# Bioinformatics Cloud Services for Life Sciences

Christophe BLANCHET

**Institut Français de Bioinformatique - IFB**
**French Institute of Bioinformatics - ELIXIR-FR**
CNRS UMS3601 - Gif-sur-Yvette - FRANCE

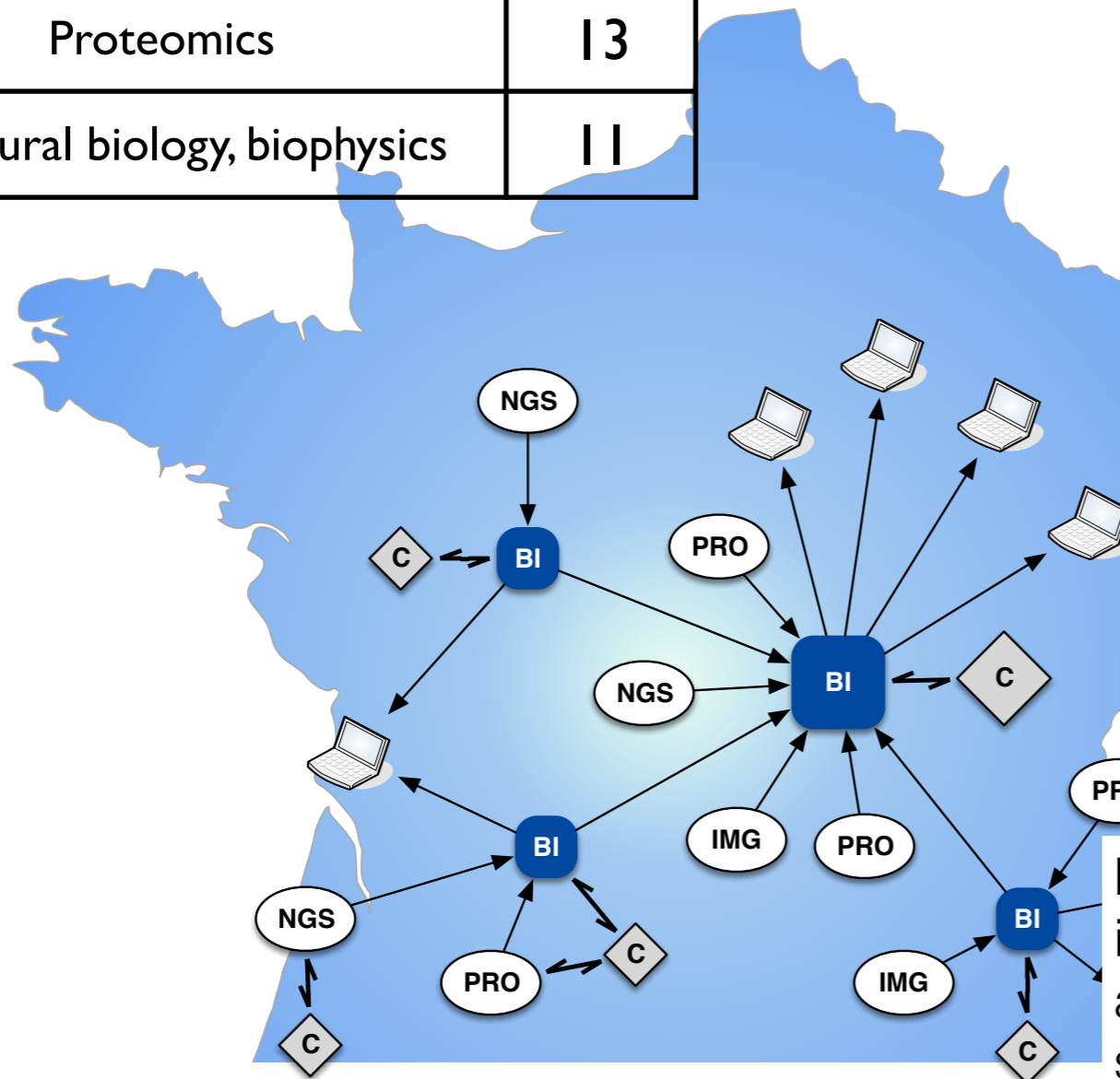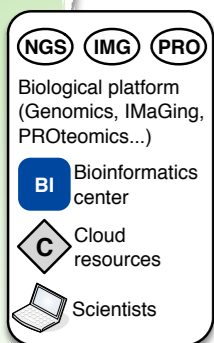Nettab Conference
15 October 2015, Bari

# Experimental data producers in life sciences (FR)

| French national platforms (GIS IBISA) | Nb |
|---|---|
| Cellular imaging | 18 |
| Genomics, transcriptomics | 16 |
| Proteomics | 13 |
| Structural biology, biophysics | 11 |

French NGS platforms

Source: omicsmaps.com

Regional centers distribute the load in terms of computing and storage, and provide better interactions with scientists

NGS  IMG  PRO
Biological platform (Genomics, IMaGing, PROteomics...)

