

Laniakea: an open solution to provide Galaxy “on-demand” instances over heterogeneous cloud infrastructures.

Marco Antonio Tangaro¹, Giacinto Donvito², Marica Antonacci², Matteo Chiara³, Pietro Mandreoli³, Graziano Pesole^{1,4}, Federico Zambelli^{1,3}

1. Institute of Biomembranes, Bioenergetics and Molecular Biotechnologies, National Research Council (CNR), Via Amendola 165/A, 70126 Bari, Italy

2. National Institute for Nuclear Physics (INFN), Section of Bari, Via Orabona 4, 70126 Bari, Italy

3. Department of Biosciences, University of Milan, via Celoria 26, 20133 Milano, Italy

4. Department of Biosciences, Biotechnologies and Biopharmaceutics, University of Bari, Via Orabona 4, 70126 Bari, Italy

Abstract

Background

Galaxy is rapidly becoming the de facto standard among workflow managers for bioinformatics. A rich feature set, its overall flexibility, and a thriving community of enthusiastic users are among the main factors contributing to the popularity of Galaxy and Galaxy based applications. One of the main advantages of Galaxy consists in providing access to sophisticated analysis pipelines, e.g., involving numerous steps and large data sets, even to users lacking computer proficiency, while at the same time improving reproducibility and facilitating teamwork and data sharing among researchers. Although several Galaxy public services are currently available, these resources are often overloaded with a large number of jobs and offer little or no customization options to end users. Moreover, there are scenarios where a private Galaxy instance still constitutes a more viable alternative, including, but not limited to, heavy workloads, data privacy concerns or particular needs of customization. In such cases, a cloud-based virtual Galaxy instance can represent a solution that overcomes the typical burdens of managing the local hardware and software infrastructure needed to run and maintain a production-grade Galaxy service.

Results

Here we present Laniakea, a robust and feature-rich software suite which can be deployed on any scientific or commercial Cloud infrastructure in order to provide a “Galaxy on demand” Platform as a Service (PaaS). Laying its foundations on the INDIGO-DataCloud middleware, which has been developed to accommodate the needs of a large number of scientific communities, Laniakea can be deployed and provisioned over multiple architectures by private or public e-infrastructures. The end user interacts with Laniakea through a front-end that allows a general setup of the Galaxy instance, then Laniakea takes charge of the deployment both of the virtual hardware and all the software components. At the end of the process the user has access to a private, production-grade, yet fully customizable, Galaxy virtual instance. Laniakea’s supports the deployment of plain or cluster backed Galaxy instances, shared reference data volumes, encrypted data volumes and rapid development of novel Galaxy flavours, that is Galaxy configurations tailored for specific tasks. As a proof of concept, we provide a demo Laniakea instance hosted at an ELIXIR-IT Cloud facility.

1 Conclusions

2 The migration of scientific computational services towards virtualization and e-infrastructures is one of the most visible
3 trends of our times. Laniakea provides Cloud administrators with a ready-to-use software suite that enables them to
4 offer Galaxy, a popular workflow manager for bioinformatics, as an on-demand PaaS to their users. We believe that
5 Laniakea can concur in making the many advantages of using Galaxy more accessible to a broader user base by removing
6 most of the burdens involved in running a private instance. Finally, Laniakea's design is sufficiently general and modular
7 that could be easily adapted to support different services and platforms beyond Galaxy.

8 Background

9 The recent improvements in our capacity to gather vast amounts of complex, multi-layered and interconnected
10 biomolecular data demand a parallel development and enhancement of the computational tools that we employ to
11 analyse and handle this wealth of information. On the other hand, the rapid proliferation of those tools can render the
12 execution of complex bioinformatics workflows cumbersome due to, among other things, incompatible data formats,
13 long and convoluted command lines, versioning, and the need to handle and store a multitude of intermediate files. In
14 turn, that not only makes harnessing the information contained in the biomolecular data unnecessarily onerous even
15 for expert bioinformaticians but represents also a significant obstacle to reproducibility of the analyses [1] as well as an
16 intimidating barrier for biologists aspiring to explore their own data in autonomy [2], students [3], and health-care
17 providers adopting *clinical bioinformatics* approaches within their medical protocols [4]. For these reasons, in the past
18 few years, a considerable effort has been put in the development of workflow manager platforms for bioinformatics
19 (see [5] for a review). They usually provide integrated interfaces that not only provide a more user-friendly work
20 environment but improve reproducibility, facilitate data sharing and enable collaborative data processing.

21 Galaxy

22 The Galaxy platform is one of the most successful examples of such workflow management systems. Indeed, by
23 providing a consistent, user-friendly, flexible, effective and customizable gateway to a vast array of bioinformatics
24 software and analysis workflows, Galaxy has attracted a vast and thriving community of users [6]. The software consists
25 of an open source server-side application accessible through a simple web interface that serves as a gateway to a wealth
26 of tools for datasets handling and analysis, workflow design, visualization and sharing of results. Although from the
27 user's perspective using a Galaxy service to run bioinformatics analyses is pretty straightforward, the deployment of a
28 production-grade Galaxy instance requires the configuration and maintenance of an extensive collection of helper
29 components (e.g. database management system, web server, load balancer, etc...) and an even more extensive
30 collection of bioinformatic tools and reference data. These issues, coupled with the need of an adequate or dedicated

1 IT infrastructure, required to properly support the Galaxy service for any non-trivial amounts of users or workloads, has
2 usually restricted the role of Galaxy service providers to institutions or groups with suitable IT facilities and the
3 appropriate technical know-how. At the time of writing, more than 125 Public Galaxy instances are available
4 (<https://galaxyproject.org/use/>), serving a vast community of users [6]. Despite being useful and popular, the public
5 nature of those services implies some hardly addressable shortcomings like, for example, limited quotas for computing
6 and storage resources, lack of customization options and potential concerns for data security and privacy, a worry
7 particularly noticeable when processing sensitive data. These considerations can, in turn, limit or outright interdict the
8 usage of Galaxy public instances for some specific applications or categories of users, e.g., analyses requiring big or huge
9 computing workloads or precision medicine researchers and operators.

10 Cloud solutions to Galaxy provision

11 The cloud computing model [7] is rapidly gaining popularity within the life sciences [8–11] and the biomedical [4,12,13]
12 communities. Among other advantages, it offers a set of solutions and features that can overcome or mitigate the
13 drawbacks of Galaxy public instances described above. At the time being, several efforts have already been put forward
14 in this regard. Globus Genomics [14] provides a Galaxy-based bioinformatics workflow platform, built on Amazon cloud
15 services, for large-scale next-generation sequencing analyses. CloudMan [15,16] allows individual researchers to deploy
16 Galaxy instances relying on arbitrarily sized compute cluster on the Amazon cloud infrastructure (it can also support
17 OpenStack and OpenNebula through custom deployments). The Genomic Virtual Laboratory (GVL) [17], offers Galaxy
18 through a middleware available on Nectar, the Australian cloud infrastructure for research, and also on the Amazon
19 cloud. PhenoMeNal [18] is a recent effort to develop a Cloud Research Environment that employs Galaxy to provide
20 domain-specific tools for metabolomics. Finally, Krieger and colleagues describe a possible configuration stack to deploy
21 Galaxy on an OpenStack based IaaS (Infrastructure as a Service) [19]. The appeal of Galaxy cloud solutions is also made
22 evident by the 2016 Galaxy update [20] reporting that over 2,400 Galaxy servers were launched on the Amazon cloud
23 in 2015 alone, pointing to strong demand for ready-to-use but private virtual Galaxy instances. The Amazon Galaxy
24 service [21] is, however, a commercial solution that can discourage researchers and research or healthcare facilities
25 from adopting it, due to funding or budget issues, ethical concerns or legal requirements (e.g., EU General Data
26 Protection Regulation). At the same time, other solutions like access to GVL through the Nectar cloud are not usually
27 within reach of European researchers and are not completely interoperable with other resources available to them
28 through European e-infrastructures. In this article, we introduce Laniakea, a software framework for the provision of
29 *on-demand* Galaxy instances. Laniakea is devised to run on existing scientific e-infrastructures by leveraging the open

4

1 and modular middleware architecture developed within the INDIGO-DataCloud H2020 project [22,23]
2 (<https://www.indigo-datacloud.eu/>) and is at the same time simple to deploy and maintain by infrastructure managers
3 and convenient for end-users.

