtual  
689 machine techniques to reproduce a given experiment in the exact same con-  
690 ditions (same OS, software versions...). Ongoing work includes considering  
691 techniques to re-execute experiments in new environments, where upgraded  
692 versions of software are considered [45].

## 7. Acknowledgement

The authors acknowledge the support of France-Grilles for providing computing resources on the French National Grid Infrastructure. This work has been performed in the context of the IBC (Institute of Computational Biology) in Montpellier, France. MM has received the support of the EU in the framework of the Marie-Curie FP7 COFUND People Programme, through the award of an AgreeSkills fellowship under grant agreement n 267196. Authors would like to thank Julien Coste (from INRIA) for his help in the OpenAlea server.

- [1] S. Cohen-Boulakia, U. Leser, Next generation data integration for life sciences, in: Proc. of the 25th Int. Conf. on Data Engineering (ICDE), IEEE, pp. 1366–1369.
- [2] H. Shen, Interactive notebooks: Sharing the code., Nature 515 (2014) 151–152.
- [3] M. Ragan-Kelley, F. Perez, B. Granger, T. Kluyver, P. Ivanov, J. Frederic, M. Bussonier, The jupyter/ipython architecture: a unified view of computational research, from interactive exploration to communication and publication., in: AGU Fall Meeting Abstracts, volume 1, p. 07.
- [4] A. Tsaregorodtsev, M. Bargiotti, N. Brook, A. C. Ramo, G. Castellani, P. Charpentier, C. Cioffi, J. Closier, R. G. Diaz, G. Kuznetsov, Y. Y. Li, R. Nandakumar, S. Paterson, R. Santinelli, A. C. Smith, M. S. Miguelez, S. G. Jimenez, Dirac: a community grid solution, Journal of Physics: Conference Series 119 (2008).
- [5] M. Hedges, A. Hasan, T. Blanke, Management and preservation of research data with irods, in: Proceedings of the ACM first workshop on CyberInfrastructure: information management in eScience, ACM, pp. 17–22.
- [6] C. Pradal, S. Dufour-Kowalski, F. Boudon, C. Fournier, C. Godin, Openalea: a visual programming and component-based software platform for plant modelling, Functional plant biology 35 (2008) 751–760.
- [7] J. Goecks, A. Nekrutenko, J. Taylor, Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences., Genome Biology 11 (2010) R86.

- 726 [8] P. Missier, S. Soiland-Reyes, S. Owen, W. Tan, A. Nenadic, I. Dunlop,  
727 A. Williams, T. Oinn, C. Goble, Taverna, reloaded, in: M. Gertz,  
728 T. Hey, B. Ludaescher (Eds.), Proceedings of the 22nd International  
729 Conference on Scientific and Statistical Database Management (SS-  
730 DBM), Heidelberg, Germany.
- 731 [9] B. Ludäscher, I. Altintas, On providing declarative design and pro-  
732 gramming constructs for scientific workflows based on process networks  
733 (2003).
- 734 [10] C. Pradal, C. Fournier, P. Valduriez, S. Cohen-Boulakia, Openalea: sci-  
735 entific workflows combining data analysis and simulation, in: Proceed-  
736 ings of the 27th International Conference on Scientific and Statistical  
737 Database Management (SSDBM), ACM, p. 11.
- 738 [11] V. Curcin, M. Ghanem, Scientific workflow systems-can one size fit all?,  
739 in: Proc. of Biomedical Engineering Conference, pp. 1–9.
- 740 [12] A. A. Lahcen, D. Parigot, A lightweight middleware for developing p2p  
741 applications with component and service-based principles, in: Compu-  
742 tational Science and Engineering (CSE), 2012 IEEE 15th International  
743 Conference on, IEEE, pp. 9–16.
- 744 [13] E. Ogasawara, J. Dias, D. Oliveira, F. Porto, P. Valduriez, M. Mattoso,  
745 An algebraic approach for data-centric scientific workflows, Proc. of  
746 VLDB Endowment 4 (2011) 1328–1339.
- 747 [14] D. Hardt, The oauth 2.0 authorization framework (2012).
- 748 [15] P. Guo, Cde: A tool for creating portable experimental software pack-  
749 ages, Computing in Science and Engineering 14 (2012) 32–35.
- 750 [16] D. Merkel, Docker: lightweight linux containers for consistent develop-  
751 ment and deployment, Linux Journal 2014 (2014) 2.
- 752 [17] A. Rajasekar, R. Moore, C.-y. Hou, C. A. Lee, R. Marciano, A. de Torcy,  
753 M. Wan, W. Schroeder, S.-Y. Chen, L. Gilbert, et al., irods primer:  
754 integrated rule-oriented data system, Synthesis Lectures on Information  
755 Concepts, Retrieval, and Services 2 (2010) 1–143.