BI  Bioinformatics center

C  Cloud resources

Scientists

Un déluge de donnée. Blanchet C. et Collin O., 2011, Biofutur, 323: 64-67

2

# A lot of bioinformatics tools

**tools**

BLAST
R
FastA
OMSSA
SSearch
ClustalW2
PeptideShaker
ARIA
HMMer
BWA
X!tandem
TopHat
samtools
Galaxy
Clustal
fastQC
Muscle
Omega

BWA 0.6.2
BWA 0.7.10
CAP3
CD-HIT 4.6.1
Clustal Omega 1.0.3
CLUSTALW 2.1
Cufflinks 2.0.2
Cutadapt 1.2.1
E-SURGE 1.9.0
Exonerate 2.2.0
eXpress 1.5.1
FastA 3.6
FastQC 0.10.1
Galaxy portal
GATK 2.3.4
HMMer 3.0
ImageJ 1.48
khmer 1.1
M-SURGE 1.8.5
MEME 4.7

ABYSS 1.3.4
ARIA 2.3
Bioconductor 2.11
biomaj
BLAST+ 2.2.27
Blat 35
Bowtie 0.12.8
Bowtie2 2.0.0-beta7

MMSEQ 0.11.2a
Mobyle
MODAL
MultAlin 5.4.1
MUSCLE 3.8.31
neo4j
Oases 0.2.08
OMSSA 2.1.9
PeptideShaker 0.18.3
phyml 3.1
PREDATOR 2.1.2
proline
python 2.7
R 2.13
R 3.1.1
R 3.1.2
R-studio
Ray 1.3
RSAT
samtools 0.1.18

Samtools 1.1
SearchGUI 1.10.4
SeqClean
Shiny
Stacks
STAR 2.4.0f1
SuMo v1
TGICL
TopHat 2.0.6
trim_galore 0.3.7
Trinity 2.0.4
U-CARE 2.3.2
VCFtools 0.1.11
Velvet 1.2.10
X!tandem 12-10-01-1
XPLOR-NIH 2.30
…

# Many interfaces

# The French Institute of Bioinformatics
# and its e-infrastructure

# History

**Since 2004, ReNaBi is the National Network of Bioinformatics platforms with an IBiSA label (Infrastructures in Biology, Health and Agronomy)**

**In 2010, call of proposals "Infrastructures in Biology and Health" from the "Investments for the Future" initiative.**

★ Project ReNaBi-IFB accepted in 2012 and endowed with 20m €

**Other national infrastructures (NIs)**

★ France Génomique : sequencing and genotyping NI
★ Profi : proteomics NI
★ Frisbi : structural biology NI
★ etc. (17 NIs all together) + 5 IHUs (Instituts Hospitaliers Universitaires) + 1 IRT (Institut de Recherche Technologique)

# IFB - Institut Français de Bioinformatique

## French distributed infrastructure for life-science information



http://www.france-bioinformatique.fr

CNRS UMS3601. Avenue de la Terrasse, Bât 21. 91190 Gif-sur-Yvette

## Mission : to make available core bioinformatics resources to the life science research community.

- To provide **support for national biology programs**
- To provide an **IT infrastructure** devoted to management and analysis of biological data
- To act as a middleman between the life science community and the bioinformatics/computer science research community
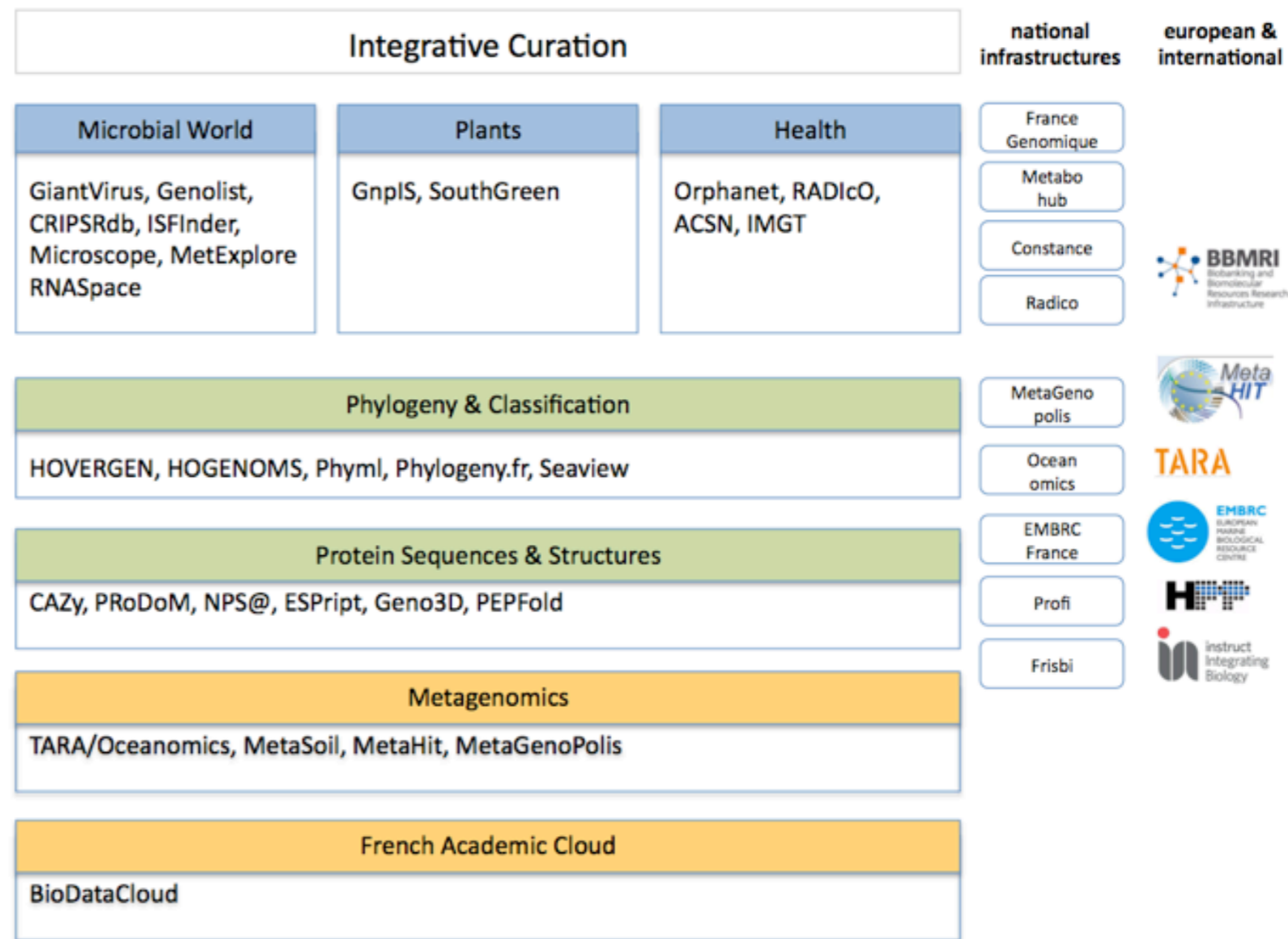
### ELIXIR French Node