4 Disambiguation of the terms “user”, “Galaxy user” and “IaaS administrator”

5
6 In this work, we use the words “user” or “end-user” to indicate anyone owning an account to access an instance of
7 Laniakea. An authenticated Laniakea user can deploy Galaxy instances and has administrator rights both over the
8 deployed instances and the virtual hardware hosting them. A single Galaxy instance can be accessed and used by many
9 people, or “Galaxy users”, and there is no need for them to be also Laniakea users. Indeed, any Galaxy instance created
10 by Laniakea receives a public IP and manages accounts in the same way of any other Galaxy instance, also supporting
11 anonymous login if the administrator so decides. Finally, we use the terms “IaaS administrator” or “IaaS manager” for
12 the person or entity running a cloud infrastructure that offers Laniakea as one of their services.

13 Methods

14

15 INDIGO-DataCloud middleware components

16 Laniakea is based on the ElectricIndigo software release of the INDIGO-DataCloud H2020 project (INDIGO from now on),
17 that officially ended in September 2017 and passed on its legacy to the EOSC-Hub, (<https://www.eosc-hub.eu/>), DEEP-
18 HybridDataCloud (<https://deep-hybrid-datacloud.eu/>) and XDC (<http://www.extreme-datacloud.eu/>) projects. INDIGO
19 aimed at making cloud e-infrastructures more accessible by scientific communities. In fact, the heterogeneity of the
20 currently available scientific cloud infrastructures often entails portability issues between the different technologies
21 employed that can result in the lack of one or more key features, e.g., automatic elasticity or multi-clouds deployments.
22 In turn, these resources usually result in less efficient solutions than commercial counterparts relying on more
23 homogeneous environments. INDIGO’s software catalogue [22] tackles e-infrastructures heterogeneity through
24 comprehensive support to open standards and solutions, e.g., adopting both OpenStack (<https://www.openstack.org>)
25 and OpenNebula (<https://opennebula.org>) as Cloud Management Platforms (CMPs) and allowing to transparently
26 deploy applications using either virtual machines (VMs) or Docker containers. This result has been achieved by
27 integrating the needs of a wide range of use cases and by taking advantage of publicly available open-source cloud
28 components, adapting or enhancing them if needed to obtain the desired functionalities and embarking in the
29 development of new software packages when strictly necessary. Fig. 1 provides a complete overview of the INDIGO

1 PaaS software components and the installation sequence needed to deploy them correctly. Table S1 (in
2 Supplementary2.docx) collects the URLs of all the INDIGO components required by Laniakea.

3 **Orchestrator and Infrastructure Manager**

4 The INDIGO PaaS Orchestrator and the Infrastructure Manager (IM) components (<https://github.com/indigo->
5 [dc/orchestrator](https://github.com/indigo-dc/orchestrator)), which are both based on the OASIS Topology and Specification for Cloud Applications (TOSCA) [24,25]
6 open standard language, are in charge of setting up the virtual infrastructure environment and deploy the software
7 framework. Cloud applications and services, together with their requirements and dependencies and regardless of the
8 underlying platform, are described using YAML syntax [26]. This approach enables the deployment of complex
9 applications from small reusable building blocks, the “node types”, and the topology of their relationships. The INDIGO
10 PaaS Orchestrator coordinates deployment of the application according to the specifications described in the TOSCA
11 template: it selects the most suitable cloud infrastructure site among those available and delegates the deployment and
12 configuration of the virtual infrastructure on the target site to the Infrastructure Manager (IM)
13 (<https://github.com/indigo-dc/im>). IM supports deployment on OpenStack, OpenNebula and commercial cloud
14 providers alike, and serves as an abstraction layer for the definition and provision of the required resources. Finally, the
15 software layout of the infrastructure’s nodes is described using Ansible roles (<http://docs.ansible.com>). They instruct
16 the automation engine on how to install and configure the end-user applications or services, like Galaxy, on bare OS
17 images. All in all, this architecture allows for the deployment of a wide range of cloud agnostic applications and services.

18 **Authorization and Authentication Infrastructure**

19 The INDIGO Identity and Access Management (IAM) is an Authentication and Authorisation Infrastructure (AAI) service
20 that manages users identities, attributes (e.g., affiliation and groups membership) and authorization policies for the
21 access to INDIGO based resources. IAM supports several token-based protocols, i.e. X.509 [27], OpenID Connect [28]
22 and SAML [29], thus allowing the federation of different infrastructures. This system allows users to connect to
23 federated PaaS and manage heterogeneous and distributed resources through a single AAI service and with just one
24 user account, easing the overhead experienced by end-users in case of multiple INDIGO services running on the same
25 or federate infrastructures.

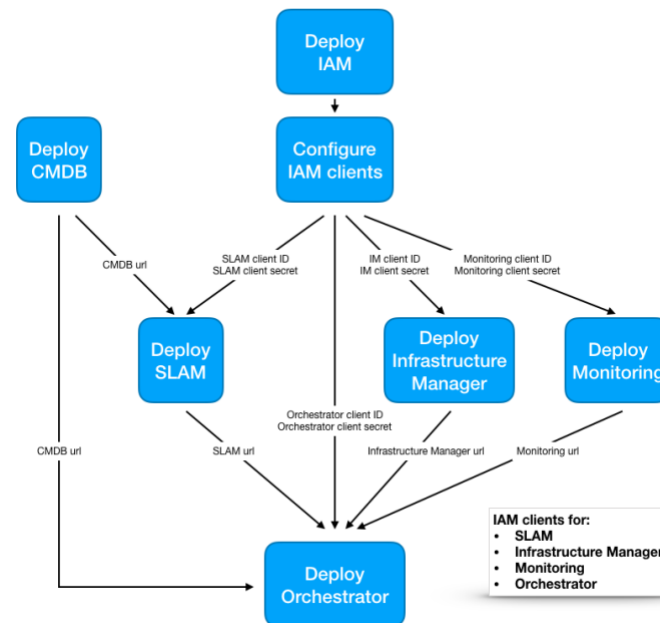
26

27

28

6

1



2

3 Fig. 1. INDIGO-DataCloud PaaS software components and deployment procedure. The first step consists in installing and
4 configuring IAM; this allows registered clients to request and receive information about authenticated end-users. Each
5 INDIGO service must authenticate to a dedicated IAM client using a public identifier, the "client id", and a private password,
6 the "client secret", provided during their configuration. In this way, IAM can verify that the connecting service is trusted
7 and grants it the rights to interact with the other PaaS and IaaS services. The CMDB service installation follows as the second
8 step; it provides detailed information on the available cloud sites and the virtual environments that they support. The
9 following steps involve the deployment of SLAM, the Service Level Agreement Manager, the IM, and the Monitoring stack.
10 Finally, the INDIGO PaaS-Orchestrator is deployed.

11

12 Other INDIGO components

13 - The **CLUES** service is an elasticity manager for HPC clusters that enables dynamic cluster resources scaling, deploying
14 and powering-on new working nodes depending on the workload of the cluster and powering-off and deleting them
15 when no longer required. Once new jobs are submitted to the cluster Resource Manager (RM) queue (e.g., SLURM [30]
16 or Torque [31]), CLUES contacts the Orchestrator and starts the provisioning of available temporary resources, thus
17 improving the overall efficiency of the infrastructure.

18 - **FutureGateway** is built on top of the Liferay open source portal framework (<https://www.liferay.com>) and provides
19 GUI based portlets that interact with the PaaS layer, allowing the authentication through the integrated INDIGO IAM
20 portlets and the customization of relevant parameters of the TOSCA templates by the user before dispatching them to
21 the Orchestrator for deployment.

22 Encryption at block device level

1 The encryption layer is based on LUKS (Linux Unified Key Setup) [32], the current standard for encryption on Linux
2 platforms. It provides robustness against low-entropy passphrase attacks using salting and iterated PBKDF2 passphrase
3 hashing. LUKS supports secure management for multiple user passwords, allowing to add, change and revoke passwords
4 without re-encryption of the whole device. Supplementary1.docx contains a detailed description of the Laniakea
5 encryption framework.

6 CernVM File System

7 CernVM File System (CVMFS) [33] is used to support common reference data repositories. CVMFS is a read-only POSIX
8 file system originally developed to facilitate the distribution of High Energy Physics analysis software through HTTP
9 protocol. Data files are hosted on any server and can be mounted concurrently on multiple compute nodes through a
10 Linux filesystem module-based client (FUSE) that loads and caches only the specific portions of the files that are
11 required, on-demand. This solution is suitable for quasi-static files that have to be shared across several geographically
12 distributed clusters, and it supports local caches to speed up read operations on the hosted data too.

