

The Genomics Virtual Laboratory

Andrew Lonie

Victorian Life Sciences Computation Initiative
University of Melbourne



What is the Genomics Virtual Lab?

**Nationally distributed platform for
genomics, built on the federal
*Research Cloud***



R D S I
Research Data Storage
Infrastructure



nectar

<http://nectar.org.au>

The Australian Research Cloud



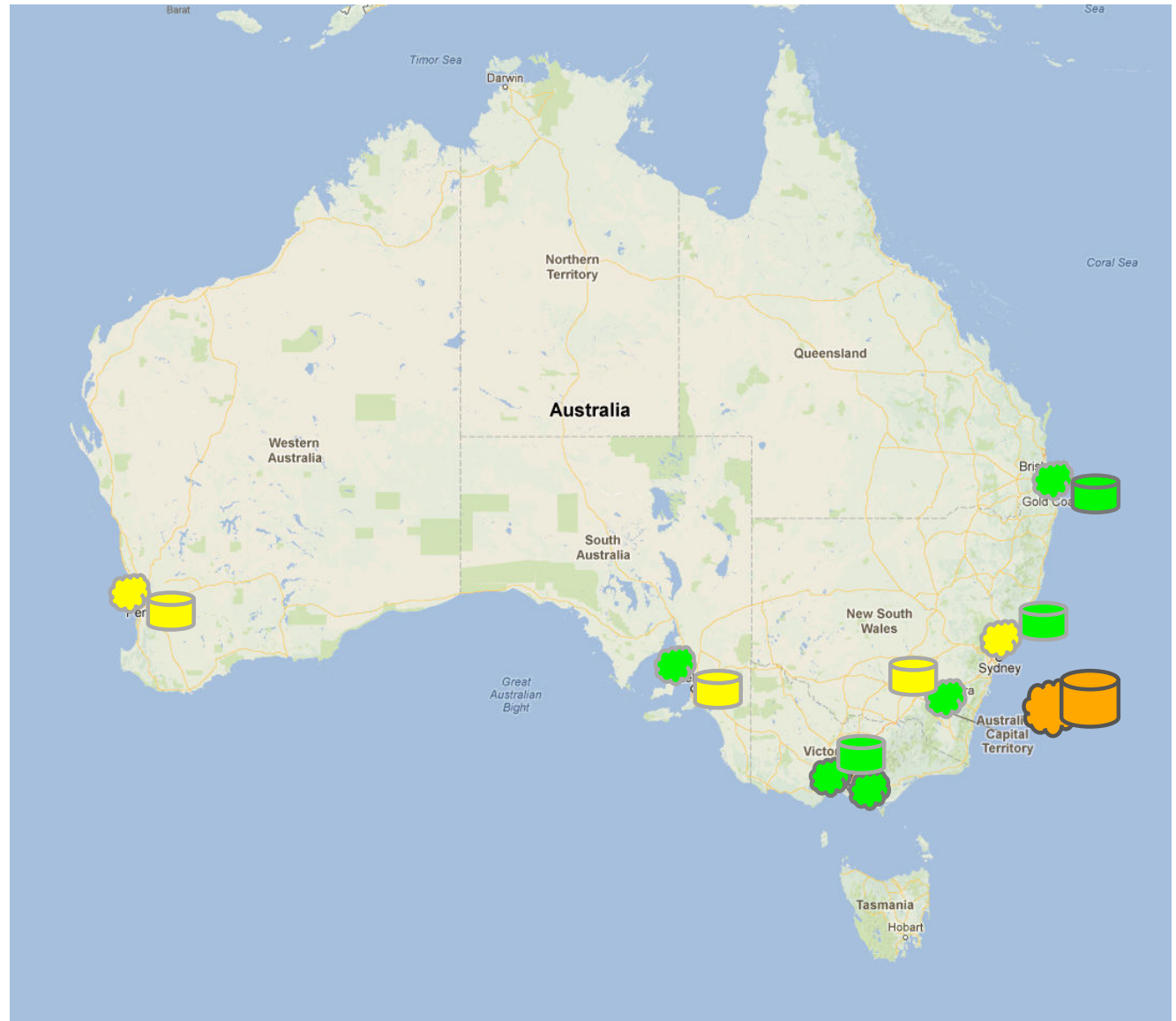
Cloud node:
3-6000 cores



Data node:
1-5 PB



Coming 2014-15



GVL: Drivers

To provide a **genomics analysis platform** with:

- 1.Reproducibility*
- 2.Accessibility*
- 3.Performance*
- 4.Flexibility*
- 5.Consistency*
- 6.Functionality*

for as many researchers as possible

GVL: Design principles

Criteria	Design Implication
<i>Accessible</i>	Minimal client-side requirements
<i>Reproducible</i>	Workflow support + software & tool management process
<i>Performance</i>	User-managed scaling of compute resources + high availability resources
<i>Flexible</i>	User configurable + administrable Multiple interaction modes
<i>Consistent</i>	Single platform from training to analysis
<i>Functional</i>	Pre-populated with suite of tools for common use cases + required reference data + visualisation options

GVL: Design implications

Criteria	Design Implication	Technical implication
Accessible	Minimal client-side requirements	<u>Web based</u> tool and management interfaces
Reproducible	Workflow support + software & tool management process	<u>Workflow platforms</u> + automated process for <u>deployable underlying environment</u>
Performance	User-managed scaling of compute resources + high availability resources	<u>Cloud-based architecture</u> + interface for managing resources
Flexible	User configurable + administrable	<u>Per-user instances</u> accessible through web and command line; user-administrable environment
Consistent	<u>Single platform from training to analysis</u>	<u>Tutorials and guides</u> for training using best practice tools + <u>scalability</u>
Functional	<u>Pre-populated with suite of tools for common use cases + required reference data</u> + visualisation options	Process for <u>building underlying images</u> Automated configuration of <u>reference datasets</u>

For as many researchers as possible...