- The European distributed infrastructure for life-science information
- To optimize the **interactions and coordination** between the national level and ELIXIR and other ESFRI infrastructures in biomedical and environmental field,
- To promote **consistency and complementarities** between the components offered by the ELIXIR French node and those of other European nodes

# Support to projects

## Support to biological, biomedical or technological projects

- **Large scale institutional projects and projects with other infrastructures**

- **Technological projects for developing services and tools**

- **Biology and biomedical research projects**

- **Services offered to industry**

| | Integrative Curation | | | national infrastructures | european & international |
|---|---|---|---|---|---|
| main themes | **Microbial World** | **Plants** | **Health** | France Genomique | |
| | GiantVirus, Genolist, CRIPSRdb, ISFInder, Microscope, MetExplore RNASpace | GnpIS, SouthGreen | Orphanet, RADIcO, ACSN, IMGT | Metabo hub | |
| | | | | Constance | BBMRI |
| | | | | Radico | |
| transversal tools & data | **Phylogeny & Classification** | | | MetaGeno polis | Meta HIT |
| | HOVERGEN, HOGENOMS, Phyml, Phylogeny.fr, Seaview | | | Ocean omics | TARA |
| | **Protein Sequences & Structures** | | | EMBRC France | EMBRC |
| | CAZy, PRoDoM, NPS@, ESPript, Geno3D, PEPFold | | | Profi | HPP |
| | **Metagenomics** | | | Frisbi | Instruct Integrating Biology |
| | TARA/Oceanomics, MetaSoil, MetaHit, MetaGenoPolis | | | | |
| | **French Academic Cloud** | | | | |
| | BioDataCloud | | | | |

## Call for new proposals in progress

# IFB e-Infrastructure

**Mission : to provide core bioinformatics resources to the life science research community.**

- To set up a **French IT infrastructure (cloud)** devoted to management and analysis of biological data
- To provide hardware, data collections and bioinformatics tools
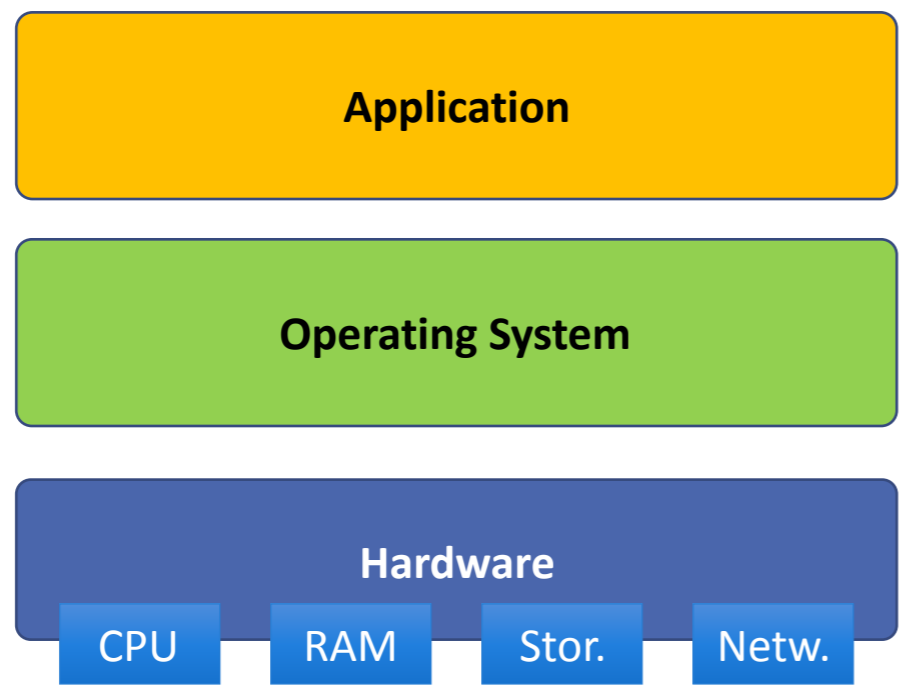- To collaborate with international infrastructure (ELIXIR)

**Current resources**

- A **national hub : IFB-core**
  IT resources hosted at CNRS IDRIS SC center
- A network of **regional centers**
  32 bioinformatics platforms - 15,000 cores - 5 PB
- 4 running clouds

➡️**Create a federation of clouds for life sciences**

# Virtualisation

With some limits…

Virtual machines
**1 … N**

| App | | App |
|---|---|---|
| Guest OS | | Guest OS |
| Virtual HW | | Virtual HW |
| Hypervisor | | |
| Hardware | | |

CPU  RAM  Stor.  Netw.

| Application |
|---|
| Operating System |
| Hardware |

CPU  RAM  Stor.  Netw.