13 Results

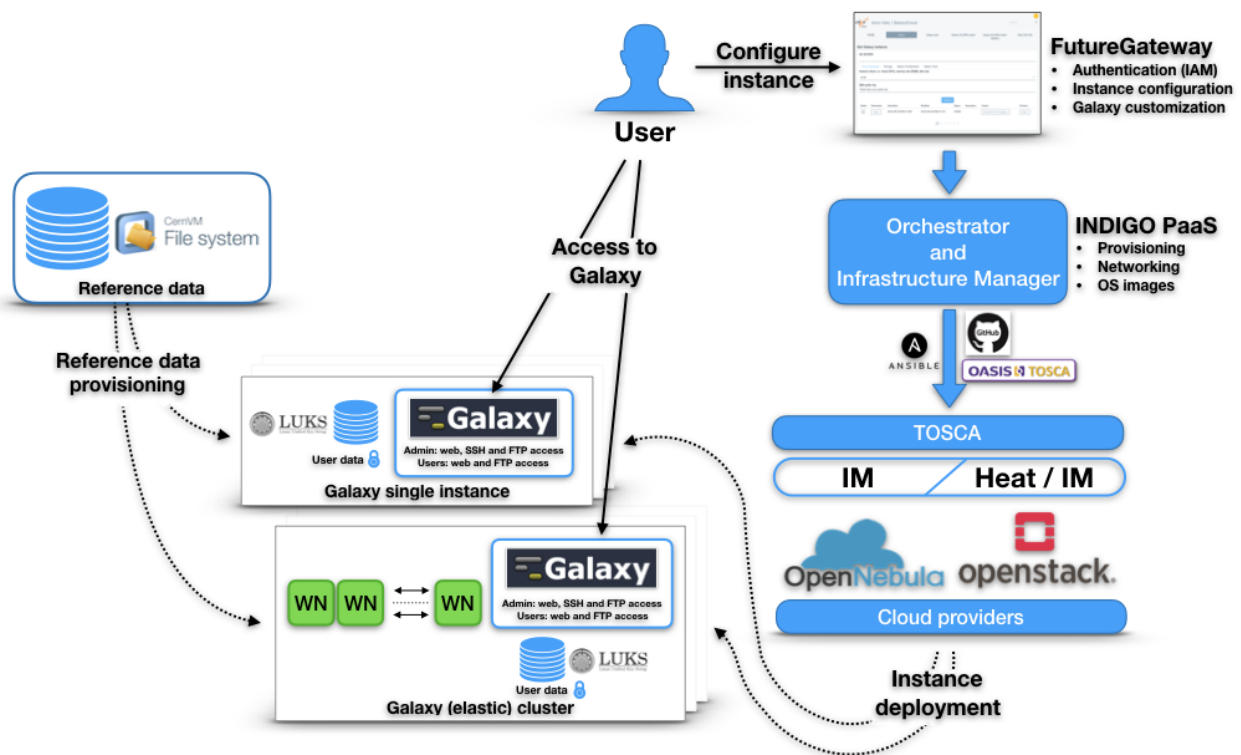
14 A wide range of Galaxy users and instance types do exist. For example, users may stretch from plain account owners to
15 instance administrators and from developers to training courses teachers. At the same time, Galaxy instances can have
16 different footprints in terms of computational resources, depending for example on typical usage or number of
17 concurrent users; can be short-lived to satisfy temporary needs (e.g., a training course) or long-lived to provide a
18 persistent work environment; public or private; needing advanced data security features in order to work on sensitive
19 human data or not, etc.... As a consequence, Galaxy PaaS providers should be able to cover the requirements of the
20 largest possible number of user/instance type combinations in order to maximize the convenience of their service.

21 Laniakea's architecture (Fig 2) has been developed to address all these heterogeneous needs swiftly, empowering end-
22 users with a wide array of customization options that can deliver from out-of-the-box, stable, production-grade Galaxy
23 instances already configured with collections of bioinformatics tools, reference data and cluster support, to blank-sheet
24 Galaxy instances, ready to be completely personalized. Another key feature of Laniakea consists in the integration, for
25 the first time at best of our knowledge, of a built-in technology to encrypt storage volumes on a Galaxy on-demand
26 platform. This functionality can be used to provide robust data protection through state-of-the-art encryption protocols:
27 the secure layer insulates stored sensitive data from unauthorized access both by malicious attackers or trusted users
28 of the same cloud infrastructure and notably including the administrator(s) of the cloud and hardware layers
29 themselves. Finally, Laniakea supports reference data sharing via CVMFS and a high degree of instance scalability

8

1 through an array of deployment configurations ranging from single node Galaxy instances to SLURM managed Galaxy
 2 clusters, including also cluster elasticity based on CLUES.































3 Table 1 provides a detailed features comparison table of Laniakea, PhenoMeNal and GVL. Table S2 (in
 4 Supplementary2.docx) collects the URLs of the software components of Laniakea. A video demo showcasing a typical
 5 user interaction routine is available at <https://goo.gl/xnWNQd>. In the following paragraphs, we describe in more detail
 6 the components and features of the software suite.






7

8 Fig. 2. Architecture and deployment flow. The FutureGateway portal provides users with the front-end to configure and
 9 manage Galaxy instances, the resulting TOSCA templates are sent to the PaaS layer that employs INDIGO services to deploy
 10 instances over the IaaS, that in turn provides virtual hardware, storage, and networking. Finally, the Galaxy instance (single
 11 or cluster backed) is configured with the requested flavour and attached to a plain or LUKS encrypted storage volume
 12 (depending on user's choice) and to the CernVM-FS shared volume with reference data. At the end of the process, the
 13 public IP address of the new Galaxy instance is provided to the user.
 14
 15

| | Laniakea / INDIGO | PhenoMeNal | Genomic Virtual Lab / CloudMan |
|-------------------------|--------------------|---------------------------|--------------------------------|
| Virtual cluster support | ✓ SLURM | ✓ Kubernetes + KubeNow | ✓ SLURM + CloudMan |
| Automatic elasticity | ✓ SLURM + CLUES | ✗ | ✓ SLURM + CloudMan |