Galaxy Main



GVL: Philosophy

Genomics Virtual Lab



Galaxy Main

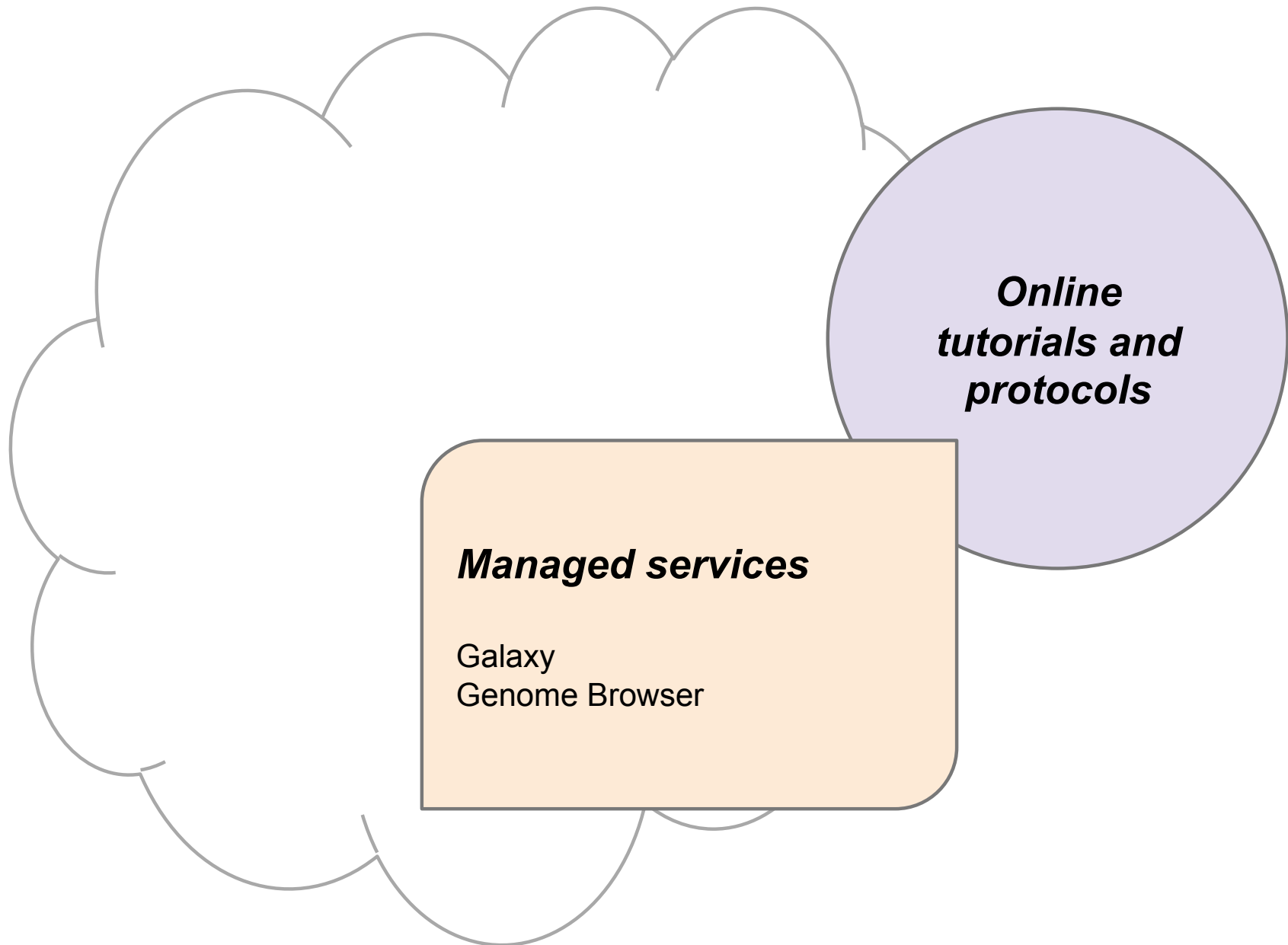


GVL: In practice

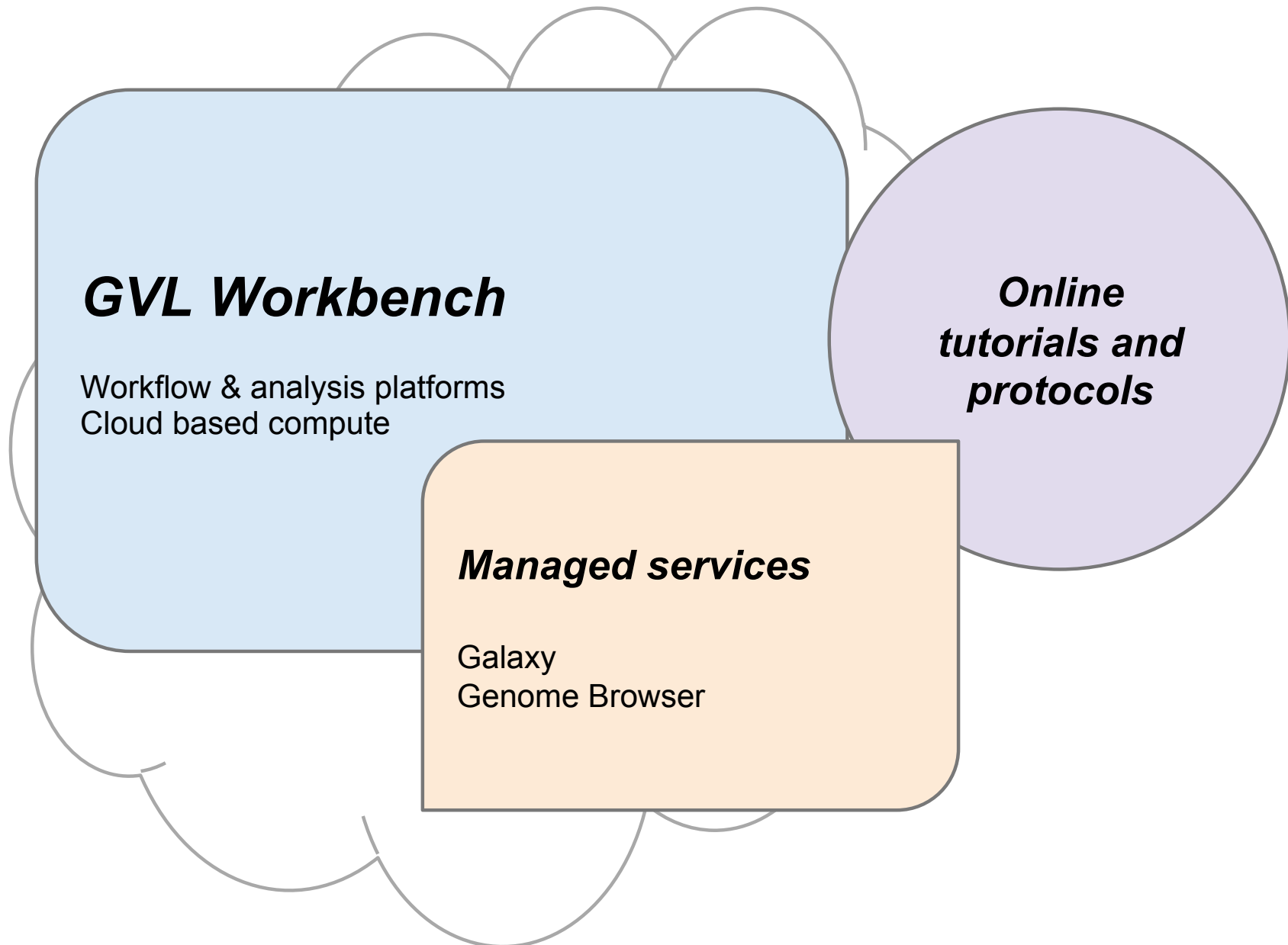


***Online
tutorials and
protocols***

GVL: In practice



GVL: In practice



GVL: <http://genome.edu.au>

GET

GVL Workbench

Workflow & analysis platforms
Cloud based compute

LEARN

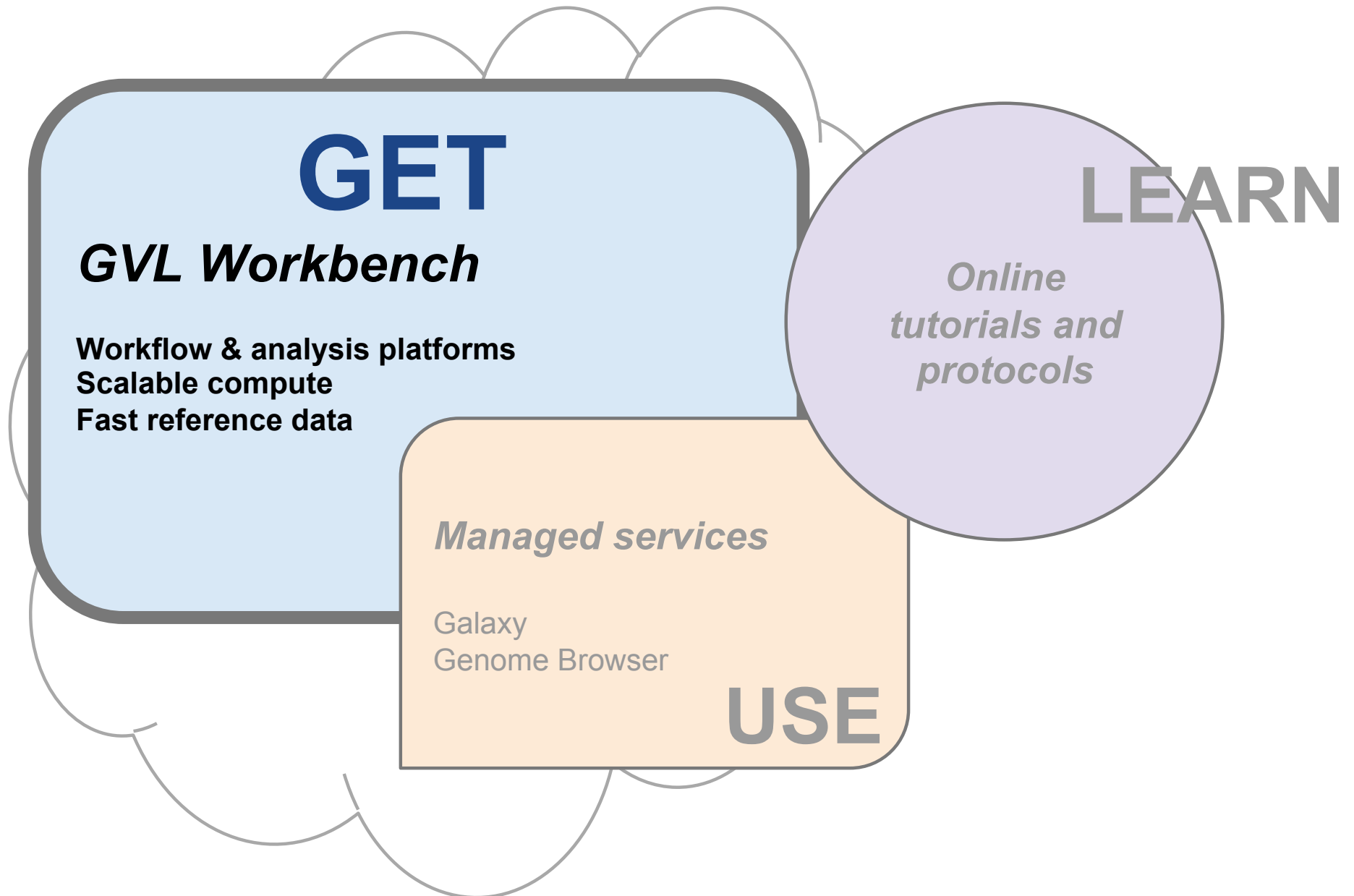
*Online
tutorials and
protocols*

Managed services

Galaxy
Genome Browser

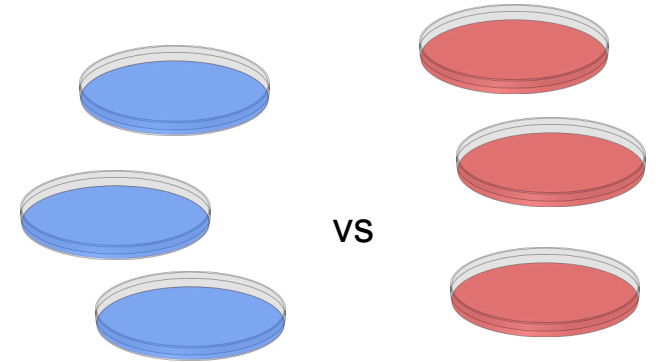
USE

GVL: Developer's perspective



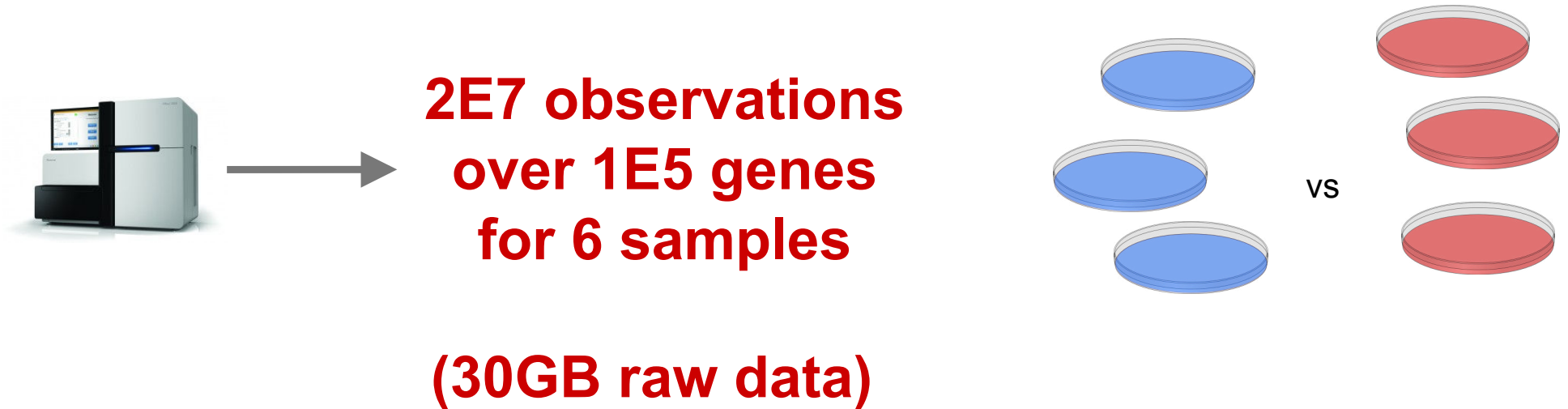
What characterises genomic analysis?

eg: Differential Gene Expression



“What genes are turned on in **blue cells** and turned off in **red cells**?”

Differential Gene Expression



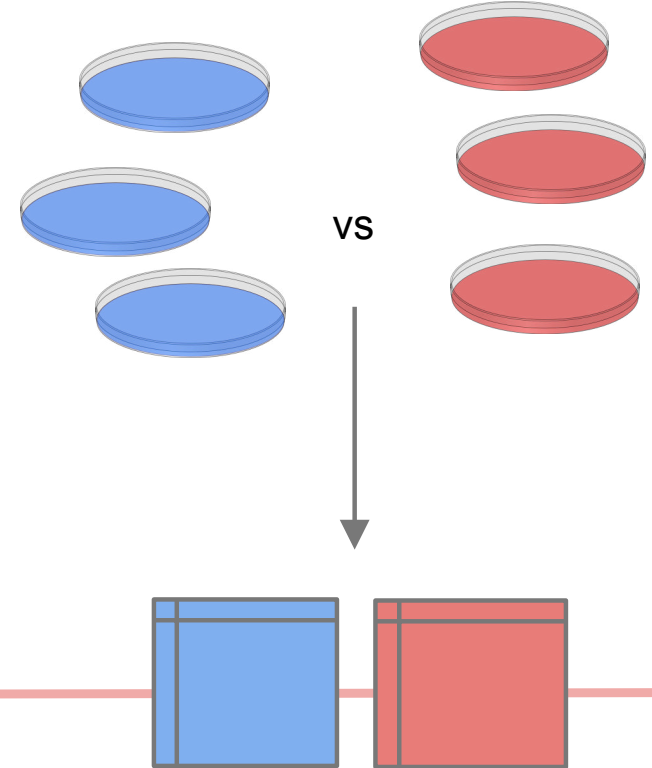
“What genes are turned on in **blue cells** and turned off in **red cells**?”

Differential Gene Expression

DATA REDUCTION

Workflows
Large reference data
High compute

**2E7 observations over
1E5 genes for 2x3
samples**



Differential Gene Expression

DATA REDUCTION

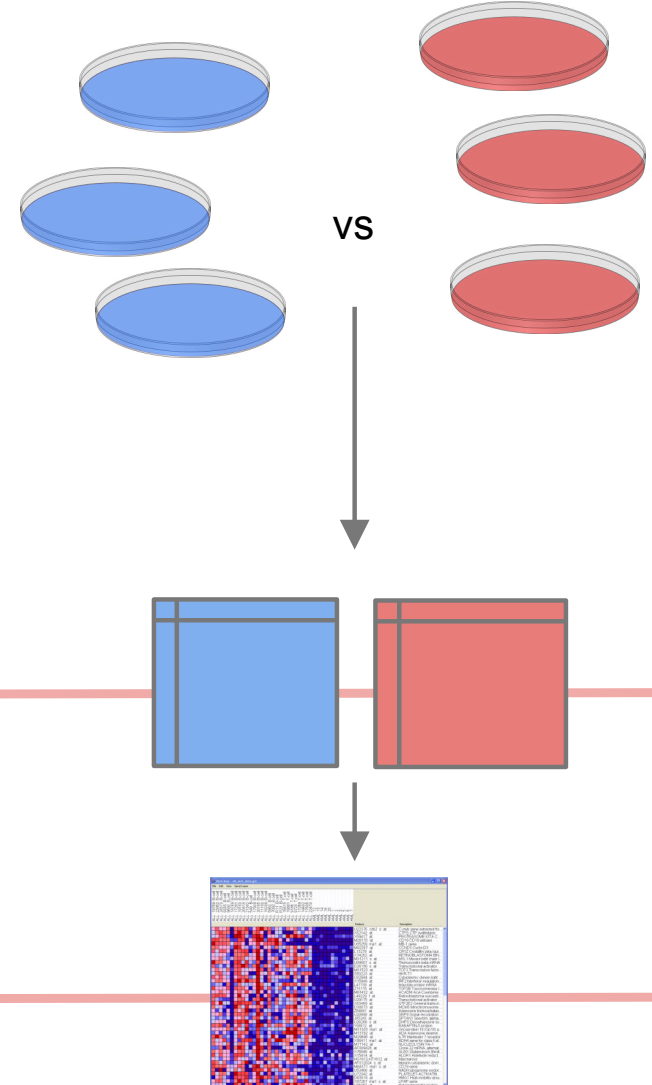
Workflows
Large reference data
High compute