Physical server

# IFB's Cloud-s-



**In IFB's premises**
- IFB-core (Gif)
- GenOuest (Rennes)

**In collaboration**
- BiLille/Univ.Lille (Lille)
- BISTRO/IPHC-EGI fedcloud (Strasbourg)

**PoC & experiments**
- URGI (Versailles)
- PRABI-LBBE (Lyon)

**=> Towards a federation**
- common identity and access management
- interoperability of VMs
- …

# IFB-core's cloud

| IFB-core | # Compute Cores | # TB Storage | # TB RAM | Max VM size | Technology | Location |
|----------|-----------------|--------------|----------|-------------|------------|----------|
| **Pilot** | **200** | **50** | **2** | **40c 256GB** | **StratusLab** | **CNRS-IDRIS, Paris** |
| *2016-S1* | *3,000* | *500* | *-* | *?144c 3TB?* | *StratusLab* | *CNRS-IDRIS, Paris* |
| *2017* | *10,000* | *2,000* | *-* | *??* | *StratusLab* | *CNRS-IDRIS, Paris* |

**Provide scientists with bioinformatics resources**
**- data and tools -**
**as cloud appliances**

# Create bioinformatics "appliances"

**tools**

BLAST  R

OMSSA

FastA

ClustalW2  SSearch

PeptideShaker

ARIA  BWA  X!tandem

TopHat  HMMer

samtools  Galaxy

Clustal  Muscle

Omega  fastQC

App | App

Guest OS | Guest OS

Virtual HW | Virtual HW

Hypervisor

Hardware

CPU | RAM | Stor. | Netw.

**Create new cloud services**

**+** Linux system **+**

**Virtual Machines**

## Bioinformatics Marketplace

**Sequences** **Structures** **Proteomics** **Galaxy** ...

### Appliance ?
- predefined virtual machine
- including tools, pipeline, recipes…
- Ready to run

## Appliance annotation
- Title
- Description (w. controlled voc.)
  - ★ Topics
  - ★ Tools
- Contact
- Developer(s) and **maintainer(s)** !

# Appliances - Topics

**Bioimaging**

**Ecology of population**

**Genomics tools**

**Mass Spectrometry**

**Molecular structural analysis**

**Multiple Sequence Alignment**

**Nucleotide and Protein sequence searching**

**Proteomics**

**Public databases**

**Sequence analysis**

**...**

# IFB's bioinformatics appliances



Scientific apps

CLI

Web

Remote desktop

Galaxy

Galaxy 'MODAL'

Galaxy AVIESAN 2015

Galaxy 'RADseq'

PhyML

bioCompute Node

RSAT

MacSyFinder

R statistics

Aria

Proteomics

Imaging

Eco Pop

SynBioWatch

BioDataCloud IGV

bioHadoop

Docker

Utiliti

bioData 'NFS'

bioData 'BioMaj'

Cassandra

BlobSeer

Ubuntu

CentOS

Base OS

Neo4j

Data mgmt

# Help developers to create appliance

## Appliances are created

- by the life science developers/experts of different domains

## Appliances in progress

- BioDataCloud-RNAseq
- ProFi
- REPET
- TriAnnot
- Clinical NGS for cancerology (CLB & CFB)
- Bacterial genomics (AGMIAL, Insygth)
- Metagenomics (iMetAMOS)
- …



manual

Linux system

**Galaxy + Tools**

automatic

Versions 1,2..n & Updates

Galaxy **NGS**

Galaxy **RAD-seq**

Galaxy **MODAL**

Galaxy **Aviesan**

Publish VMs

**Cloud Marketplace**

# Docking bioinformatics tools

docker.com

**IFB's docker hub**
(@GenOuest)

Registry
of images

**Container
layer**

| App | App |
|---|---|
| Docker Engine | Docker Engine |
| Guest OS | Guest OS |
| Virtual HW | Virtual HW |
| Hypervisor | |
| Hardware | |

CPU   RAM   Stor.   ima   Netw.   IB)

pull

push

**IFB's
Cloud**

docker
virtual machine

docker
virtual machine

**Developer**

**User**

abyss
blast
blast+
bowtie
bwa
clustalw2
cufflinks
fasta36
fastqc
gor4
hmm
meme
mmseq
cufflinks

0        275        550

# Managing biological data

## Collections of reference data

- Databases updates and index built in IFB-core (BioMAJ)
- Transfers from IFB-core to regional PFs

## Experimental data: archiving (and treatment)

- Regional desks: deposit
- Replicate to IFB-core (iRODS?)

## User data: distribution, optimisation, security

- Object storage (replication)
- Multi-site noSQL (distribution)
- Multi-site workflow (optimisation)
- Biomedical data (end-to-end security)



GenOuest

IFB-core

PRABI

# Move VMs rather than data



IFB's marketplace & VMs repository for life sciences

VMs

NGS · IMG · PRO
Biological platform (Genomics, IMaGing, PROteomics...)