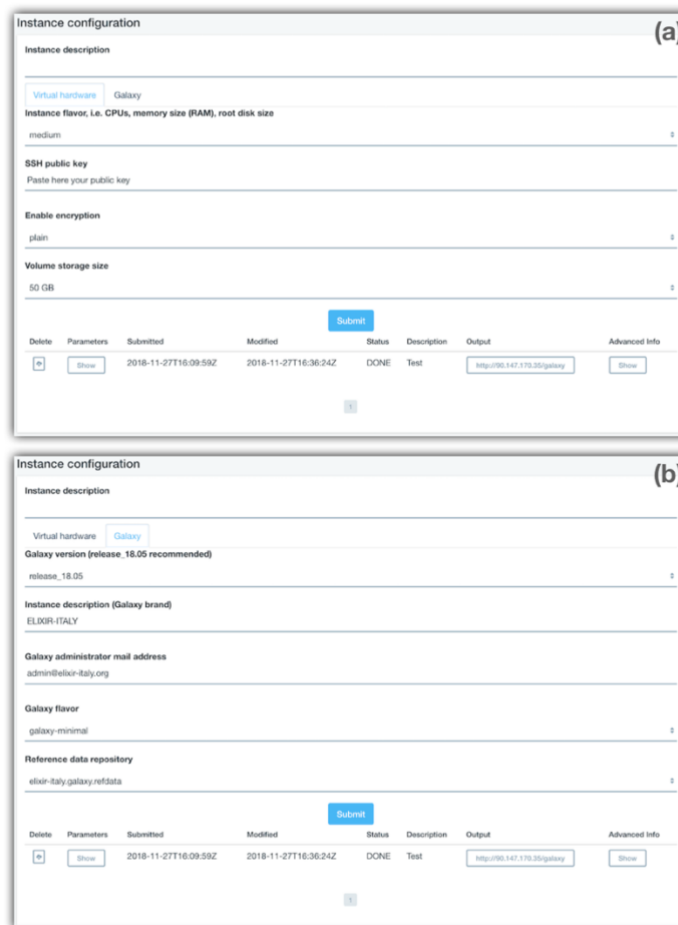
| | | | |
|--|--|---|---|
| User data protection |  Support to on-demand, on-the-fly user data encryption |  |  |
| Reference data sharing |  CernVM File System |  |  CernVM File System |
| Production-grade Galaxy instance (Nginx, PostgreSQL, proFTPd, uWSGI) |  Based on CENTOS 7 (supports also Ubuntu 16.04) |  Based on Ubuntu 14.04 |  Based Ubuntu 14.04/Ubuntu 16.04 |
| Administrator access to Galaxy |  Http, ssh, ftp. |  Http, ftp. Ssh access is allowed to the kubernetes VMs. Then, the administrator has to enter the Galaxy Docker container. At least a basic knowledge of the platform is needed. |  Http, ssh, ftp. |
| Galaxy instance build process |  Image Snapshot / Bare VMs / Docker It is possible to deploy pre-configured VMs or start a fresh Galaxy installation over a bare CentOS image. |  Only Docker containers available. Docker with tools need to be centrally updated. The process to integrate or update tools requires a good knowledge of kubernetes and Docker engines. |  Pre-configured VMs or Docker containers available. Image and Docker with tools need to be centrally updated or installed after Galaxy deployment. |
| Galaxy flavours |  Galaxy Minimal Galaxy CoVaCS Galaxy GDC Somatic Variant Galaxy Epigen Galaxy RNA Workbench <i>laaS administrator can easily add New flavours.</i> |  Only one flavour, specific for metabolomics analysis. |  One general purpose flavour + BioDocketlet Exome-Seq Chemical Toolbox RNA Structural Analysis RNA Workbench |
| Installation of new Galaxy tools and usage as tools development environment |  New tools can be installed through the ToolShed. Tool can be normally developed and installed in Galaxy. |  All new tools have to be deployed on the kubernetes cluster using Docker containers. The process to integrate or update tools requires a good knowledge of kubernetes and Docker engines. |  New tools can be installed through ToolShed. Tool can be normally developed and installed in Galaxy. |
| Virtual technology |  Virtual machines / Docker (in progress) |  Docker |  Virtual machines / Docker |
| Authentication and Authorization | INDIGO IAM (include support for Google Accounts). ELIXIR-AAI integration is ongoing. | ELIXIR-AAI (include support for Google accounts) | Google / GitHub / LinkedIn / Facebook accounts. |
| Access to virtual resources |  The process is completely transparent to the end user. He is automatically redirected to the right laaS after login into the portal. Only the login credentials are required. |  After the login into the portal the user can select among a set of supported cloud providers and has to provide the corresponding credentials to access. |  After the login into the portal the user can select among a set of supported cloud providers and has to provide the corresponding credentials to access. Each available GVL application could have a different set of available cloud providers. |
| Supported Cloud Provider Managers | INDIGO is open source and currently supports OpenStack, OpenNebula, Amazon AWS, MS Azure. Working on more (e.g. Google Cloud). | Terraform and kubernetes are exploited to deploy Galaxy. Openstack, Amazon AWS, Google Cloud are currently supported. | Custom open source solution, supporting OpenStack and Amazon AWS. |
| e-Infrastructures federation |  It is possible to federate multiple laaS (both private and public), through IAM, Elixir AAI, etc. Different policies can be put in place, allowing (groups of) user to deploy in a certain e-infrastructure. |  The user has to select the e-infrastructure where deploy Galaxy, thus providing the corresponding credentials. |  The user has to select the e-infrastructure where deploy Galaxy, thus providing the corresponding credentials. |

| | | | |
|---|---|--|--|
| Deployment of the full software stack to local IaaS |  The INDIGO software stack, including Laniakea, can be installed on any compatible e-infrastructure. The procedure is documented. Moreover, Ansible roles and Docker containers to automatically deploy each PaaS component are provided. |  The phenomenal portal is centrally linked to the supported e-infrastructures. The documentation to locally install the portal is available. |  GVL can be installed on any compatible e-infrastructure or the portal can be linked to the supported e-infrastructure. The procedure is documented. |
|---|---|--|--|

1 Table 1. Comparative features table of Laniakea, PhenoMeNal and GVL. Virtual environment (light yellow background),
 2 Galaxy (light green background) and IaaS specific (light blues background) features are compared. Data for PhenoMeNal
 3 and GVL have been gathered from available documentation using our best effort. Only features related to the Galaxy on-
 4 demand platform have been considered, leaving aside GVL's and Phenomenal's support for other software (e.g., RStudio,
 5 Jupyter).
 6

7 The web front-end

8 The web front-end of Laniakea (Fig. 3), is built upon the INDIGO FutureGateway component and becomes accessible
 9 after authentication through INDIGO IAM. The front-end represents the entry point to the PaaS from a user's point of
 10 view. It provides the configuration interface for initialization and deployment of all the virtual hardware/Galaxy instance
 11 type combinations available through separate configuration panels for Galaxy single VM or cluster deployments.



12

13 Fig 3. Web front-end. The web interface is organised using different tabs for each configuration task. The "Virtual hardware"
 14 tab (a) allows the selection of the virtual hardware, e.g., the number of CPUs and quantity of memory, storage size and to

1 submit the user SSH public key that will grant access to the VM at the end of the deployment. For cluster deployments, the
2 same tab also allows configuring the number of worker nodes required and their hardware configuration. The “Galaxy” tab
3 (b) provides the software configuration options: Galaxy version, instance description, the mail address of the administrator,
4 the flavour and Reference data repository from the available ones. The list of active instances and their status is displayed
5 at the bottom of the page.
6

7 The configuration front-end is organized in two panels to guide users through the design of the new Galaxy instance:

- 8 • The Virtual Hardware configuration tab exposes the array of available hardware setups, differing for the
9 number of virtual CPUs and quantity of RAM. Users are required to provide a valid SSH public key at this point.
10 The key will be used to grant access to the Virtual Machine, once deployed. The storage volume size and type,
11 i.e., plain or encrypted can also be selected at this stage. The Virtual Cluster deployment panel allows for the
12 setup of both the front-end and worker nodes.
- 13 • The Galaxy Configuration tab allows the selection of the Galaxy release version among the ones currently
14 supported, the reference dataset, the e-mail of the instance administrator and finally the *Galaxy flavour*. The
15 Galaxy flavour is the pre-configured set of bioinformatics tools that will be installed and thus ready to use
16 straight away after deployment (more on this in the “Galaxy flavours” section).

17 The settings defined through the interface are in turn submitted to the INDIGO Orchestrator that manages the
18 deployment and configuration processes. When the procedure terminates, a public IP address becomes available which
19 provides access to the freshly minted Galaxy server. From that moment onwards, the user has full administrator
20 privileges over the instance and can access the underlying Virtual Machine through SSH using the private key
21 corresponding to the public one that was provided. Public key authentication has been preferred over using passwords
22 since this approach avoids most of the drawbacks of using and safeguarding passwords both for users and the IaaS
23 administrator.

24 The Galaxy environment

25 The deployment of a Galaxy multi-user production environment is a complex task that requires several auxiliary
26 software components that in turn must be properly installed and configured. Laniakea automates this lengthy and error-
27 prone procedure by spawning on-the-fly virtual environments. Their standard configuration (Table 2) is as follows:
28 CentOS7 as the operating system (OS), PostgreSQL as the database engine, Nginx, uWSGI, and ProFTPD as the web,
29 application and FTP servers respectively. Apart from the OS, for which there are no official recommendations, this
30 configuration is rooted in the guidelines for production environments issued by the Galaxy Project itself
31 (<https://docs.galaxyproject.org/en/latest/admin/production.html>). CentOS 7 has been selected as the reference OS for

1 Laniakea virtual environments due to its well-known adherence to standards and the long span of the foreseen official
2 support for this release (updates until June 30, 2024); however, Ubuntu 16 is also supported. Laniakea fully supports
3 virtual machines while docker containers are currently in the final stage of testing (containers available at
4 <https://hub.docker.com/u/laniakeacloud/>) and will be soon available, leaving the choice of the virtual environment
5 solution to the preferences of the IaaS administrator. The on-the-fly deployment of the virtual environment offers a
6 higher degree of agnosticism over pre-configured virtual machines and at the same time ensures that each new instance
7 of Galaxy makes use of the latest available release of software components. However, this approach suffers from two
8 possible limitations: first, the installation procedure takes time, about five hours on our test IaaS and second, there is a
9 chance of some software components (mainly bioinformatics tools) failing to install due to any temporary unavailability
10 of the online repositories needed for their installation. To overcome these limitations, Laniakea also provides a backup
11 procedure to instantiate pre-configured images of Galaxy instances that the IaaS administrator can easily create from
12 the on-the-fly ones previously described and provide to the user through the “Galaxy Express” configuration panel of
13 Laniakea. In this latter case, the procedure will automatically adapt the configuration of the image to make it compatible
14 with the virtual hardware and any other parameter (e.g., the instance administrator credentials) selected by the user.
15 Apart from being more resistant to third parties repositories unavailability, this procedure allows faster deployments,
16 i.e., less than an hour on our test IaaS, at the cost of a lesser degree of IaaS agnosticism.