DATA ANALYSIS

Interactive
Flexibility

**2E7 observations over
1E5 genes for 2x3
samples**

**Test difference:
[1,2,3] vs [4,5,6]**



Differential Gene Expression

DATA REDUCTION

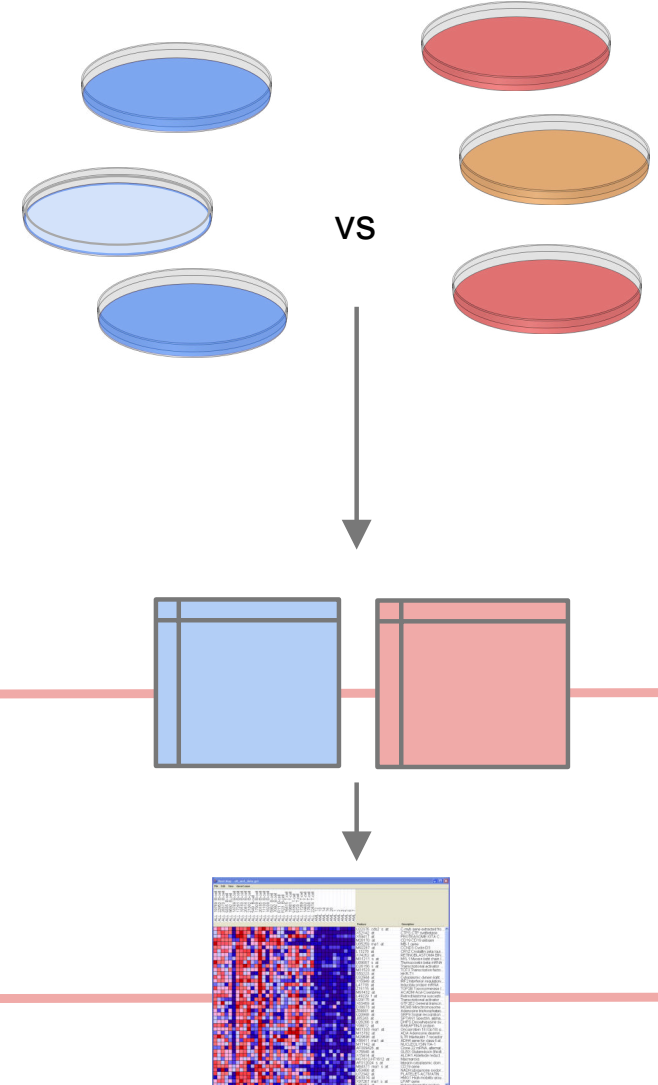
Workflows
Large reference data
High compute

DATA ANALYSIS

Interactive
Flexibility

**2E7 observations over
1E5 genes for 2x3
samples**

**Test difference:
[1,2,3] vs [4,5,6]**



Differential Gene Expression

DATA REDUCTION

Workflows
Large reference data
High compute

DATA ANALYSIS

Interactive
Flexibility

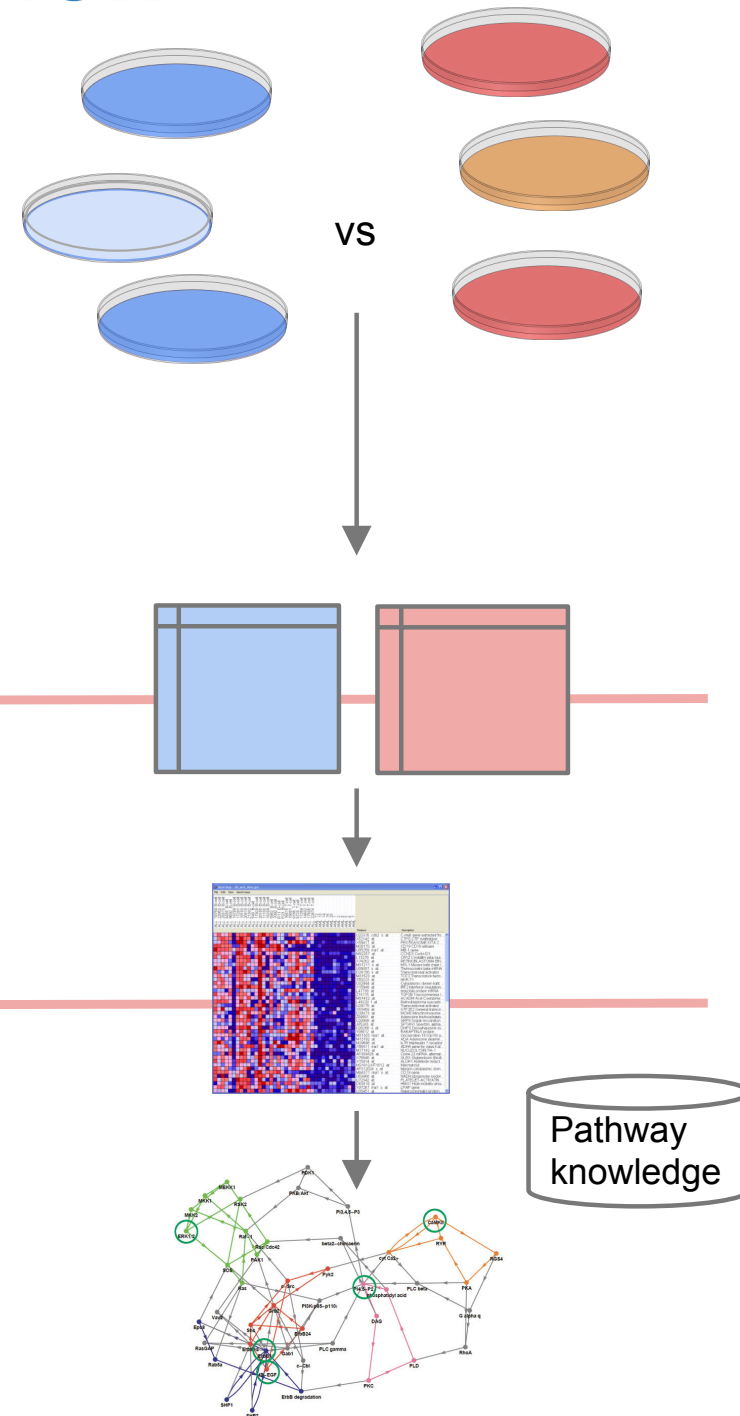
DATA INTERPRETATION

External services
Visualisation tools

**2E7 observations over
1E5 genes for 2x3
samples**

**Test difference:
[1,2,3] vs [4,5,6]**

**Biological pathways
affected? Mechanism?**



Differential Gene Expression

Characteristic

S

DATA REDUCTION

Workflows

Large reference data

High compute

DATA ANALYSIS

Interactive

Flexibility

DATA INTERPRETATION

External services

Visualisation tools

**2E7 observations over
1E5 genes for 2x3
samples**

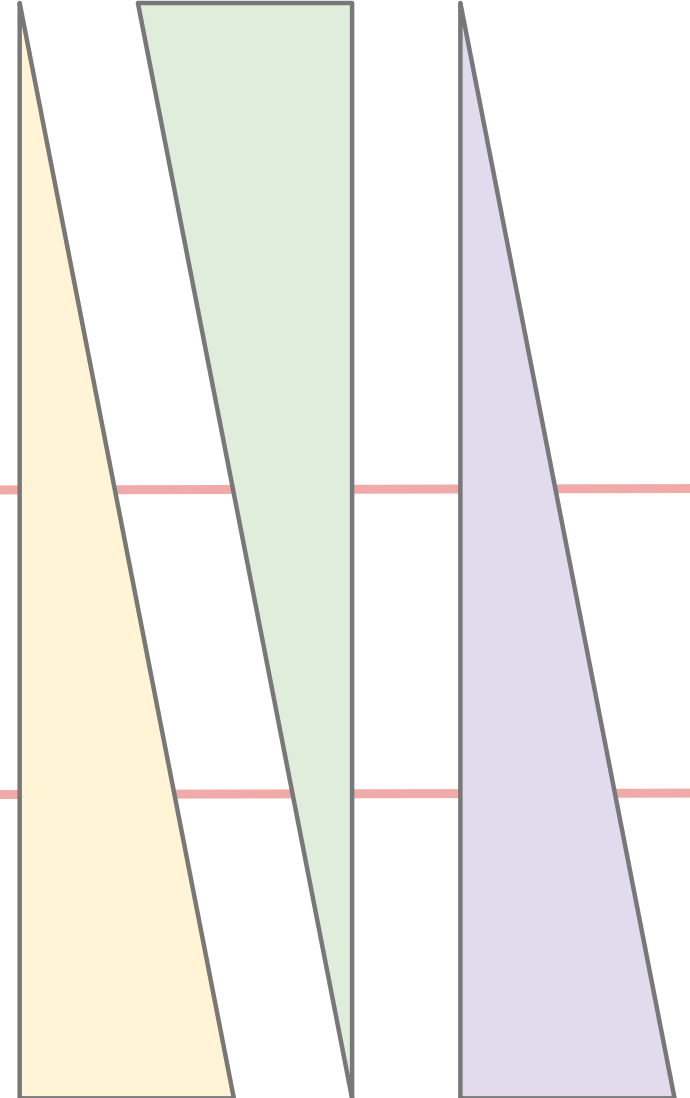
**Test difference:
[1,2,3] vs [4,5,6]**

**Biological pathways
affected? Mechanism?**

Tools

Data

Domain Context



What characterises genomics?

- Very large experimental datasets per user/group
- I/O intensive high compute initial analysis
 - ‘data reduction’: raw data to sample summaries
- Large suite of data analysis tools, interactive
 - a bit subjective
- Complex context for interpretation, external tools
 - more subjective, domain knowledge
- Little modelling/simulation

GVL Workbench: Requirements

A web-based *per-user* workbench providing:

- access to multiple tools
- on a scalable back end compute cluster
- with fast access to large reference data,
- user administrable and configurable
- with multiple modes of interaction
- and a mechanism for reproducible workflows

all highly available and accessible

i.e. with a minimal cost of entry to the user

Why per-user?

Managed service: objective

A short time later...



Why per-user?

Managed service: objective



A short time later...



