



HAL
open science

Multi-Cloud deployment for microbial genomes analysis

Jonathan Lorenzo, Bryan Brancotte, Thomas Lacroix, Mohamed Bedri,
Jean-Francois Gibrat, Christophe Blanchet

► To cite this version:

Jonathan Lorenzo, Bryan Brancotte, Thomas Lacroix, Mohamed Bedri, Jean-Francois Gibrat, et al.. Multi-Cloud deployment for microbial genomes analysis. JOBIM 2017 - Journées Ouvertes Biologie Informatique Mathématiques, Jul 2017, Lille, France. , pp.234, 2017, JOBIM 2017- LILLE. hal-01605864

HAL Id: hal-01605864

<https://hal.science/hal-01605864>

Submitted on 2 Jun 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - ShareAlike 4.0 International License

Multi-Cloud deployment for microbial genomes analysis



project H2020 644925

Jonathan Lorenzo 1, Bryan Brancotte 1, Thomas Lacroix2, Mohammed Behri 1, Jean-François Gibrat 1 et Christophe Blanchet 1

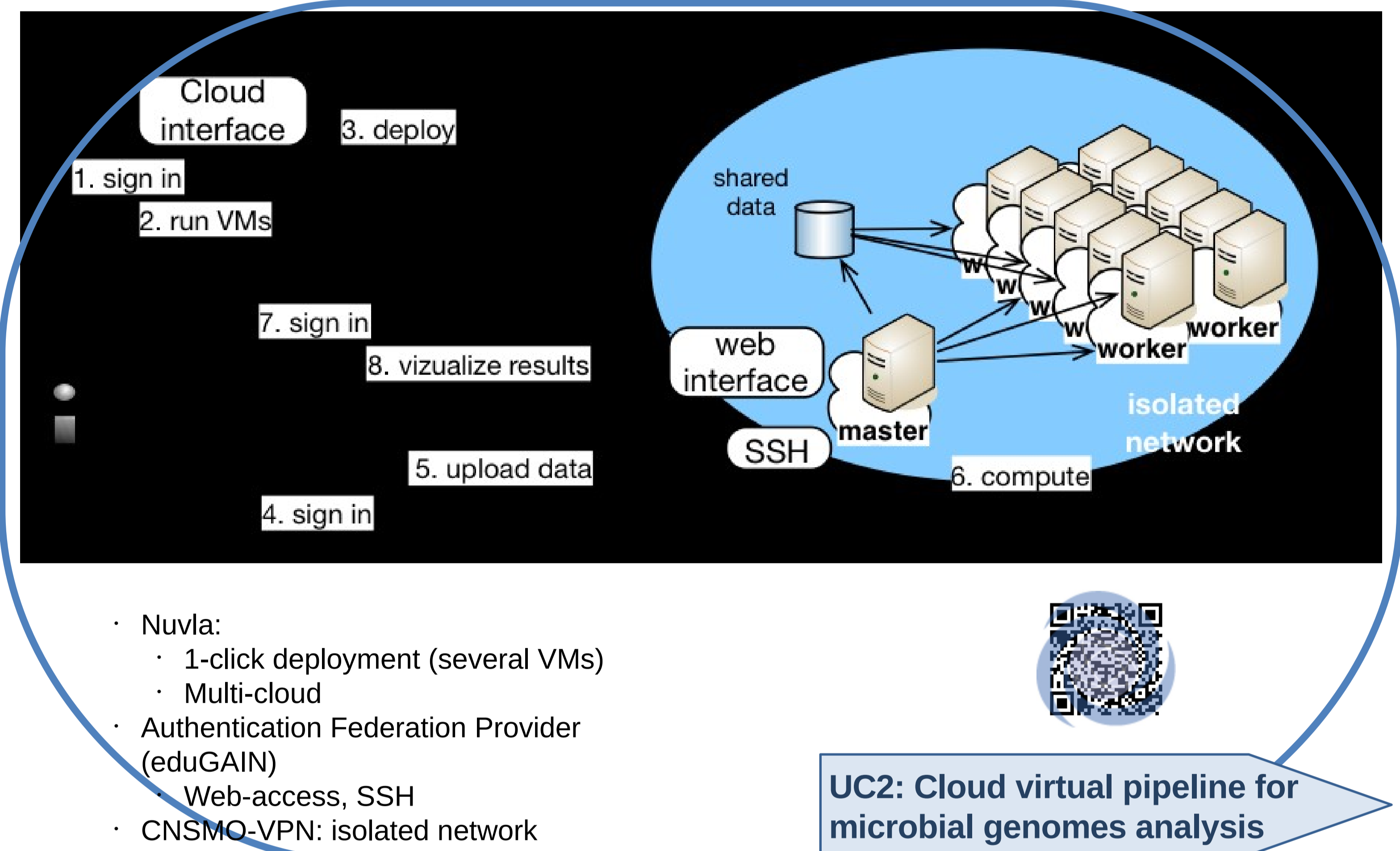
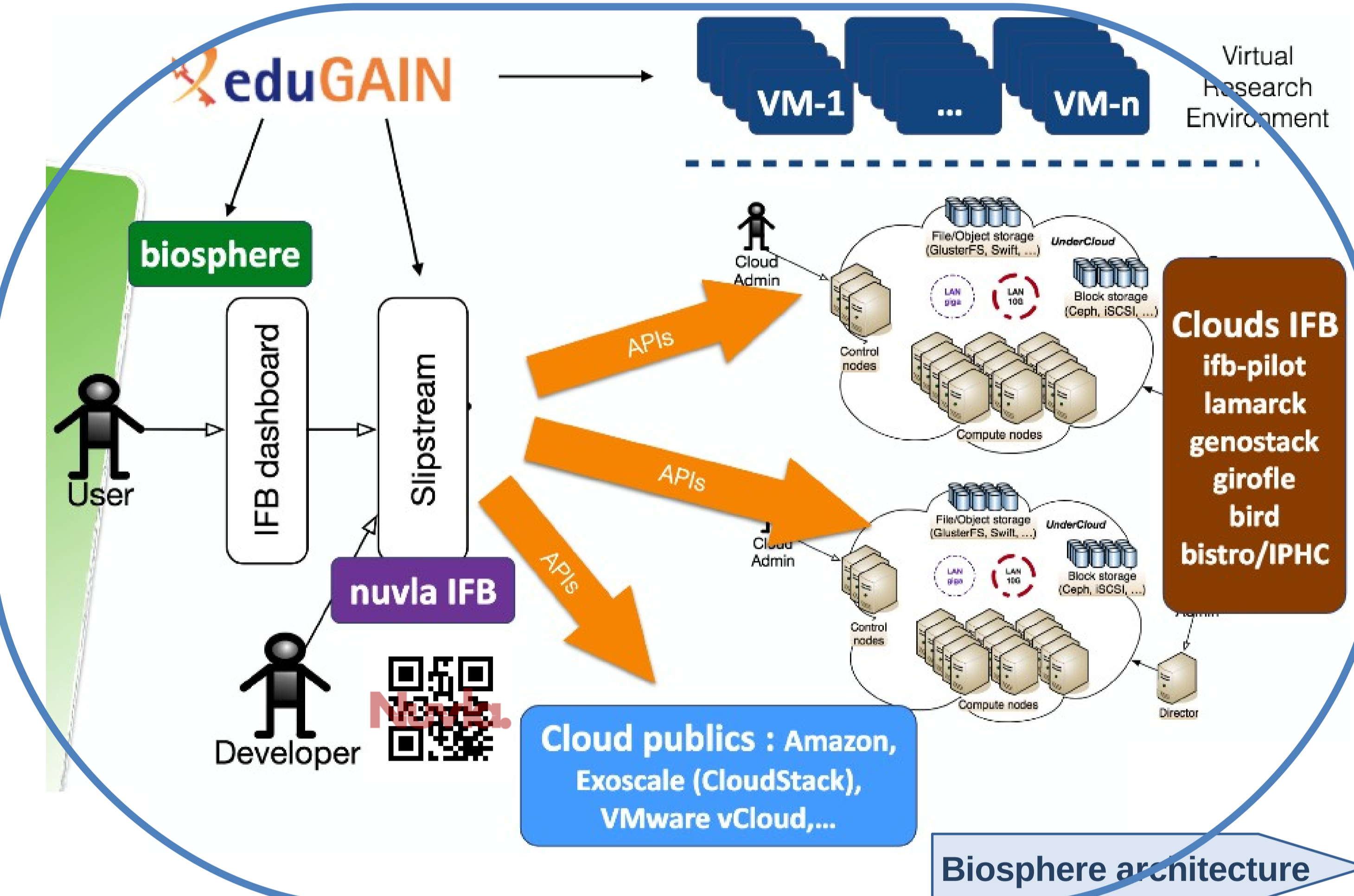
1 CNRS, UMS 3601 ; Institut Français de Bioinformatique, IFB-core, Avenue de la Terrasse, F-91190 Gif-sur-Yvette, France
2 MaIAGE, INRA, Université Paris-Saclay, 78350 Jouy-en-Josas, France