**BI** Bioinformatics centre

**C** Cloud resources

Researchers

20

# IFB's bioinformatics cloud services

# A cloud for Bioinformatics

# A cloud driven through a web dashboard



http://cloud.france-bioinformatique.fr/cloud

# Browse the marketplace and run an App !

# RAINBio : Registry of bioinformatics tools and VMs

Prototype

**Query :**
**topic ? tool ?**
**VM ?**

**IFB's Cloud Marketplace**

**ELIXIR's Services Registry**

VMs

Tools

**Life science researcher**



**RAINBio**

**Graph DB (Neo4J)**

| topic | VMs | tools |
|---|---|---|
| Sequence comparison | BIO compute node - 3.1 | ClustalW2 |
| Gene expression | Galaxy MODAL - 1.0 | MPAgenomics |
| Bioinformatics | BIO compute node - 3.1, Galaxy MODAL - 1.0 | ABySS, MPAgenomics |
| Statistics | Galaxy MODAL - 1.0 | MPAgenomics |
| Phylogeny | PhyML - 0.2 | PhyML |
| Sequence search | BIO compute node - 3.1 | NCBI BLAST |
| SNP | Galaxy MODAL - 1.0 | MPAgenomics |
| Data search, query and retrieval | BIO Data - 1.2 | BioMAJ |
| Sequence assembly | BIO compute node - 3.1, Galaxy MODAL - 1.0 | ABySS |

# App R Statistical Computing



## R software environment for statistical computing and graphics

- include common bioinformatics module
- Biobase, BiocGenerics, BiocInstaller, GenomeInfoDb…

## RStudio IDE

- integrated development environment (IDE) for R
- features: console, syntax-highlighting editor …

## Shiny web framework

- powerful web framework for building web applications using R.
- without requiring HTML, CSS, or JavaScript knowledge.
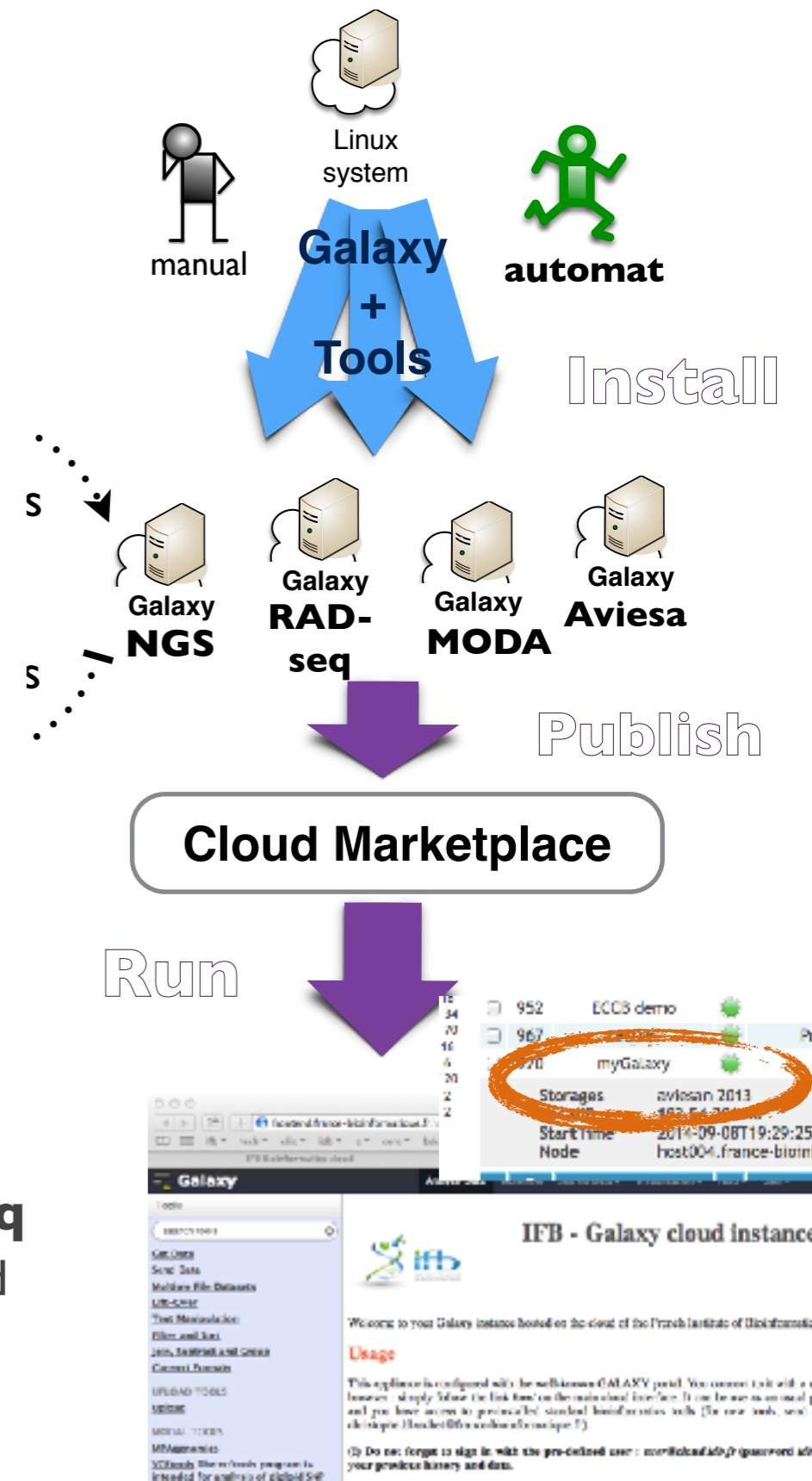
*Contact: Stéphane Delmotte (IFB PRABI-LBBE)*