GVL: Philosophical assertion



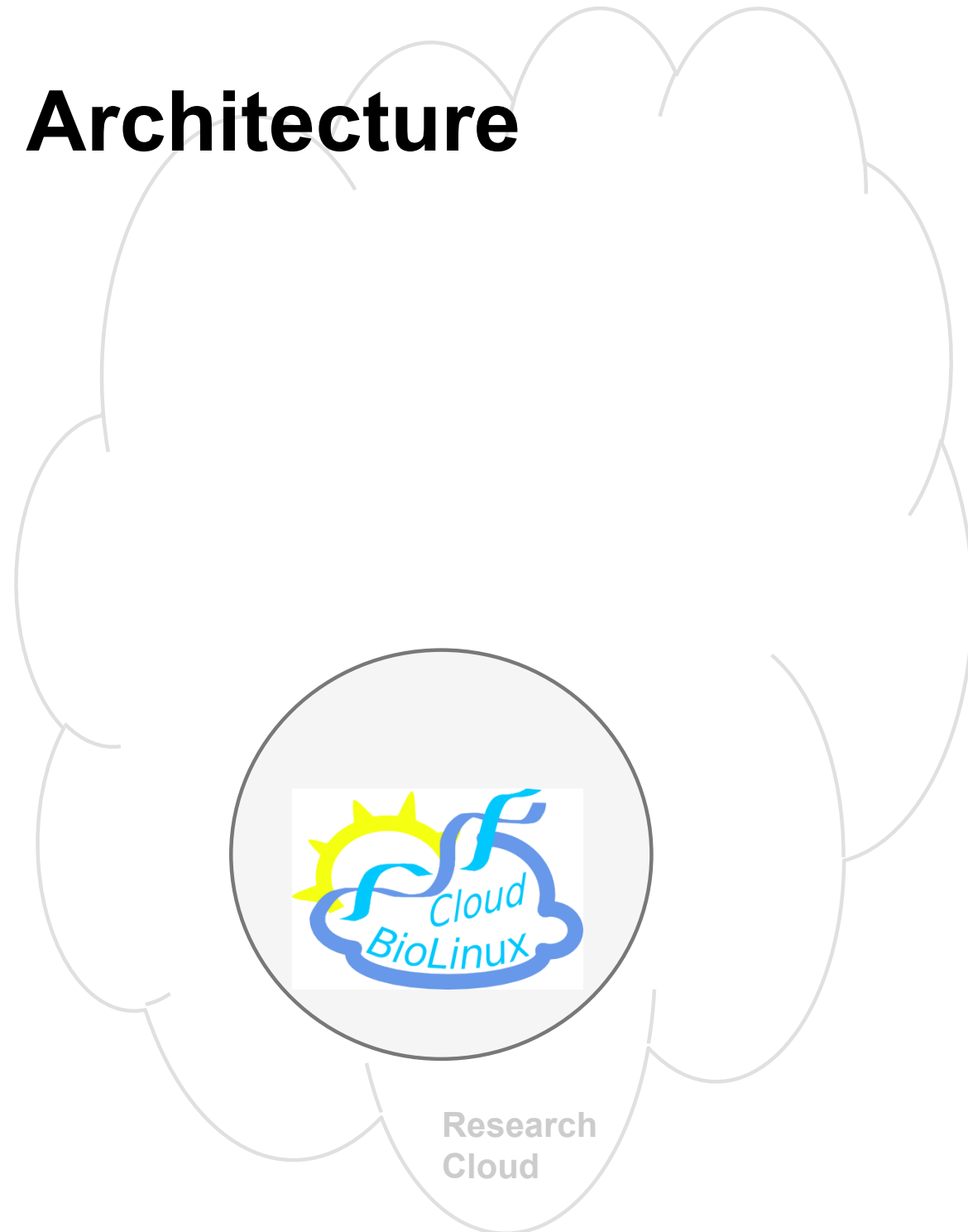
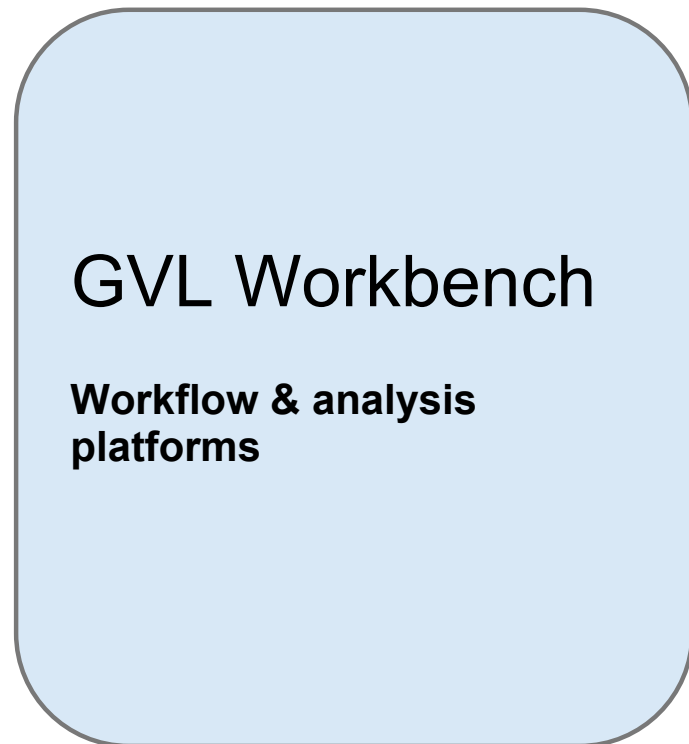
GET a GVL

<http://genome.edu.au> → GET

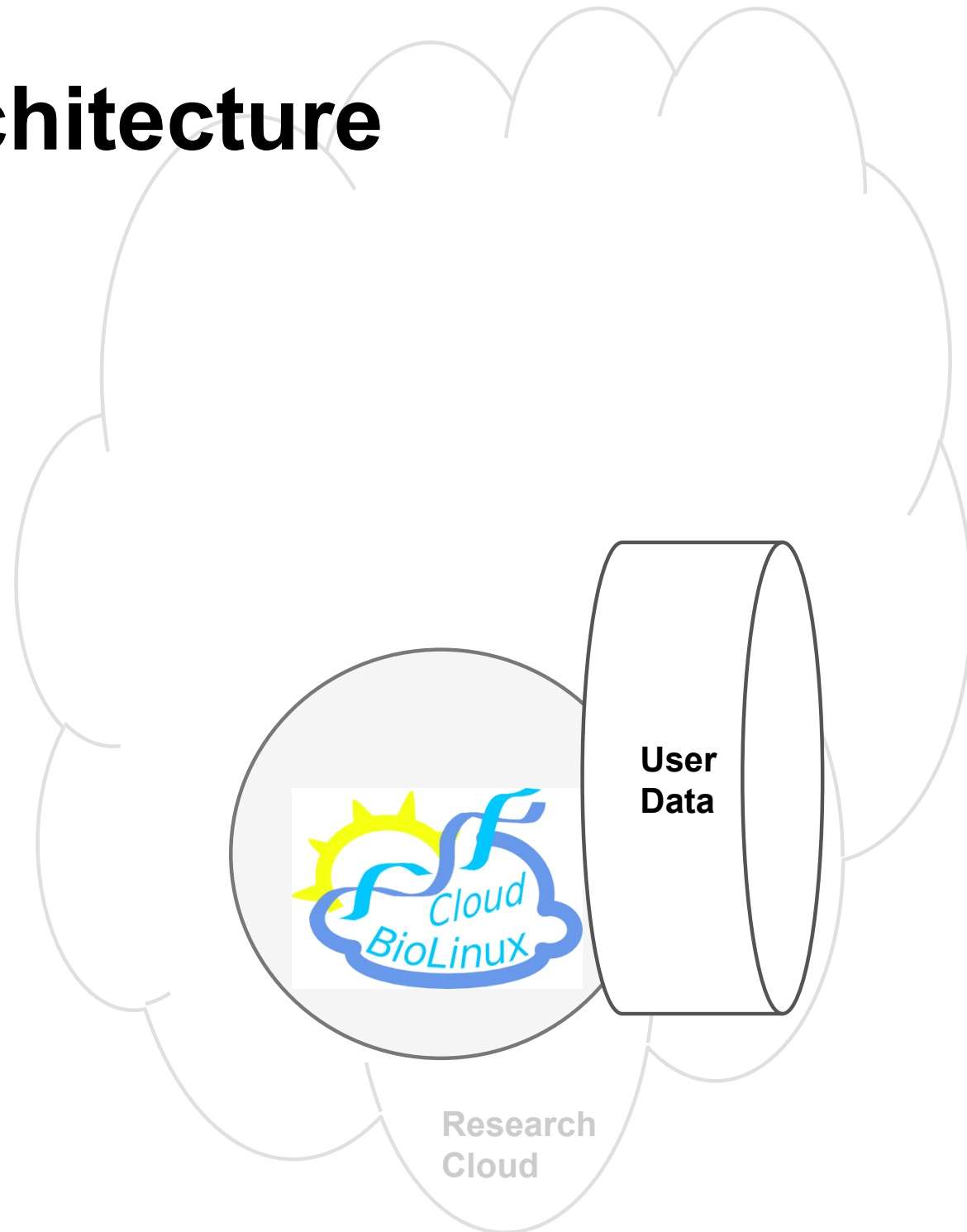
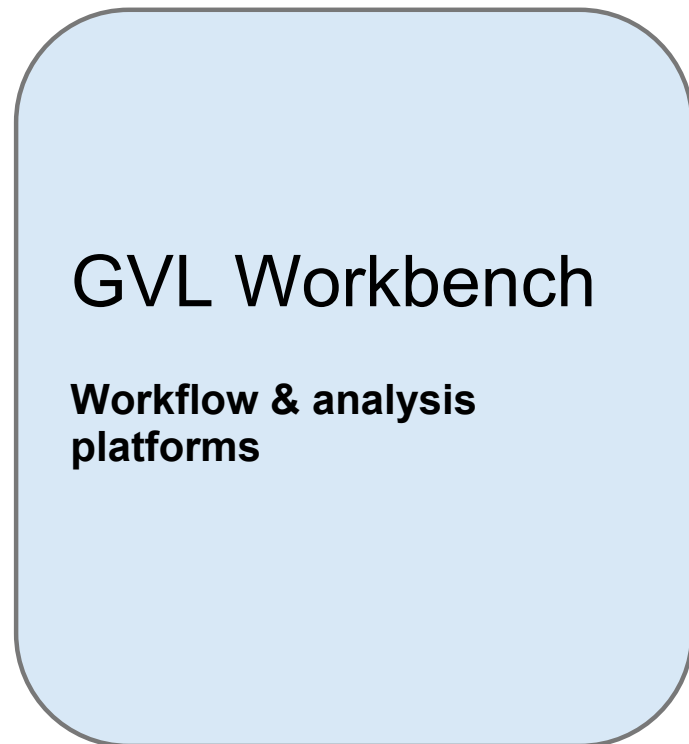
Building (deploying and running) a GVL instance:

- 1. Create a CloudBioLinux server VM*
- 2. Download and install a preconfigured Galaxy*
- 3. Attach pre-populated indexed genomes data*
- 4. Start Galaxy*
- 5. Add extra compute nodes as required*

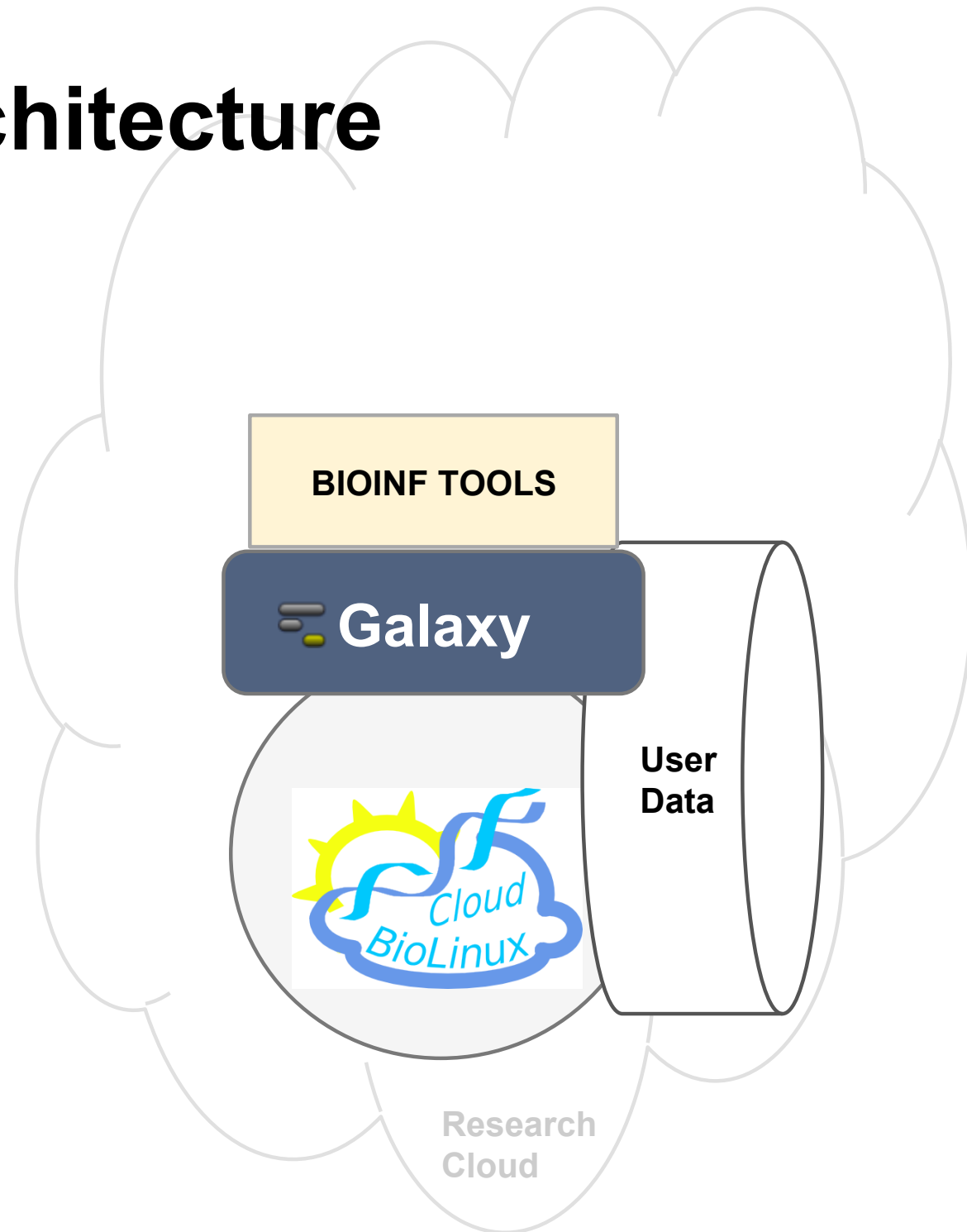
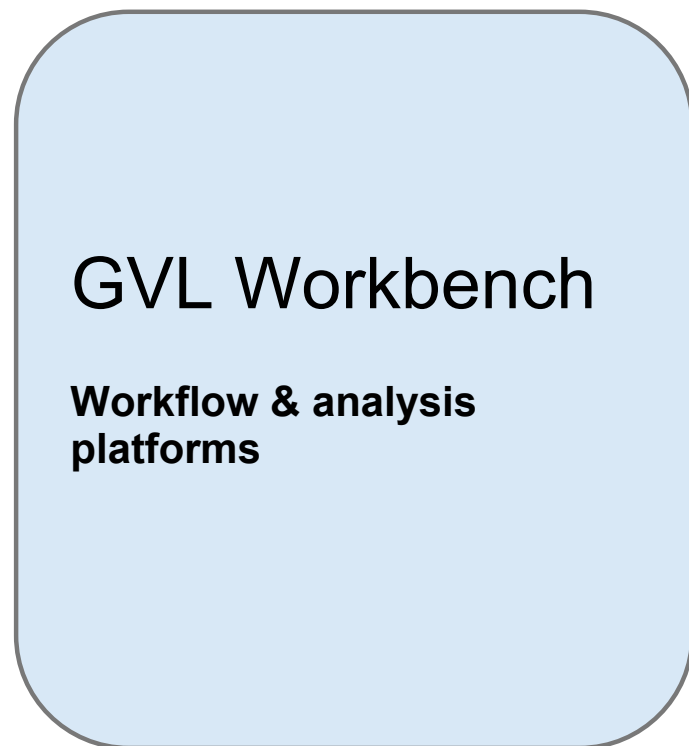
GVL Workbench: Architecture