| Software component | Version installed by Laniakea |
|--------------------|-----------------------------------|
| Operative System | CentOS 7 (supported Ubuntu 16.04) |
| Galaxy | release_17.05 and release_18.05 |
| PostgreSQL | 9.6 |
| NGINX | 1.12.2 |
| uWSGI | 2.0.17.1 |
| PROFTPD | 1.3.5e |

17 Table 2. Laniakea’s Galaxy production environment components summary.

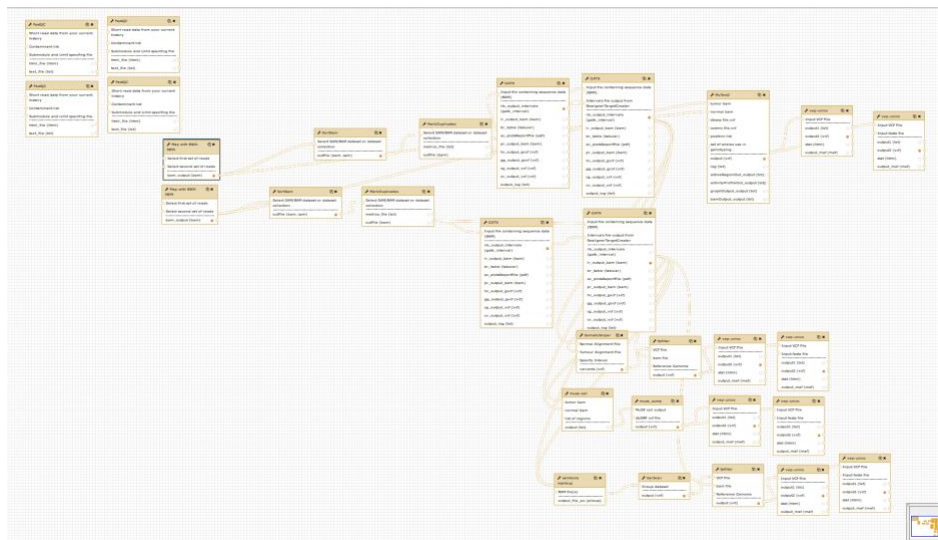
18

19 Galaxy flavours

20 The Galaxy Tool Shed (<https://toolshed.g2.bx.psu.edu/>) is the official Galaxy tools repository and provides a convenient
21 system for the administration and customization of a Galaxy web-server with tools and workflows tailored for its
22 expected typical use. However, despite the significant improvements brought by the recent adoption of Conda
23 (<https://conda.io>) for the management of packages and dependencies, the concurrent installation and configuration of
24 many tools and their respective dependencies and reference datasets can still be burdensome and sometimes it requires
25 a low-level understanding of the Galaxy environment and of the tools to be installed. For this reason, Laniakea provides
26 a handy set of domain-specific *Galaxy flavours*, that is Galaxy instance enriched with cured collections of tools already

1 installed, configured, tested, organised in workflows and ready to be used out-of-the-box. To access a target Galaxy
2 instance and install the required bioinformatics tools Laniakea employs the official Galaxy Project library Ephemeric
3 (<https://github.com/galaxyproject/ephemeris>), and YAML recipes are used to list the required packages. A helper
4 Ansible role is used for the fine-tuning of packages configuration and to solve dependencies that, for any reason, are
5 missing or malfunctioning upon installation with Conda. All in all, this approach allows IaaS administrators to easily
6 prepare additional flavours by just creating new YAML recipes and eventually tweak the resulting configuration using
7 the helper Laniakea Ansible role.

8 Currently, Laniakea provides five proof-of-concept flavours that are named “galaxy-minimal”, “galaxy-epigen”, “rna-
9 workbench”, “GDC_Somatic_Variant” and “CoVaCS”. The first flavour is a bare instance providing just the default tools
10 embedded in any Galaxy installation and provides a clear foundation for users who need a blank sheet ready to be
11 customized to their own needs or for developers of new tools and flavours. The “galaxy-epigen” flavour is based upon
12 the layout of the Italian Epigen Project’s (<http://www.epigen.it/>) Galaxy server (<http://www.beaconlab.it/epigalaxy>)
13 and provides a selection of tools for the analysis of ChIP-Seq and RNA-Seq data. The “rna-workbench” flavour is based
14 on [34] and includes more than 50 tools dedicated to RNA-centric analyses including, i.e., alignment, annotation,
15 secondary structure profiling, target prediction, etc... The “GDC_Somatic_Variant” flavour is a porting of the Genomic
16 Data Commons (GDC) pipeline for the identification of somatic variants on whole exome/genome sequencing data
17 (<https://gdc.cancer.gov/node/246>), the resulting Galaxy workflow is shown in Fig. 4. In this case, Galaxy wrapper
18 s for several tools, which were not previously available, have been developed from scratch. Finally, the CoVaCS flavour
19 implements the homonymous workflow (Fig. S1 in Supplementary2.docx), described in [35], and set of tools for
20 genotyping and variant annotation of whole genome/exome and target-gene sequencing data. All in all, we believe that
21 these examples provide an adequate proof-of-concept of the relative easiness of generating Galaxy flavours for
22 Laniakea.



1

2 Fig. 4. The GDC Somatic Variant analysis pipeline implemented as a Galaxy workflow for the corresponding flavour. The
3 workflow design interface of Galaxy is a powerful instrument to elaborate complex workflows that link together the output
4 and the input of different tools in an intuitive fashion.

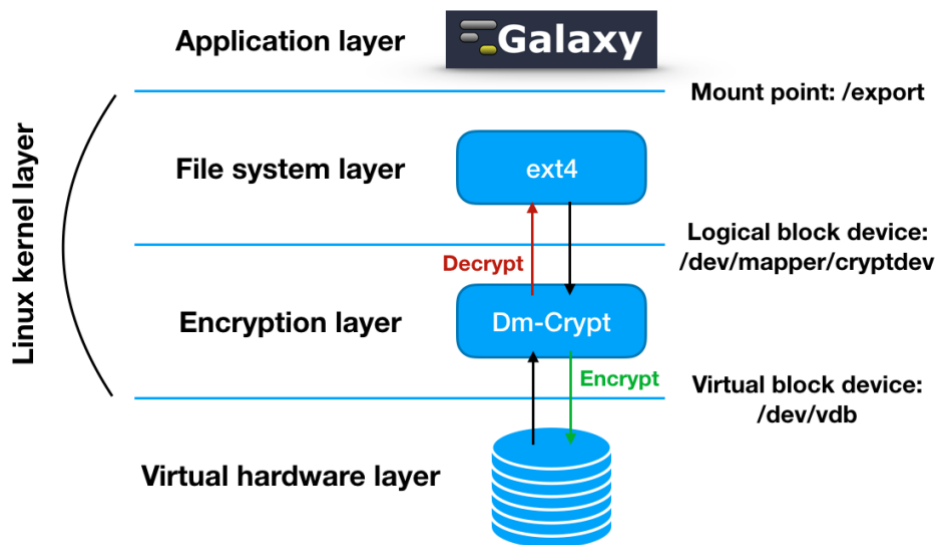
5

6 IaaS wide Reference data availability

7 To avoid useless replication of reference data and facilitate reproducibility of analyses, Laniakea employs a read-only
8 CVMFS volume shared among all the instances of the same IaaS. The IaaS manager can choose to provide its own CVMFS
9 customized reference dataset or to mirror or link a remote one, e.g., the one made available by the Galaxy Project. In
10 the latter case, the required reference data need to be transferred from a remote service on-the-fly. To streamline the
11 creation of new CVMFS repositories by the IaaS manager, we have made available a suitable Ansible role and
12 documented the procedure. As a proof of concept, we provide three different repositories. The first is a manually
13 curated reference dataset maintained by us, which contains the latest releases of the human, mouse, yeast, fruit fly and
14 *A. thaliana* genomes and their corresponding indexes for Bowtie [36], Bowtie2 [37] and BWA [38,39]. The second is a
15 repository tailored for variant calling in human, to be used with the CoVaCS and GDC Galaxy flavours, and provides
16 access to a collection of publicly available variants, provided as VCF format files, derived from a selection of large-scale
17 human genome resequencing projects [40–42], along with curated human genome assemblies and indexes obtained
18 from the GATK bundle repository [43]. Finally, we provide a mirror of the Galaxy project “by hand” reference repository.
19 In principle, this approach could be extended to optimize the use of storage resources and performances in case of a
20 single user or group using several geographically distributed IaaS resources. Galaxy administrators needing additional
21 reference data not included in these repositories will still be able to add them using the storage resources assigned to
22 their instance.