Contact : jonathan.lorenzo@france-bioinformatique.fr, Thomas.lacroix@inra.fr

Introduction

In the post-NGS area, sequencing bacterial genomes is very cheap (few hundreds €). Most of the time, users are no longer content to analyse a single genome; they want to compare large collections of related genomes (strains). This entails that biologists have to pay too much attention and dedicate their time to sequence the genomes, instead of thoroughly analysing the genomic data. Thus, this brings light to the increasing need for automating the annotation of bacterial genomes and carrying out efficient data mining.

In that context, IFB hub and the IFB-MIGALE platform developed a virtual environment (appliance), based on virtual machines, called "bacterial genomics" that aims to provide biologists and bioinformaticians access to suitable resources via the cloud. For example, Prokka [1] is a software tool for the rapid annotation of prokaryotic genomes. Insyght [2] developed by IFB-MIGALE is a tool for the visualization of the synteny (local conservation of the gene order along the genomes) and the exploration of the landscape of both conserved and idiosyncratic genomic regions across multiple genomes. The platform automatically launches a set of bioinformatics tools (e.g. BLAST, HMMER, Prodigal...) to analyse the data and stores the results in a relational database (PostgreSQL). These tools use several public reference data collections. A web interface allows the user to browse the results. Setting up the platform requires solid skills in system administration since many bioinformatics tools with different dependences need to be installed as well as a relational database management system, a web server and servlet container, etc. Moreover, performing the analysis of a large number of genomes requires large computing resources and the use of parallel computing.



- Nuvla:
 - 1-click deployment (several VMs)
 - Multi-cloud
- Authentication Federation Provider (eduGAIN)
 - Web-access, SSH
- CNSMQ-VPN: isolated network

UC2: Cloud virtual pipeline for microbial genomes analysis

How to run Bacterial Genomics appliance for an analysis

1 Search a Bacterial Genomics appliance in RAINBio, the Biosphere catalogue

2 Select an appliance and run it in one-click (*)

4 Run analysis

```

[ifbuser@master-1 ~] exit
logout
Connection to 192.54.201.220 closed.
MaxBowl-Pre-3-4 ssh ifbuser@192.54.201.220
Warning: Permanently added '192.54.201.220' (ECDSA) to the list of known hosts.
Last login: Tue Jun 20 14:32:19 2017 from 157.136.202.14
ifbuser created for launch job
ghost : show the status of Sun Grid Engine hosts, queues, jobs
qstat -f : show the status of Sun Grid Engine jobs and queues
[ifbuser@master-1 ~]
    
```

3 Direct Access on Insyght portal (**)

5 Shutdown appliance at the end of analysis

(*) You can also adjust resources and target cloud
(**) Access by eduGAIN after authentication in the federation identity (no ssh-key needed) or by ssh-key for developer

Conclusion

The goal is to deploy the "appliance" in one click over one or more cloud infrastructures. To achieve this, new features to automate deployment of complex application were added to the IFB's cloud portal through the connection to the SlipStream cloud broker. Developed by SixSq, SlipStream is a multi-cloud application management platform. It automates the full application management lifecycle, within Infrastructure as a Service (IaaS) cloud infrastructures. Such complex application deployments can be done over several cloud infrastructures and provide scientists with high-level cloud features such as the dynamic allocation of a dedicated network for the isolation of the virtual machines, with the replication of the user data and with a direct link from the cloud portal to the Insyght web portal.

[1] Seemann T. Prokka: rapid prokaryotic genome annotation. Bioinformatics. 2014 Jul 15;30(14):2068-9. PMID:24642063
[2] Lacroix T., Loux V., Gendraud A., Hoebeke M., and Gibrat J.F. Insyght: navigating amongst abundant homologues, synteny and gene functional annotations in bacteria, it's that symbol! Nucleic Acids Res. 2014 Dec 1; 42(21): e162. doi: 10.1093/nar/aku867 · PMID: PMC4245967