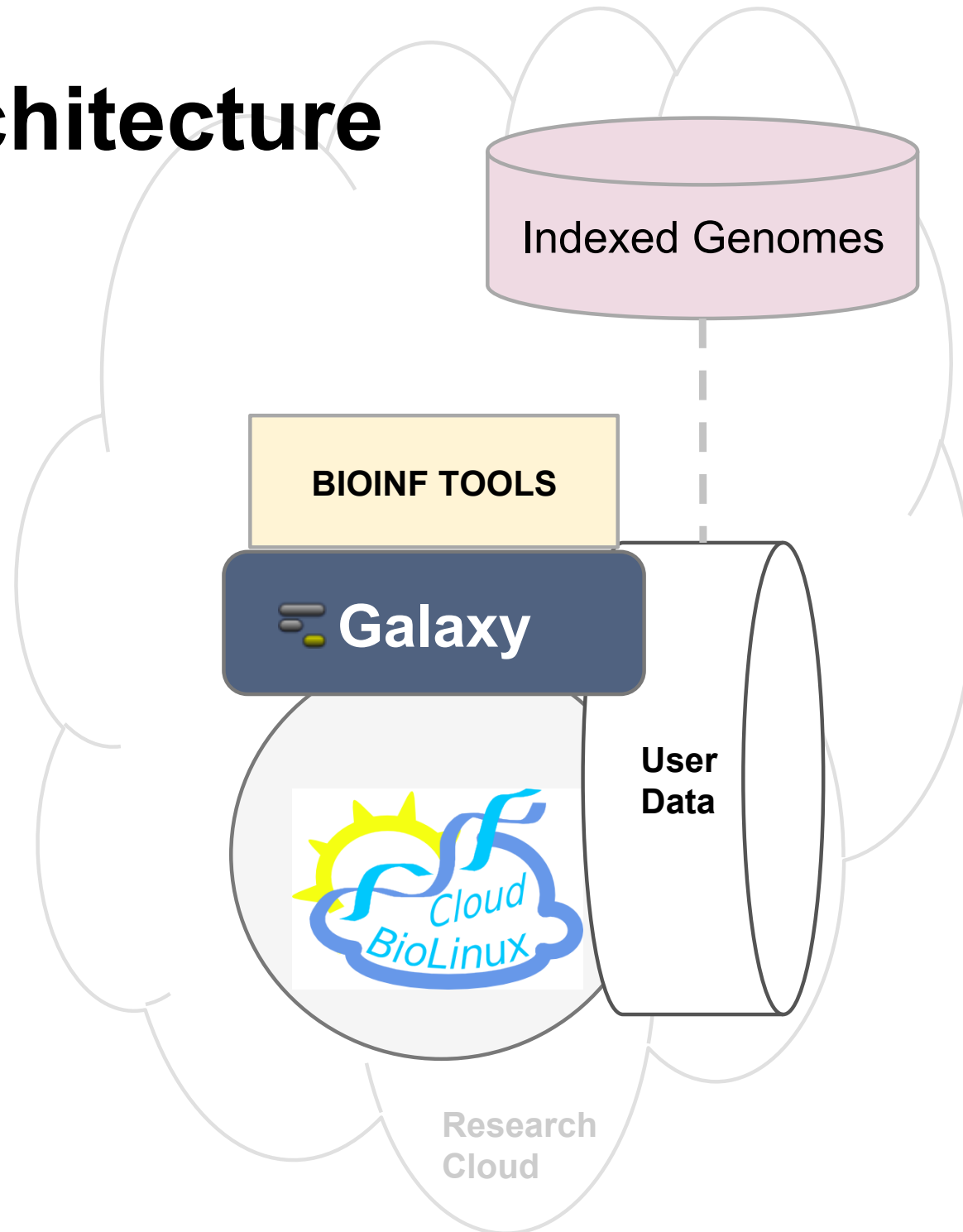
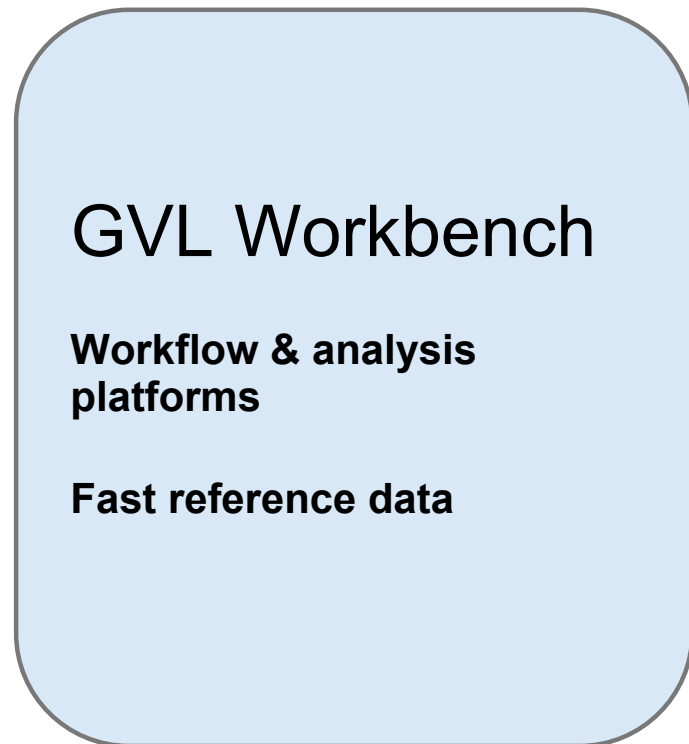
Workbench: Architecture