1 **Data protection and isolation**

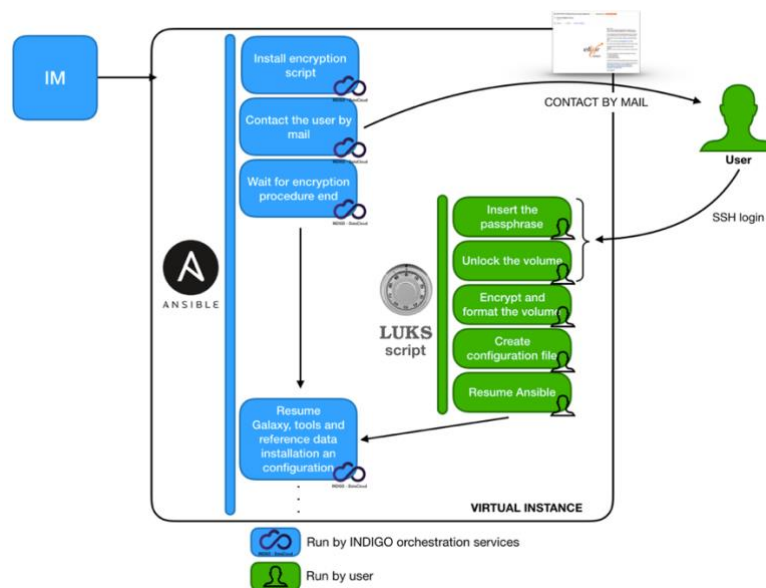
2 Unless proper countermeasures are in place, data stored on a VM could be exposed to anyone with legitimate or
3 illegitimate access to the underlying IaaS and physical hardware [44]. This issue poses technical, ethical and legal
4 concerns, and in particular when sensible human genetic data are involved. These problems are exceptionally relevant
5 for health operators and researchers involved in clinical bioinformatics or similar scenarios. We tackle this issue by
6 providing to Laniakea users a security layer that seamlessly encrypts the storage volume using file system level
7 encryption based on the *dm-crypt* Linux kernel's subsystem coupled with LUKS encryption strategy based on aes-xts-
8 plain64 cipher (see Table S3 in Supplementary2.docx for the detailed configuration parameters). The storage volume is
9 encrypted using a *key stretching* approach: a randomly generated master key is encrypted using the user passphrase
10 through PBKDF2 key derivation. This procedure makes both brute force and *rainbow tables* [45] based attacks more
11 computationally expensive, and at the same time allows for multiple passphrases and passphrase change or revocation
12 without re-encryption. Finally, the LUKS *anti-forensic splitter* feature protects data against recovery after volume
13 deletion. The resulting instance layout consists of Galaxy running on top of a standard file system but transparently
14 using the encrypted volume for storing data as long as it is unlocked and mounted (Figure 5).



15

16 Fig. 5. The relationship between Galaxy, the file system, and dm-crypt. Data are encrypted and decrypted on-the-fly when
17 writing and reading through dm-crypt. The underlying disk encryption layer is completely transparent to Galaxy that
18 employs a specific mount point in order to store and retrieve files from the volume.

19



1

2 Fig. 6. Storage encryption workflow. The user receives an e-mail with instructions to connect through SSH to the virtual
3 Galaxy instance being created. The user is then requested to enter a passphrase that will be used to encrypt the volume
4 and requested each time the encrypted volume will need to be unlocked, e.g. during a reboot of the VM hosting the Galaxy
5 instance.

6

7 The encryption procedure is coordinated by the IM which installs the encryption package and sends to the user an e-
8 mail containing the information needed to log in to the newly created VM, together with a brief description of the
9 encryption procedure and a detailed step by step how-to (Figure 6). Users are thus required by this procedure to insert
10 the passphrase for file system encryption manually through an SSH connection to the Galaxy instance being deployed;
11 a similar system will allow mounting the encrypted volume each time the encrypted instance is re-booted. This two-
12 steps solution has been devised to separate the orchestration of services from the encryption procedure, ensuring that
13 the encryption passphrase is never being exchanged as plain text during the deployment procedure and avoiding any
14 interaction with the IaaS administrator(s). As a result, the user of Laniakea is the only one holding the passphrase and
15 thus able to unlock the encrypted volume.

16 To validate our data encryption strategy, we simulated two different attack scenarios. In the first scenario, the attacker
17 obtains unauthorized access to the unmounted encrypted volume, while the second simulates the improper use of
18 administrator IaaS privileges when the LUKS volume is already unlocked and in use by a running Galaxy instance.

19 For the first scenario we compared two identical volumes, one encrypted and the other not, both attached to the same
20 Galaxy instance, with the same set of permissions and each containing a copy of the same plain text file. Once detached,
21 we created a binary image file of each volume and tried to access the data structure through *hex dump*. We were able

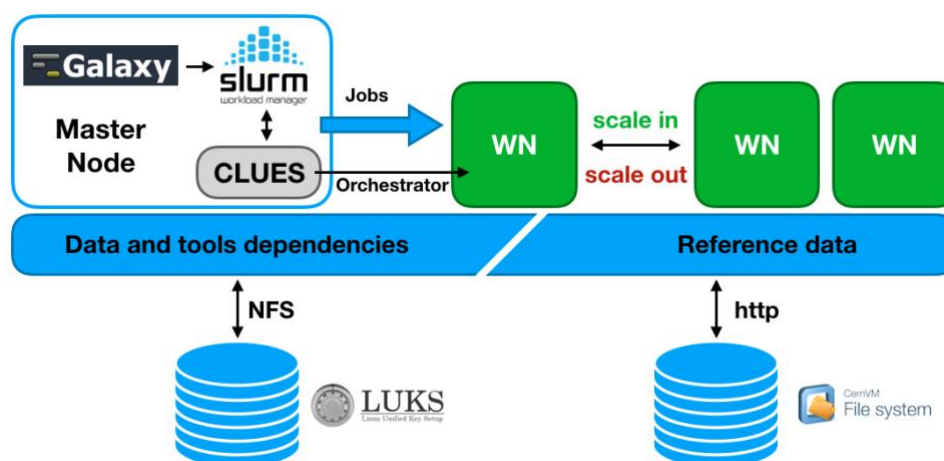
1 to quickly retrieve the original content of the text file from the non-encrypted volume while the hex dump of the
2 encrypted volume did not contain the original text in any discernible form.

3 For the second scenario, we tried to read data from the volume already mounted on a running Galaxy instance, using
4 the OpenStack cloud controller. We were not able to gain access to the LUKS encrypted device by any means without
5 providing the correct passphrase. All in all, we believe that the approach described here can be safely used to insulate
6 any data uploaded to an encrypted Galaxy instance from malicious access as long as the instance itself and the
7 encryption keys remain uncompromised.

8 Cluster support

9 From the point of view of the IaaS administrator, the option to offer static and/or elastic cluster support to users
10 provides the alternative between guaranteeing a constant pool of resources to those instances attached to a static
11 cluster, or greater control over the efficient usage of the available resources, for those instances that instead rely on an
12 elastic cluster. The latter solution (Fig. 7), in fact, dynamically scales the number of cluster nodes available to a Galaxy
13 instance depending on its workload. From the user point of view, both solutions enable straightforward access to
14 computational resources beyond those assigned to the Galaxy instance virtual hardware, enabling a higher number of
15 simultaneous users, more rapid execution of jobs and additional room for computationally intensive analyses.

16



17

18 Fig. 7. Galaxy elastic cluster architecture. Initially, only the master node, that hosts Galaxy, SLURM, and CLUES, is deployed.
19 The SLURM queue is monitored by CLUES, and new worker nodes are deployed to process pending jobs up to the maximum
20 number set during cluster configuration, thus adapting resources availability to the current workload. The user home
21 directory and persistent storage are shared among master and worker nodes through the Network File System (NFS)
22 enabling the sharing of CONDA tools dependencies. If and when the dependencies are not satisfied by CONDA, the required
23 packages are installed during deployment on each worker node. The CVMFS shared volume is also mounted on each worker
24 node so to ensure that tools have access to reference data.