- 756 [18] G. Sandve, A. Nekrutenko, J. Taylor, E. Hovig, Ten simple rules for  
757 reproducible computational research, *PLoS comp. biology* 9 (2013)  
758 e1003285.
- 759 [19] L. Moreau, P. Missier, *Prov-dm: The prov data model* (2013).
- 760 [20] J. D. Hunter, *Matplotlib: A 2d graphics environment*, *Computing In*  
761 *Science & Engineering* 9 (2007) 90–95.
- 762 [21] E. Navon, O. Miller, A. Averbuch, Color image segmentation based on  
763 adaptive local thresholds, *Image and vision computing* 23 (2005) 69–85.
- 764 [22] S. Sural, G. Qian, S. Pramanik, Segmentation and histogram generation  
765 using the hsv color space for image retrieval, in: *Image Processing. 2002.*  
766 *Proceedings. 2002 International Conference on*, volume 2, IEEE, pp. II–  
767 589.
- 768 [23] D. Comaniciu, P. Meer, Mean shift: A robust approach toward fea-  
769 ture space analysis, *Pattern Analysis and Machine Intelligence, IEEE*  
770 *Transactions on* 24 (2002) 603–619.
- 771 [24] J. Liu, E. Pacitti, P. Valduriez, M. Mattoso, A survey of data-intensive  
772 scientific workflow management, *Journal of Grid Computing* (2015) 1–  
773 37.
- 774 [25] I. Foster, Y. Zhao, I. Raicu, S. Lu, Cloud computing and grid computing  
775 360-degree compared, in: *Grid Computing Environments Workshop,*  
776 2008. GCE’08, Ieee, pp. 1–10.
- 777 [26] J. Yu, R. Buyya, A taxonomy of workflow management systems for grid  
778 computing, *Journal of Grid Computing* 3 (2005) 171–200.
- 779 [27] G. C. Fox, D. Gannon, *Workflow in grid systems*, Wiley Interscience,  
780 2006.
- 781 [28] M. T. Özsu, P. Valduriez, *Principles of distributed database systems*,  
782 Springer Science & Business Media, 2011.
- 783 [29] F. Tardieu, R. Tuberosa, Dissection and modelling of abiotic stress  
784 tolerance in plants, *Current opinion in plant biology* 13 (2010) 206–212.

- 785 [30] R. T. Furbank, M. Tester, Phenomics—technologies to relieve the phe-  
786 notyping bottleneck, *Trends in plant science* 16 (2011) 635–644.
- 787 [31] D. Houle, D. R. Govindaraju, S. Omholt, Phenomics: the next challenge,  
788 *Nature Reviews Genetics* 11 (2010) 855–866.
- 789 [32] F. Fiorani, U. Schurr, Future scenarios for plant phenotyping, *Annual*  
790 *review of plant biology* 64 (2013) 267–291.
- 791 [33] S. Dhondt, N. Wuyts, D. Inzé, Cell to whole-plant phenotyping: the  
792 best is yet to come, *Trends in plant science* 18 (2013) 428–439.
- 793 [34] D. Turi, P. Missier, C. Goble, D. De Roure, T. Oinn, Taverna work-  
794 flows: Syntax and semantics, in: *e-Science and Grid Computing*, IEEE  
795 *International Conference on*, IEEE, pp. 441–448.
- 796 [35] P. M. Kelly, P. D. Coddington, A. L. Wendelborn, Lambda calculus  
797 as a workflow model, *Concurrency and Computation: Practice and*  
798 *Experience* 21 (2009) 1999–2017.
- 799 [36] J. Brandt, M. Bux, U. Leser, A functional language for large scale  
800 scientific data analysis, in: *BeyondMR, ICDT/EDBT Workshop*.
- 801 [37] J. Montagnat, B. Isnard, T. Glatard, K. Maheshwari, M. B. Fornarino,  
802 A data-driven workflow language for grids based on array programming  
803 principles, in: *Proceedings of the 4th Workshop on Workflows in Sup-*  
804 *port of Large-Scale Science*, ACM, p. 7.
- 805 [38] D. Rogers, I. Harvey, T. Truong Huu, K. Evans, T. Glatard, I. Kallel,  
806 I. Taylor, J. Montagnat, A. Jones, A. Harrison, Bundle and pool archi-  
807 tecture for multi-language, robust, scalable workflow executions, *Journal*  
808 *of Grid Computing (JOGC)* 11 (2013) 457–480.
- 809 [39] S. Cohen-Boulakia, U. Leser, Search, adapt, and reuse: the future of  
810 scientific workflows, *ACM SIGMOD Record* 40 (2011) 6–16.
- 811 [40] Z. Bao, S. Cohen-Boulakia, S. B. Davidson, A. Eyal, S. Khanna, Differ-  
812 encing provenance in scientific workflows, in: *Data Engineering*, 2009.  
813 *ICDE’09. IEEE 25th International Conference on*, IEEE, pp. 808–819.

- 814 [41] P. Groth, M. Luck, L. Moreau, A protocol for recording provenance in  
815 service-oriented grids, in: *Principles of Distributed Systems*, Springer,  
816 2005, pp. 124–139.
- 817 [42] O. Biton, S. Cohen-Boulakia, S. B. Davidson, C. S. Hara, Querying  
818 and managing provenance through user views in scientific workflows, in:  
819 *Data Engineering, 2008. ICDE 2008. IEEE 24th International Confer-*  
820 *ence on*, IEEE, pp. 1072–1081.
- 821 [43] S. Cohen-Boulakia, J. Chen, P. Missier, C. Goble, A. Williams,  
822 C. Froidevaux, Distilling structure in Taverna scientific workflows: a  
823 refactoring approach, *BMC Bioinformatics* 15 (2014) S12.
- 824 [44] A. Chapman, H. Jagadish, Issues in building practical provenance sys-  
825 tems., *IEEE Data Eng. Bull.* 30 (2007) 38–43.
- 826 [45] F. S. Chirigati, D. Shasha, J. Freire, Reprozip: Using provenance to sup-  
827 port computational reproducibility., in: *TaPP, International Workshop*  
828 *on Theory and Practice of Provenance*.