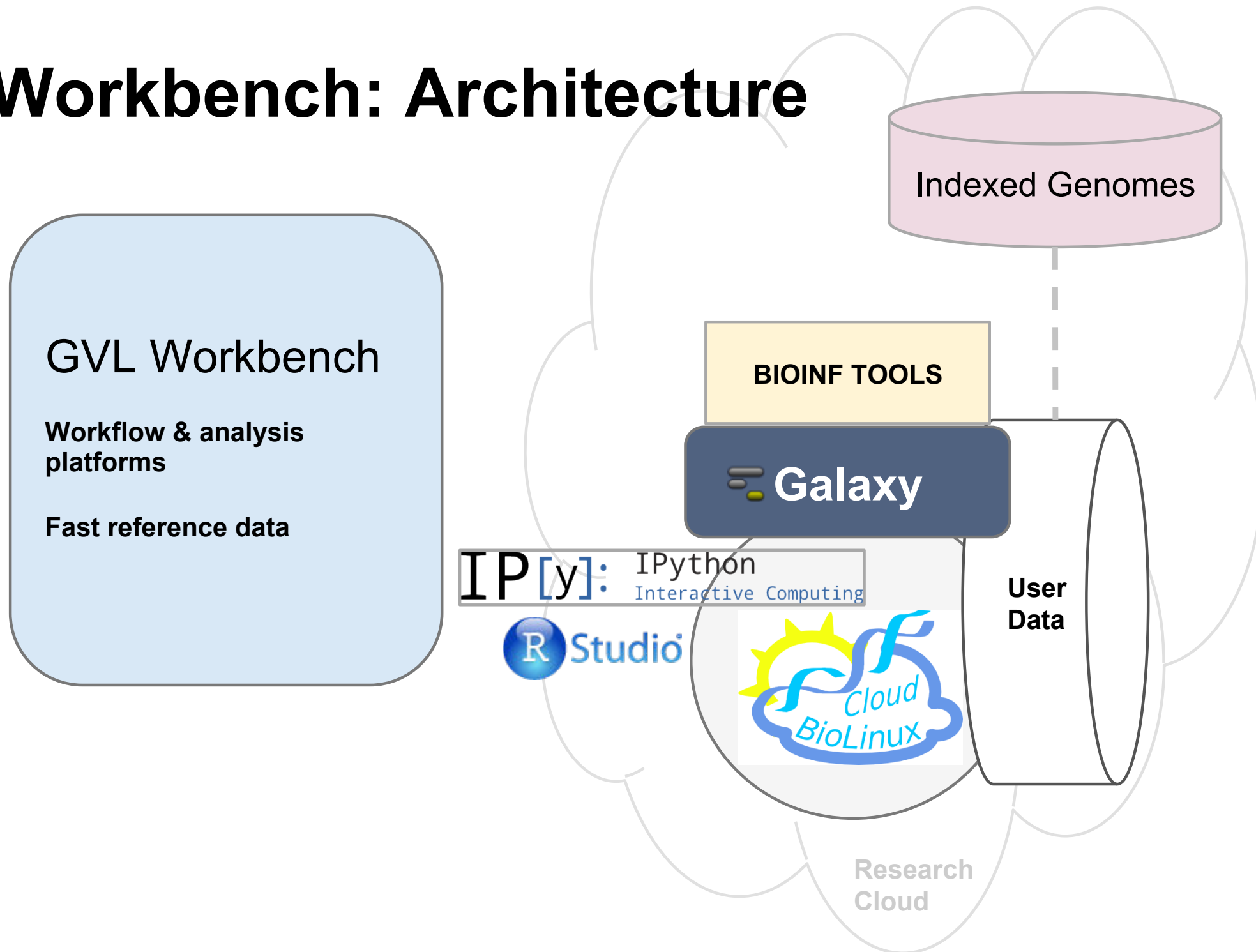
Workbench: Architecture



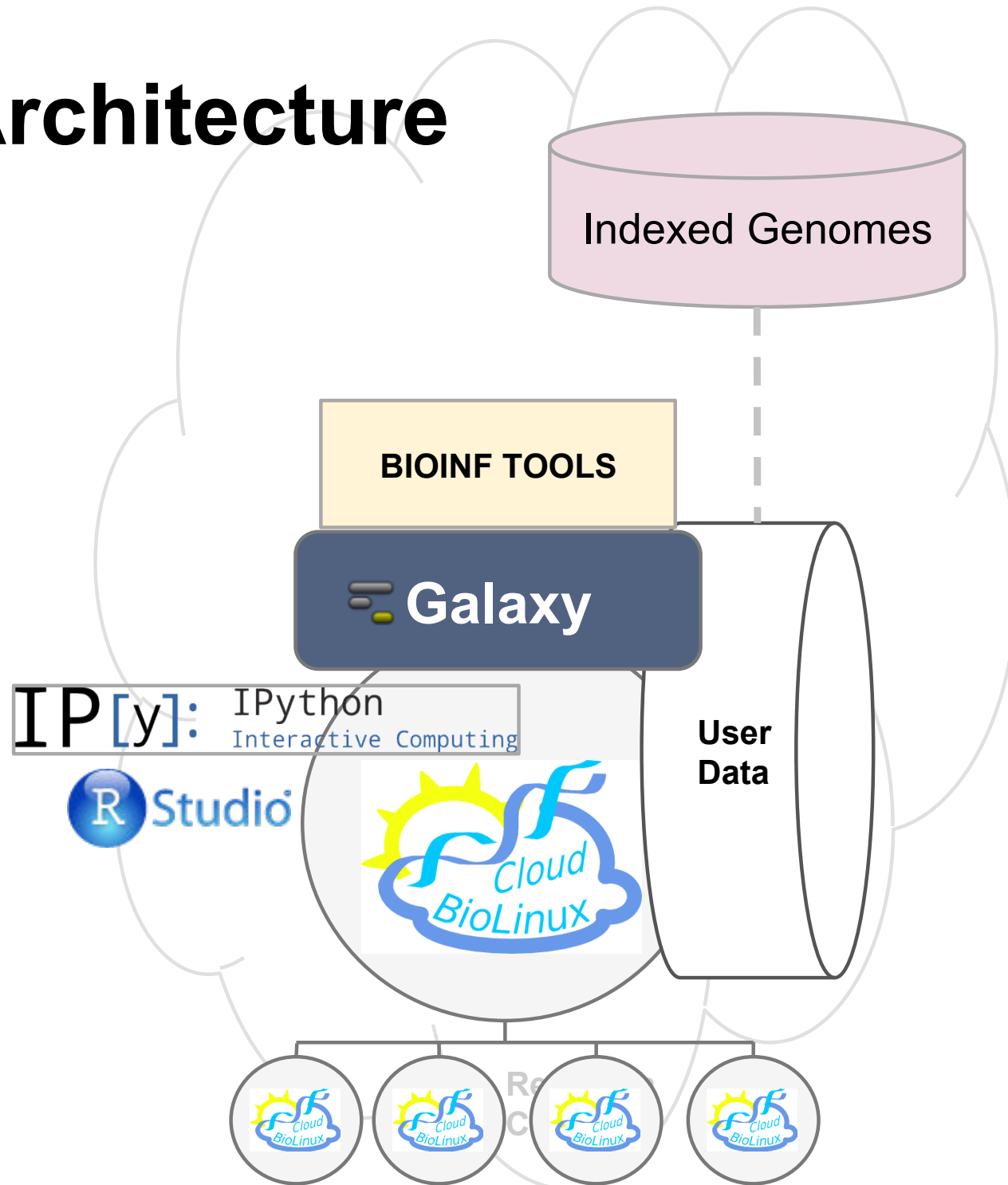
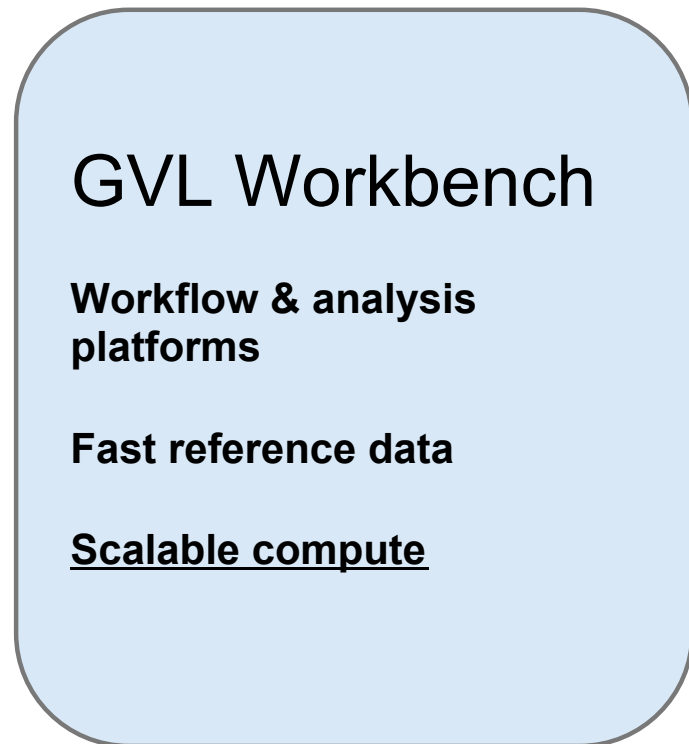
Workbench: Architecture



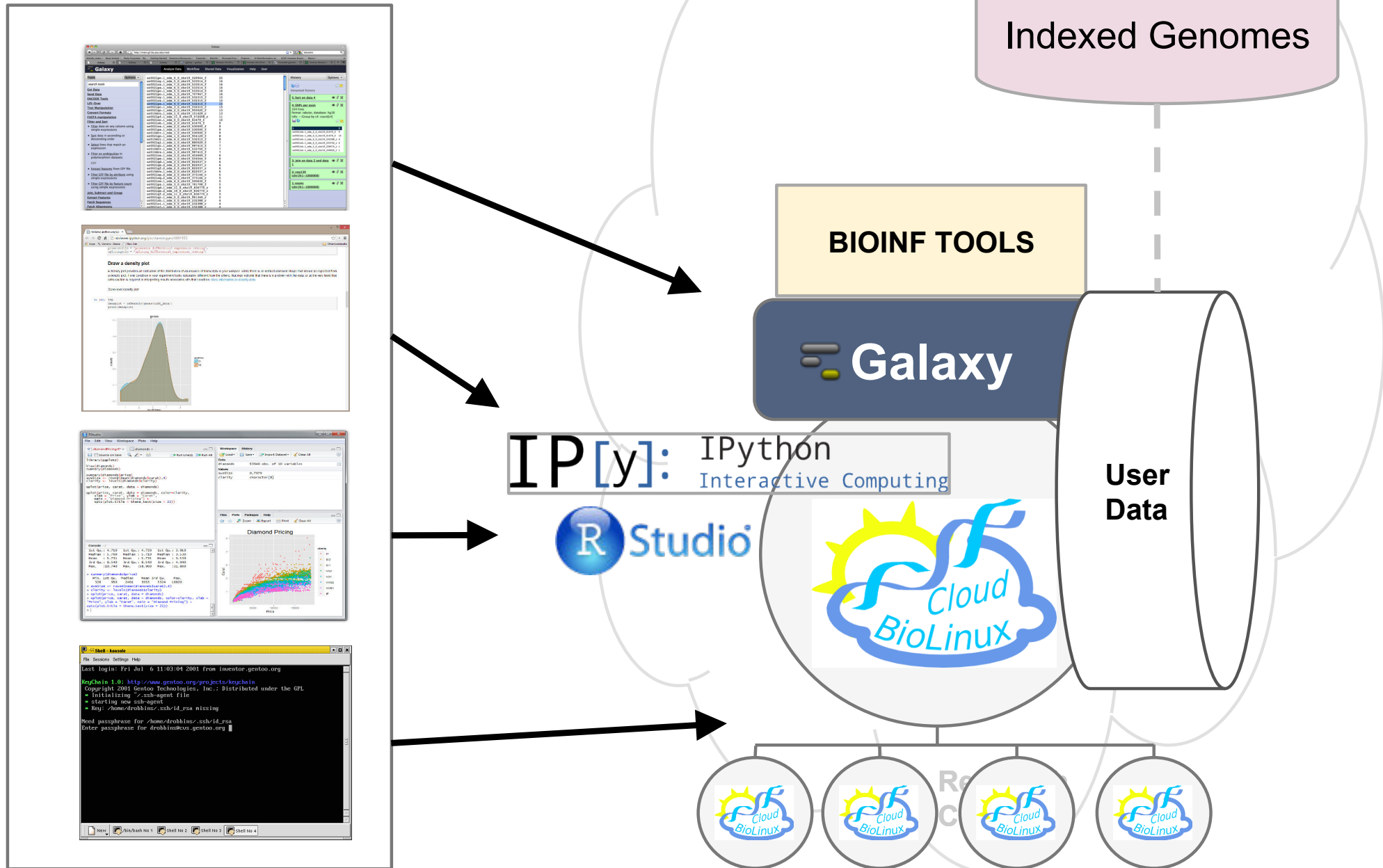
Workbench: Architecture



Workbench: Architecture



Workbench: Architecture



Engineering: Deploying and running a GVL

<http://launch.genome.edu.au>

Cloudman = Middleware for building, distributing and managing cloud-based platforms, especially Galaxy



Afgan et al. *BMC Bioinformatics* 2012, **13**:315
<http://www.biomedcentral.com/1471-2105/13/315>



SOFTWARE

Open Access

CloudMan as a platform for tool, data, and analysis distribution

Enis Afgan^{1,3,4} Brad Chapman² and James Taylor^{3,4}

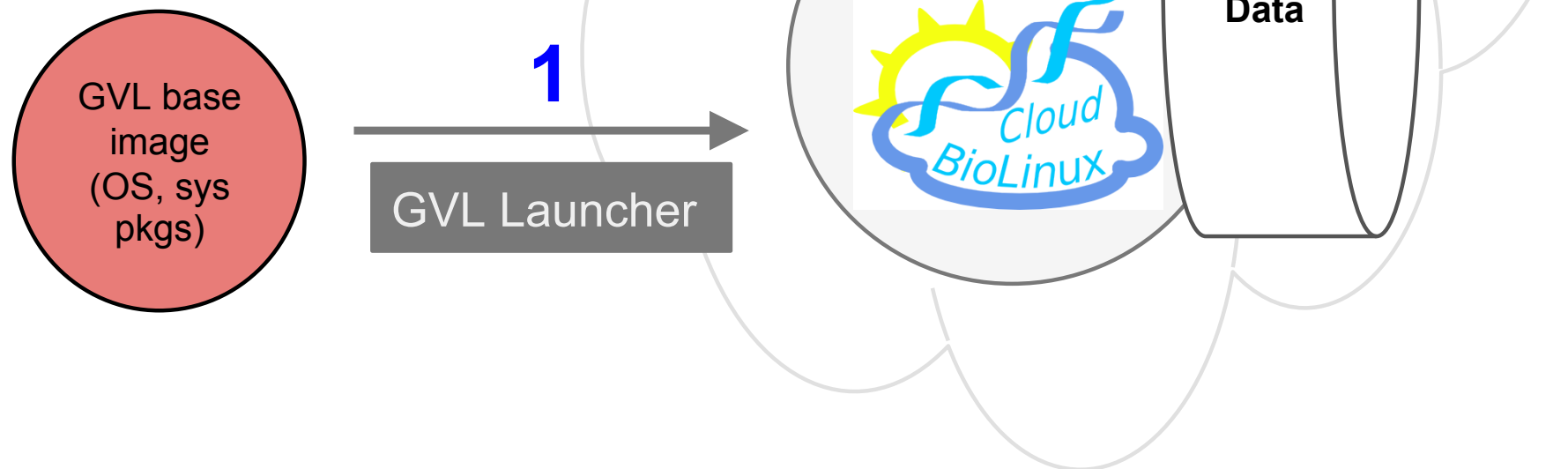
Abstract

Background: Cloud computing provides an infrastructure that facilitates large scale computational analysis in a scalable, democratized fashion. However, in this context it is difficult to ensure sharing of an analysis environment and associated data in a scalable and precisely reproducible way.