25

1 Laniakea pilot instance

2
3 In order to test and demonstrate Laniakea, we have deployed a pilot service over an OpenStack (Mitaka release) IaaS
4 hosted at the ELIXIR-IT ReCaS-Bari facility [46]. The current INDIGO PaaS layer configuration of this pilot service is
5 reported in Table S4 (Supplementary2.docx). During the first closed beta program, starting Dec 2018, up to 128 CPUs,
6 256 GB of RAM and 10 TB of disk storage will be reserved for the service. Resources will be gradually increased over the
7 next few months in order to make this prototype instance grow into a full-fledged service provided by the Italian Node
8 of ELIXIR (<https://www.elixir-europe.org/>). The service front-end is available at elixir-italy-laniakea.cloud.ba.infn.it.

9 Discussion

10 Laniakea offers a feature-rich solution to include a Galaxy on-demand service within the portfolio of public and private
11 cloud providers. This result is achieved leveraging the INDIGO middleware that, having been designed to support a vast
12 array of scientific services, may already be present and supported within the same cloud infrastructure or, as an
13 alternative, can also be used to pilot locally available computational resources from a remote INDIGO deployment.
14 Laniakea's scalability encompasses a variety of possible Galaxy setups that span from small instances to serve, e.g., small
15 research groups, developers and didactic purposes to production grade instances with (elastic) cluster support enabling
16 multiple concurrent users and computationally demanding analyses. New Galaxy flavours, allowing to implement and
17 make readily available future or existing data analysis pipelines, can be quickly deployed and shared through Ansible
18 recipes to ease the error-prone and lengthy routines of tools installation. We are confident that the Galaxy PaaS
19 delivered with Laniakea can effectively mitigate the need to host and maintain local hardware and software
20 infrastructures in several different scenarios, favouring a more efficient use of the available resources, harnessing the
21 improved reliability offered by cloud environments and also helping to enhance the reproducibility of bioinformatics
22 analyses. Finally, the data security layer of Laniakea addresses relevant issues for the analysis of sensitive data in human
23 genetics research and clinical settings. Future developments will be aimed at improving the compatibility of Laniakea
24 with a broader array of existing cloud setups, and at extending cluster support to other resource managers (e.g.,
25 TORQUE [31], HTCondor [47], etc...). Finally, we plan to expand the compatibility of Laniakea with Docker in order to
26 offer Galaxy containers as long-running services, exploiting the great number of dockerized Galaxy flavours already
27 hosted at Docker Hub (<https://hub.docker.com/>) and maintained by the Galaxy community.

1 Conclusions

2 Laniakea represents a clear example of the current trend of services virtualization, following the direction set forth, e.g.,
3 by the European Open Science Initiative (EOSC) Declaration and enables researchers and scientific services providers to
4 implement several recommendations therein outlined. Indeed, Laniakea offers a platform-agnostic, cloud-based service
5 that can be almost effortlessly kept up to date; this, in turn, facilitates the provision of software and services for the Life
6 Science field and beyond. We believe that by easing the barriers posed by the software and hardware layers required
7 to deploy and maintain Galaxy instances, and thus by enabling better access to cutting-edge technology to a broader
8 audience of researchers and other stakeholders, the approach adopted by Laniakea can contribute significantly to the
9 efficient employment of computational resources in the coming years. Finally, Laniakea design could be easily adapted
10 to deploy other tools beyond Galaxy in order to provide additional instruments for the Life Science community or to
11 serve different scientific communities.

12 Acknowledgements

13 This work has been supported by the European Commission H2020 research and innovation program under grant
14 agreement RIA 653549.

15 References

- 16 1. Piccolo SR, Frampton MB. Tools and techniques for computational reproducibility. *Gigascience*.
17 2016;5:1–13.
- 18 2. Kumar S, Dudley J. Bioinformatics software for biologists in the genomics era. *Bioinformatics*.
19 2007;23:1713–7.
- 20 3. Attwood TK, Blackford S, Brazas MD, Davies A, Schneider MV. A global perspective on evolving
21 bioinformatics and data science training needs. *Brief Bioinform* [Internet]. 2017;1–7. Available
22 from: [http://academic.oup.com/bib/article/doi/10.1093/bib/bbx100/4096809/A-global-](http://academic.oup.com/bib/article/doi/10.1093/bib/bbx100/4096809/A-global-perspective-on-evolving-bioinformatics)
23 [perspective-on-evolving-bioinformatics](http://academic.oup.com/bib/article/doi/10.1093/bib/bbx100/4096809/A-global-perspective-on-evolving-bioinformatics)
- 24 4. Beckmann JS, Lew D. Reconciling evidence-based medicine and precision medicine in the era of
25 big data: Challenges and opportunities. *Genome Med* [Internet]. *Genome Medicine*; 2016;8:1–11.
26 Available from: <http://dx.doi.org/10.1186/s13073-016-0388-7>
- 27 5. Cohen-Boulakia S, Belhajjame K, Collin O, Chopard J, Froidevaux C, Gaignard A, et al. Scientific
28 workflows for computational reproducibility in the life sciences: Status, challenges and
29 opportunities. *Futur Gener Comput Syst* [Internet]. Elsevier B.V.; 2017;75:284–98. Available from:
30 <http://dx.doi.org/10.1016/j.future.2017.01.012>
- 31 6. Afgan E, Baker D, Batut B, van den Beek M, Bouvier D, Čech M, et al. The Galaxy platform for
32 accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic Acids Res*
33 [Internet]. 2018;46:W537–44. Available from: <http://dx.doi.org/10.1093/nar/gky379>
- 34 7. Mell P, Grance T. The NIST Definition of Cloud Computing Recommendations of the National
35 Institute of Standards and Technology. *Nist Spec Publ*. 2011;145:7.
- 36 8. Langmead B, Nellore A. Cloud computing for genomic data analysis and collaboration. *Nat Rev*
37 *Genet* [Internet]. Nature Publishing Group; 2018;19:208–19. Available from:
38 <http://www.nature.com/doi/10.1038/nrg.2017.113>