# AppS Cloud Galaxy Portal

## Galaxy portal is widely used in the community

- analyse NGS data (mainly but not only)
- connected to community knowledge: data and indexes, tools, workflows

## Cloud advantages :

- User is **administrator of his/her own Galaxy** instance: he/she can install data and tools
- Preserve **workflows and results in cloud storage**
- Help the integration of monthly updates and new tools
- Different appliances can be available at the same time:
    - ★ a basic one with common tools for NGS
    - ★ specific ones for a domain or a set of tools
      e.g. Galaxy-MODAL, Galaxy-RADseq, **EBA-ChIP-Seq**
    - ★ or for training: create a special appliance with dedicated datasets, tools or workflows
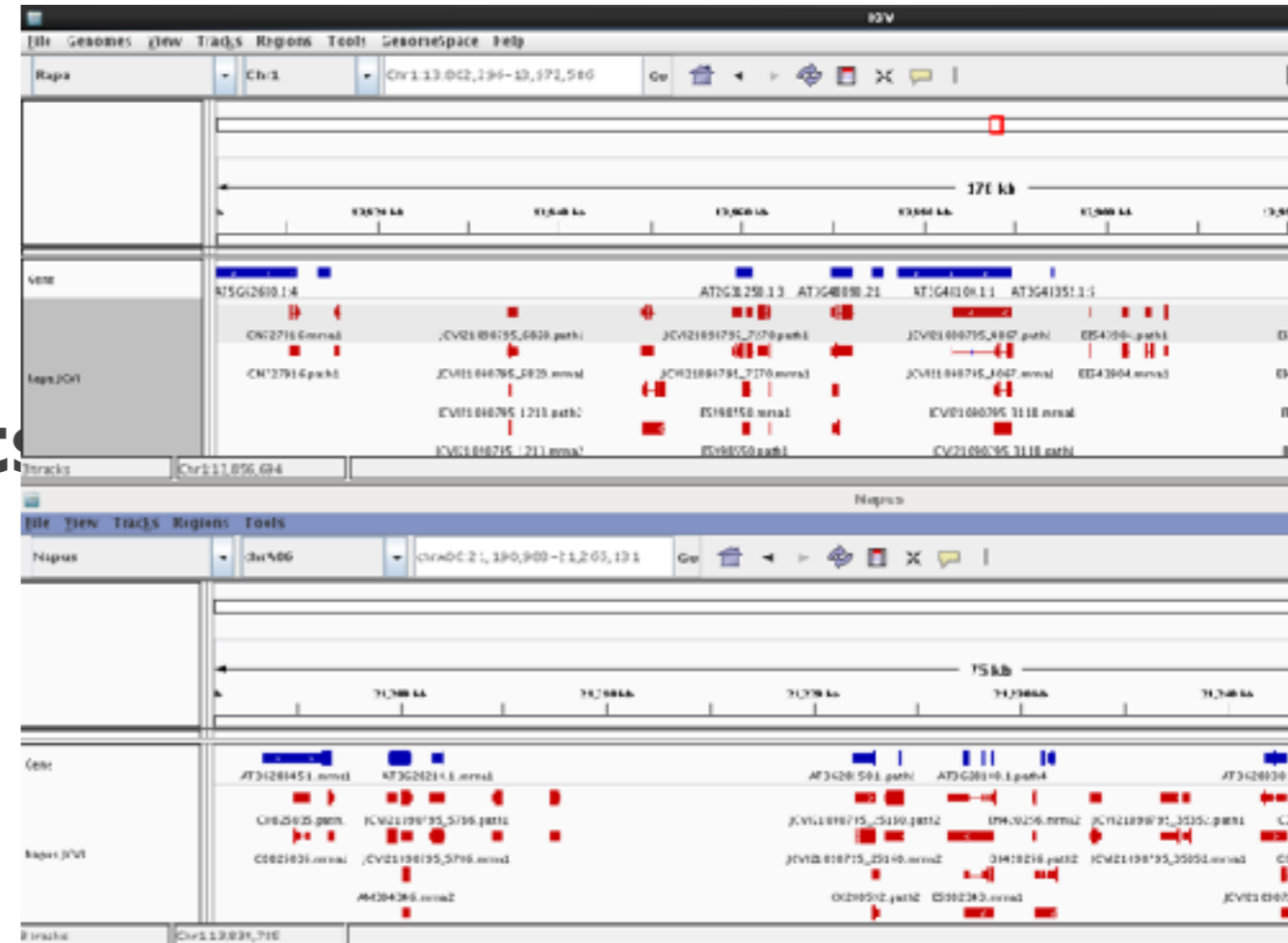      e.g. AVIESAN school 2015