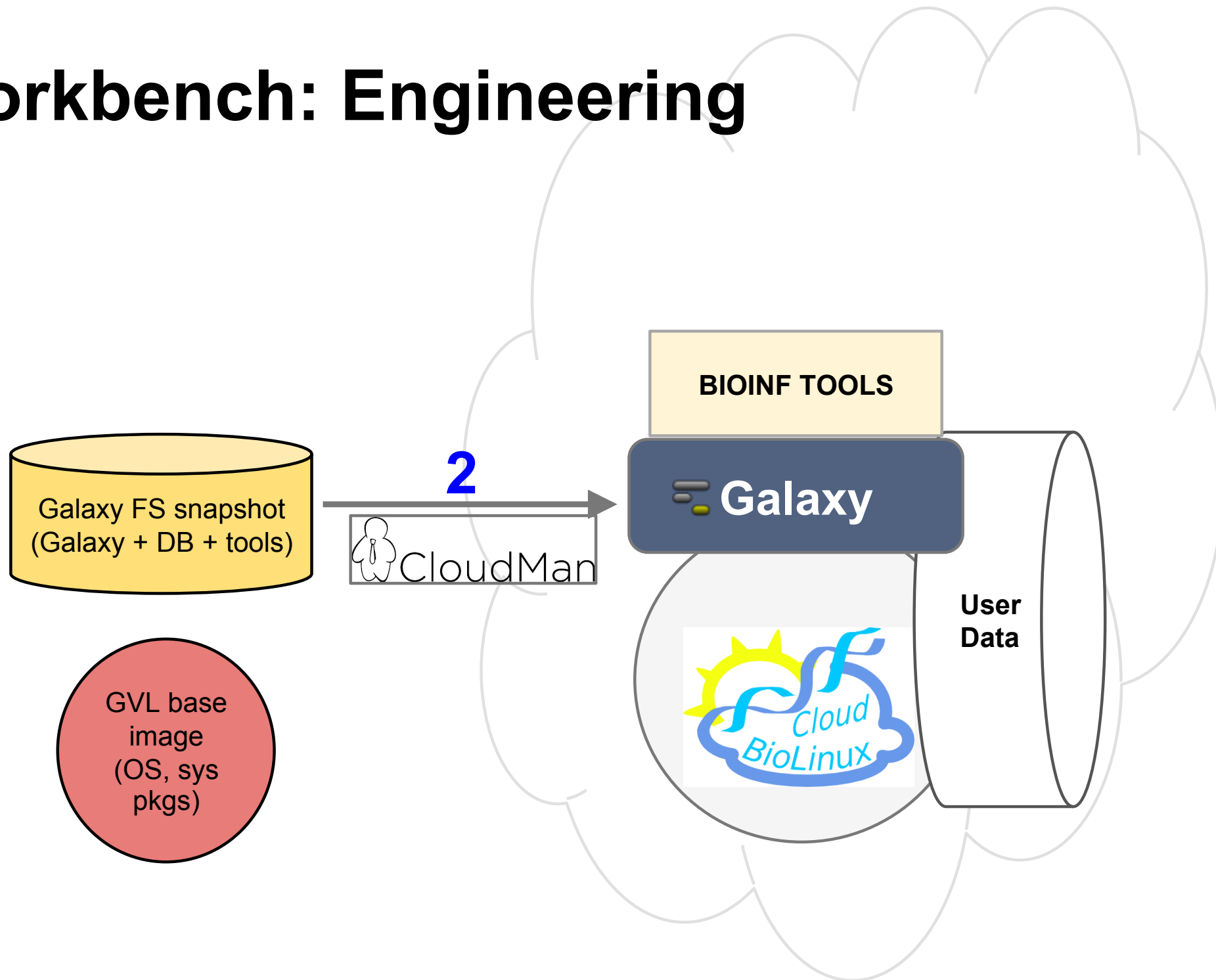
Results: CloudMan (usecloudman.org) enables individual researchers to easily deploy, customize, and share their entire cloud analysis environment, including data, tools, and configurations.

Conclusions: With the enabled customization and sharing of instances, CloudMan can be used as a platform for collaboration. The presented solution improves accessibility of cloud resources, tools, and data to the level of an