- 1 9. Karim R, Michel A, Zappa A, Baranov P, Sahay R, Rebholz-schuhmann D. Improving data
2 workflow systems with cloud services and use of open data for bioinformatics research. *Brief*
3 *Bioinform* [Internet]. 2017;1–16. Available from:
4 http://fdslive.oup.com/www.oup.com/pdf/production_in_progress.pdf
- 5 10. Pavlovich M. Computing in Biotechnology: Omics and Beyond. *Trends Biotechnol* [Internet].
6 Elsevier Ltd; 2017;35:479–80. Available from: <http://dx.doi.org/10.1016/j.tibtech.2017.03.011>
- 7 11. Warth B, Levin N, Rinehart D, Teijaro J, Benton HP, Siuzdak G. Metabolizing Data in the Cloud.
8 *Trends Biotechnol* [Internet]. Elsevier Ltd; 2017;35:481–3. Available from:
9 <http://dx.doi.org/10.1016/j.tibtech.2016.12.010>
- 10 12. Griebel L, Prokosch HU, Köpcke F, Toddenroth D, Christoph J, Leb I, et al. A scoping review of
11 cloud computing in healthcare. *BMC Med Inform Decis Mak*. 2015;15:1–16.
- 12 13. Bellazzi R. Big Data and Biomedical Informatics : A Challenging Opportunity Big Data : Why
13 Bother ? Big Data : Must-have or. *Yearb Med Inform*. 2014;9:8–13.
- 14 14. Liu B, Madduri RK, Sotomayor B, Chard K, Lacinski L, Dave UJ, et al. Cloud-based bioinformatics
15 workflow platform for large-scale next-generation sequencing analyses. *J Biomed Inform*
16 [Internet]. Elsevier Inc.; 2014;49:119–33. Available from:
17 <http://dx.doi.org/10.1016/j.jbi.2014.01.005>
- 18 15. Afgan E, Baker D, Coraor N, Chapman B, Nekrutenko A, Taylor J. Galaxy CloudMan: Delivering
19 cloud compute clusters. *BMC Bioinformatics*. 2010;11:2–7.
- 20 16. Afgan E, Chapman B, Taylor J. CloudMan as a platform for tool, data, and analysis distribution.
21 *BMC Bioinformatics* [Internet]. *BMC Bioinformatics*; 2012;13:1. Available from: *BMC*
22 *Bioinformatics*
- 23 17. Afgan E, Sloggett C, Goonasekera N, Makunin I, Benson D, Crowe M, et al. Genomics Virtual
24 Laboratory: A practical bioinformatics workbench for the cloud. *PLoS One*. 2015;10:1–20.
- 25 18. Peters K, Bradbury J, Bergmann S, Capuccini M, Aauri P De, Ebbels TMD, et al. PhenoMeNal :
26 Processing and analysis of Metabolomics data in the Cloud. 2018;
- 27 19. Krieger MT, Torreno O, Trelles O, Kranzlmüller D. Building an open source cloud environment
28 with auto-scaling resources for executing bioinformatics and biomedical workflows. *Futur Gener*
29 *Comput Syst* [Internet]. Elsevier B.V.; 2017;67:329–40. Available from:
30 <http://dx.doi.org/10.1016/j.future.2016.02.008>
- 31 20. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Čech M, et al. The Galaxy
32 platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic*
33 *Acids Res*. 2016;44:W3–10.
- 34 21. Afgan E, Baker D, Coraor N, Goto H, Paul IM, Makova KD, et al. Harnessing cloud computing
35 with Galaxy Cloud. *Nat Biotechnol*. 2011;29:972–4.
- 36 22. Campos DSI, Marco LGJ, Solagna DLP, Matyska JGL, Hardt PFM, Dutka GDL, et al. INDIGO-
37 DataCloud : a Platform to Facilitate Seamless Access to E-Infrastructures. *J Grid Comput* [Internet].
38 2018; Available from: <https://link.springer.com/article/10.1007%2Fs10723-018-9453-3>
- 39 23. Salomoni D, Campos I, Gaido L, Donvito G, Antonacci M, Fuhrman P, et al. INDIGO-Datacloud:
40 foundations and architectural description of a Platform as a Service oriented to scientific
41 computing. 2016;1–31. Available from: <http://arxiv.org/abs/1603.09536>
- 42 24. Lipton P (Ca T, Moser S (Ibm), Palma D (Vnomic), Spatzier T (Ibm). Topology and Orchestration
43 Specification for Cloud Applications - PRIMER. 2013;1–114. Available from: [http://docs.oasis-](http://docs.oasis-open.org/tosca/TOSCA/v1.0/cs01/TOSCA-v1.0-cs01.html)
44 [open.org/tosca/TOSCA/v1.0/cs01/TOSCA-v1.0-cs01.html](http://docs.oasis-open.org/tosca/TOSCA/v1.0/cs01/TOSCA-v1.0-cs01.html)
- 45 25. OASIS. TOSCA Simple Profile in YAML Version 1 . 0 Committee Specification Draft 04 / Public
46 Review Draft 01. 2015; Available from: [http://docs.oasis-open.org/tosca/TOSCA-Simple-Profile-](http://docs.oasis-open.org/tosca/TOSCA-Simple-Profile-YAML/v1.0/csprd01/TOSCA-Simple-Profile-YAML-v1.0-csprd01.pdf)
47 [YAML/v1.0/csprd01/TOSCA-Simple-Profile-YAML-v1.0-csprd01.pdf](http://docs.oasis-open.org/tosca/TOSCA-Simple-Profile-YAML/v1.0/csprd01/TOSCA-Simple-Profile-YAML-v1.0-csprd01.pdf)
- 48 26. Ben-Kiki O, Evans C, Ingerson B. YAML Ain't Markup Language (YAML™) Version 1.2. Language

- 1 (Baltim) [Internet]. 2009;1–100. Available from: <http://www.yaml.org/spec/1.2/spec.html>
- 2 27. Housley R, Polk W, Ford W, Solo D. Internet X.509 Public Key Infrastructure Certificate and
- 3 Certificate Revocation List (CRL) Profile. United States: RFC Editor; 2002.
- 4 28. OpenID Foundation. OpenID Connect Discovery 1.0 incorporating errata set 1. 2014;311376.
- 5 Available from: http://openid.net/specs/openid-connect-discovery-1_0.html
- 6 29. Cantor S, Hodges J, Hirsch F, Philpott R, Security RS a, Hughes J, et al. Profiles for the OASIS
- 7 Security Assertion Markup Language (SAML). Language (Baltim) [Internet]. 2005;16:66. Available
- 8 from:
- 9 [http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:Profiles+for+the+OASIS+Security+Assertion+Markup+Language+\(SAML\)+V2.0#0](http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:Profiles+for+the+OASIS+Security+Assertion+Markup+Language+(SAML)+V2.0#0)
- 10
- 11 30. Yoo AB, Jette MA, Grondona M. SLURM: Simple Linux Utility for Resource Management. In:
- 12 Feitelson D, Rudolph L, Schwiegelshohn U, editors. Job Sched Strateg Parallel Process. Berlin,
- 13 Heidelberg: Springer Berlin Heidelberg; 2003. p. 44–60.
- 14 31. Staples G. TORQUE Resource Manager. Proc 2006 ACM/IEEE Conf Supercomput [Internet].
- 15 New York, NY, USA: ACM; 2006. Available from: <http://doi.acm.org/10.1145/1188455.1188464>
- 16 32. Fruhwirth C. New methods in hard disk encryption. Inst Comput Lang Theory Log ... [Internet].
- 17 2005; Available from: http://git.dyne.org/tomb/plain/doc/New_methods_in_HD_encryption.pdf
- 18 33. Buncic P, Aguado Sanchez C, Blomer J, Franco L, Harutyunian A, Mato P, et al. CernVM - A
- 19 virtual software appliance for LHC applications. J Phys Conf Ser. 2010;219.
- 20 34. Grüning BA, Fallmann J, Yusuf D, Will S, Erxleben A, Eggenhofer F, et al. The RNA workbench:
- 21 Best practices for RNA and high-throughput sequencing bioinformatics in Galaxy. Nucleic Acids
- 22 Res. 2017;45:W560–6.
- 23 35. Chiara M, Gioiosa S, Chillemi G, D’Antonio M, Flati T, Picardi E, et al. CoVaCS: a consensus
- 24 variant calling system. BMC Genomics [Internet]. BMC Genomics; 2018;19:120. Available from:
- 25 <https://bmcbgenomics.biomedcentral.com/articles/10.1186/s12864-018-4508-1>
- 26 36. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short
- 27 DNA sequences to the human genome. Genome Biol. 2009;10.
- 28 37. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods.
- 29 2012;9:357–9.
- 30 38. Li H, Li H, Durbin R, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
- 31 transform. Bioinformatics [Internet]. 2009;25:1754–60. Available from:
- 32 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2705234%5C&tool=pmcentrez%5C&rendertype=abstract%5Cnpapers2://publication/doi/10.1093/bioinformatics/btp324>
- 33
- 34 39. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform.
- 35 Bioinformatics. 2010;26:589–95.
- 36 40. Sherry ST. dbSNP: the NCBI database of genetic variation. Nucleic Acids Res [Internet].
- 37 2001;29:308–11. Available from: [https://academic.oup.com/nar/article-](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/29.1.308)
- 38 [lookup/doi/10.1093/nar/29.1.308](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/29.1.308)
- 39 41. Mills RE, Pittard WS, Mullaney JM, Farooq U, Creasy TH, Mahurkar AA, et al. Natural genetic
- 40 variation caused by small insertions and deletions in the human genome. Genome Res.
- 41 2011;21:830–9.
- 42 42. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti A, et al. A global
- 43 reference for human genetic variation. Nature. 2015;526:68–74.
- 44 43. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome
- 45 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.
- 46 Genome Res. Cold Spring Harbor Lab; 2010;
- 47 44. Yuchi X, Shetty S. Enabling security-aware virtual machine placement in IaaS clouds. Proc - IEEE
- 48 Mil Commun Conf MILCOM. 2015;2015–Decem:1554–9.

- 1 45. Oechslin P. Making a Faster Cryptanalytic Time-Memory Trade-Off. 2003;617–30. Available
- 2 from: http://link.springer.com/10.1007/978-3-540-45146-4_36
- 3 46. Antonacci M, Bellotti R, Cafagna F, de Palma M, Diacono D, Donvito G, et al. The ReCaS Project:
- 4 The Bari Infrastructure. High Perform Sci Comput Using Distrib Infrastructures Results Sci Appl
- 5 Deriv from Ital PON ReCaS Proj. World Scientific; 2017. p. 17–33.
- 6 47. Thain D, Tannenbaum T, Livny M. Distributed computing in practice: The Condor experience.
- 7 *Concurr Comput Pract Exp.* 2005;17:323–56.
- 8