manual

Linux system

**Galaxy + Tools**

automat

Install

S

S

Galaxy **NGS**

Galaxy **RAD-seq**

Galaxy **MODA**

Galaxy **Aviesa**

Publish

**Cloud Marketplace**

Run

# App Multi-genomes browser



Integrative Genomics Viewer

**Based on IGV**

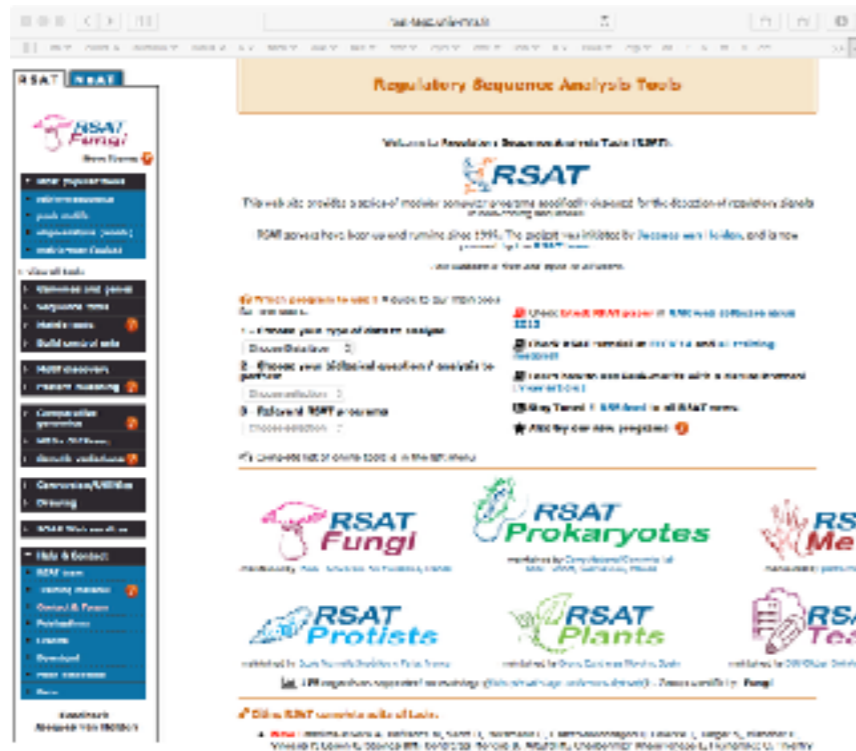**Ready to deploy in the cloud close to the datasets**

**Remote virtual desktop**

- transfer only graphical visualization
- based on NX protocol

*Contact: Marie-Laure Franchinard (IFB MIGALE)*
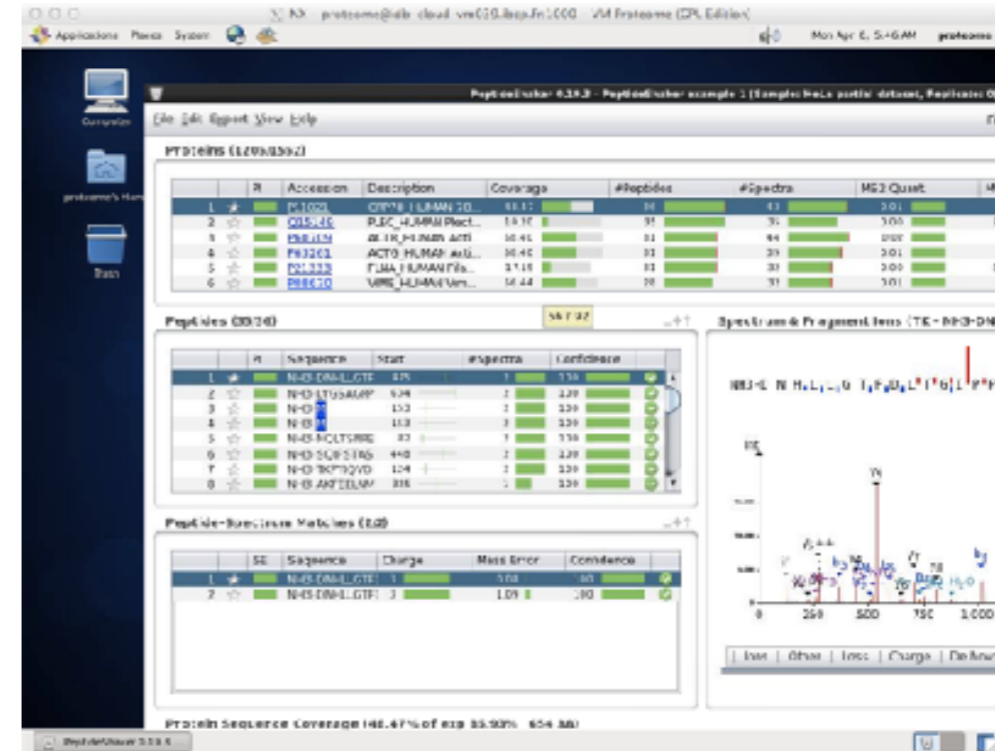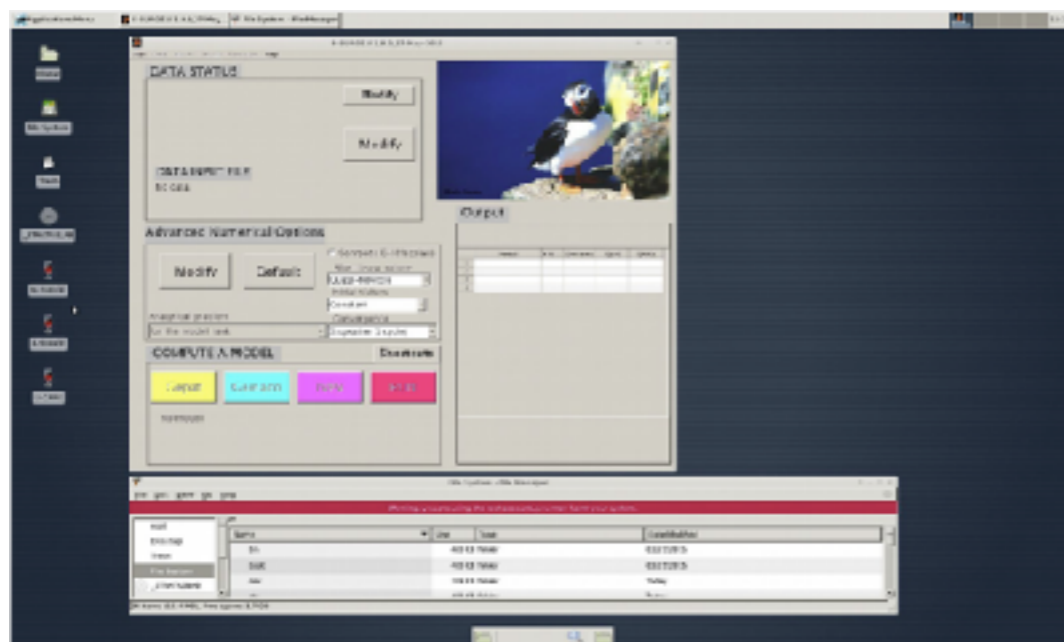
*Funded by the French BIODATACLOUD project.*

# And other apps ...

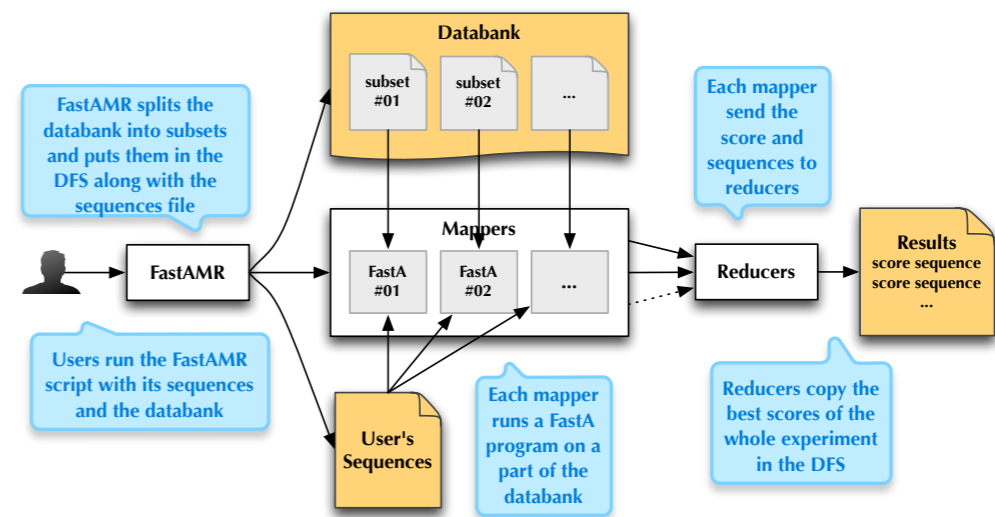**RSAT**

**Proteomics**

**Ecology of populations**

**Hadoop**

*etc.*

# Conclusion: IFB's cloud today

## 22 bioinformatics appliances already available

- **+ 10 in progress** by the experts of the different life sciences domains
  - ★ BioDataCloud-**RNAseq**, **ProFi**, **Clinical NGS for cancerology** (2x), REPET, TriAnnot, Galaxy **RAD-seq**, **Bacterial genomics**, **iMetAMOS**…
- IFB supports different domain-specific developments
  - ★ Microbial Bioinformatics, Evolutionary bioinformatics, Plant bioinformatics, Structural Biology, NGS data processing, biomedical data analysis…
  - ★ Call for new proposals in progress

## Scientific production - 239 users (October 2015)

- opened to members of IFB (standard allocated resources)
- opened to partners, academic and industry, infrastructures and projects: e.g. BioDataCloud, ProFi, MetaboHub, …
- extra resources allocation according to scientific and financial criteria

## Training

- Scientific school "Cumulo NumBio - Cloud Computing for Life Sciences" (Aussois, June 2015)
- IFB's tutorials for cloud end-users and appliance developers
- tutorial at ECCB'14 about 'Analysis of Cis-Regulatory Motifs from High-Throughput Sequence Sets'
- Bioinformatics Masters in Marseille (2014) and Rouen (2015)
- Scientific school about Genomics with Galaxy (2015)

# Questions ?

## Acknowledgments

- **IFB** members
  - IFB hub: **Patricia**, **Awa, Jean-François, Mohamed, Jonathan, Maxime, Dominique**
    *Alumni : **Marie, Quentin***
    ➡ *we are hiring !*
  - Working group IFB-GRISBI (co-chair with Olivier Collin)

- **Appliances developers**
  **Samuel** Blanck (Inria Lille), **Jacques** van Helden (TAGC), **Stéphane** Delmotte (PRABI-LBBE), **Bruno** Spataro (PRABI-LBBE), **Marie-Laure** Franchinard (MIGALE), **Anis** Djari (BioinfoGenoToul), **Bertrand** Néron (Institut Pasteur), **Adrien** Josso (MicroScope), **Thomas** Lacroix (MIGALE), **Christian** Baudet (CLB), **Germain** Paimparay & **Baptiste** Brault (CFB)**…**

- **CNRS IDRIS**: R. Medeiros, C. Gauthey and staff
- **StratusLab** members

- IFB is funded by **French programs PIA INBS 2012, BioDataCloud**
- **EU H2020 projects, CYCLONE (644925) and EGI-Engage (6541**