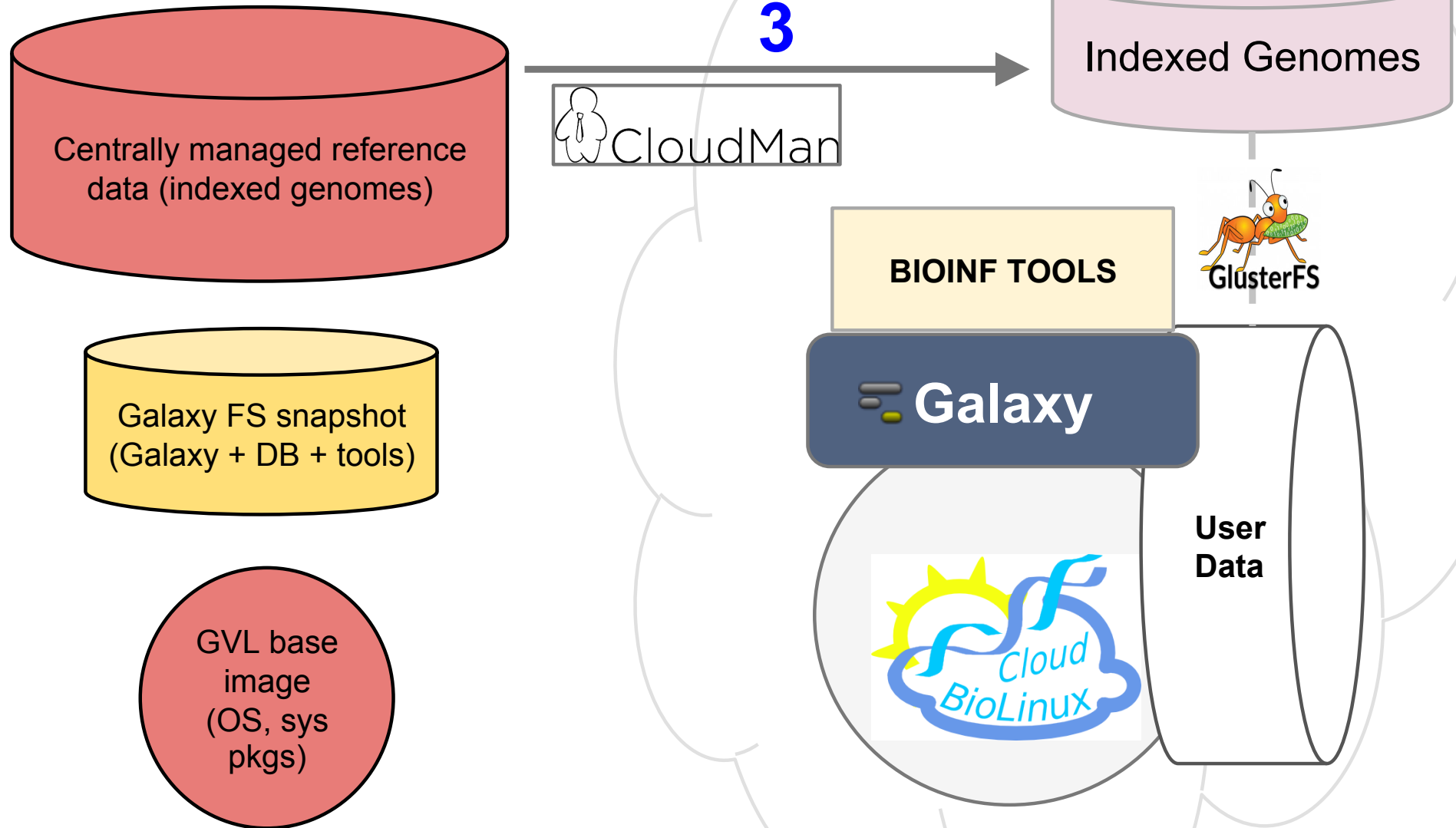
Workbench: Engineering



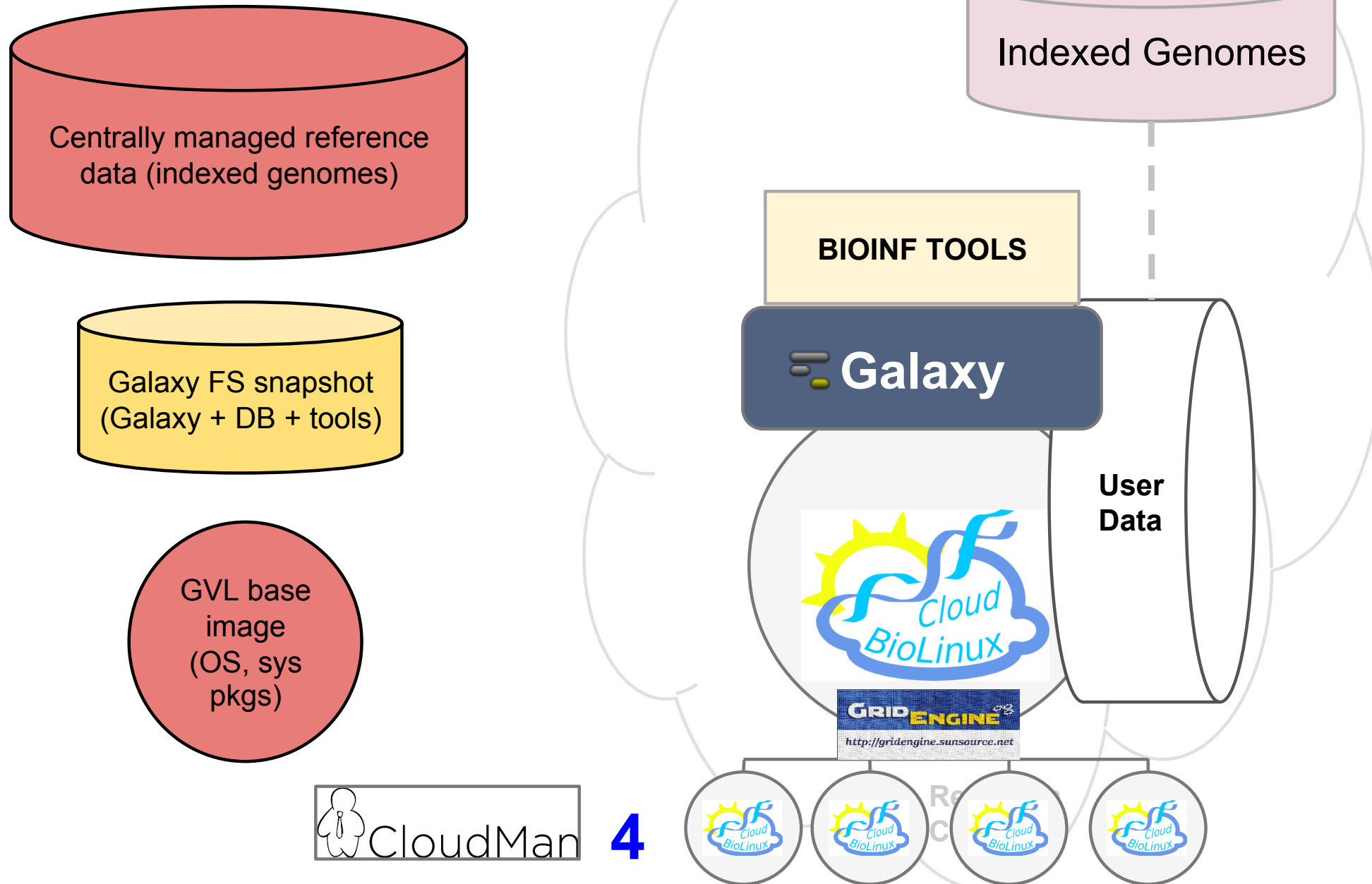
Workbench: Engineering



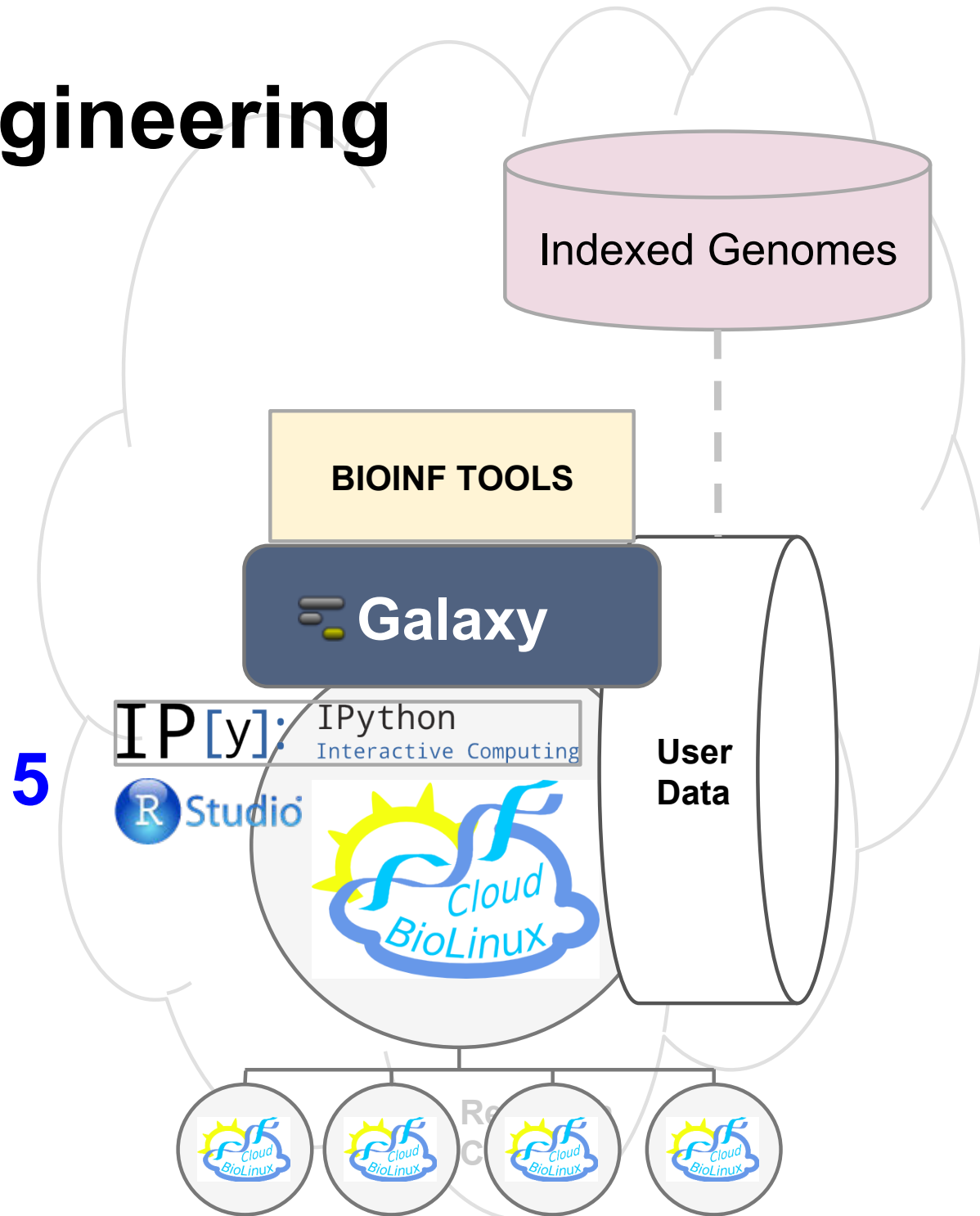
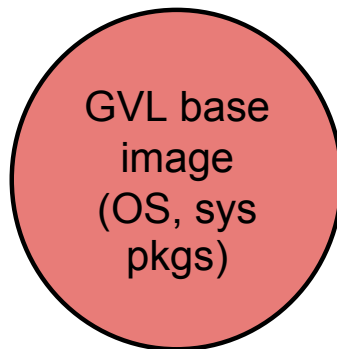
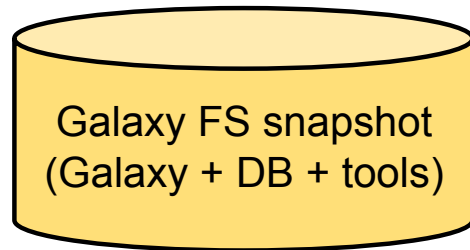
Workbench: Engineering



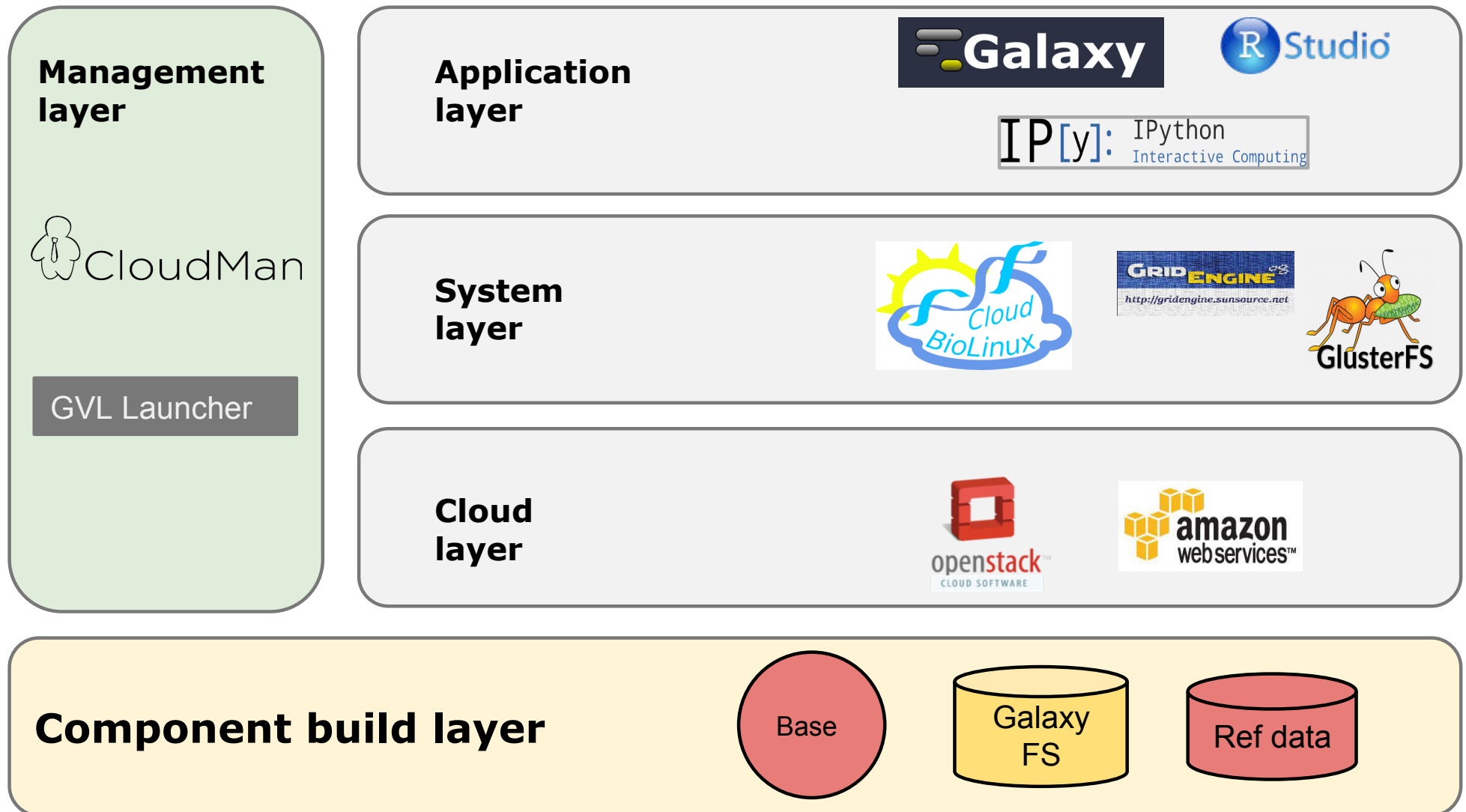
Workbench: Engineering



Workbench: Engineering



Workbench: All components



GVL: Does it work?

Technically?

Practically?

<http://genome.edu.au> → GET



	Personal GVL	Server GVL	Cluster GVL
<i>Suitable for</i>	Single user	Single user Small group/lab	Large groups Institutions
<i>Storage</i>	60GB	100-5000GB	TBs
<i>Compute</i>	2 cores	8-64* cores	>50 cores
<i>Requires</i>	NeCTAR account	NeCTAR allocation: Compute and Volume storage	Large NeCTAR allocation of compute + user-provided fast storage
<i>Runs on</i>	Any Research Cloud node	RC nodes with volumes	RC nodes co-located with fast file system
<i>Setup</i>	<u>Automatic via website</u>	<u>Automatic via website</u>	Collaboration with GVL team
<i>Configuration</i>	No configuration required	Some configuration to tune analyses	Dedicated management



Lessons?

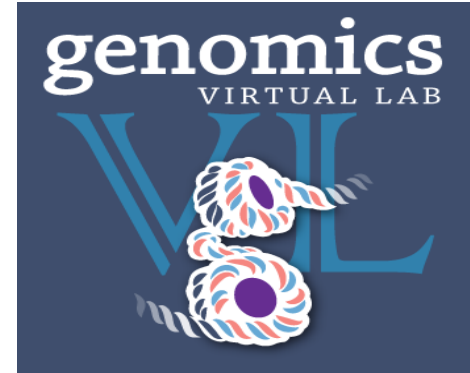
Defining and maintaining a set of tools is challenging

Providing per-user performance is challenging

The cloud is only so scalable!

Not all cloud nodes are equal

Geography matters



Lessons?

Defining and maintaining a set of tools is challenging

Providing per-user performance is challenging

The cloud is only so scalable!

Not all cloud nodes are equal

Geography matters

Resourcing
is key!

What's next for GVL?

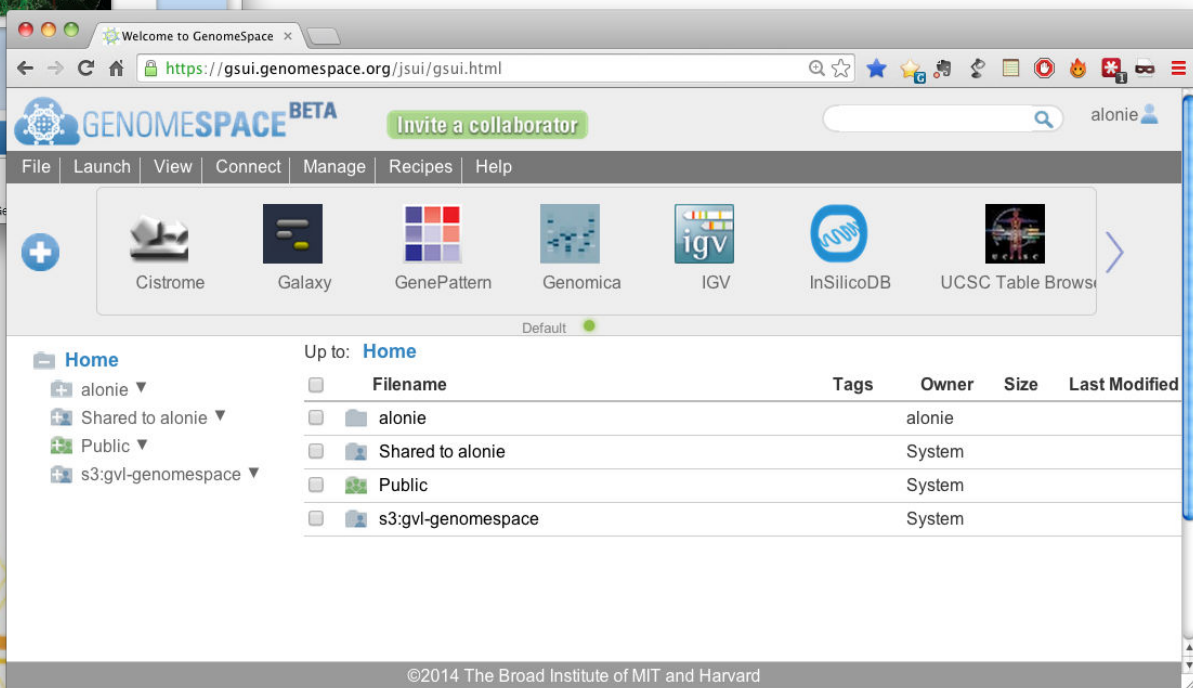
<http://genome.edu.au>



Moving data around is a problem

Whole genomes: 300GB raw data

We need to remove the desktop and USB sticks from the process!



Making the GVL possible

Go8 Universities

- [The University of Queensland](#)
- [The University of Melbourne](#)
- [Monash University](#)
- [The University of Sydney](#)
- [The University of Western Australia](#)

Medical Research Institutes

- [The Garvan Institute of Medical Research](#)
- [Victor Chang Cardiac Research Institute](#)
- [Baker IDI Heart and Diabetes Institute](#)
- [Peter MacCallum Cancer Centre](#)

eResearch Agencies

- [Queensland Facility for Advanced Bioinformatics](#) (QFAB)
- [Queensland Cyber Infrastructure Foundation](#) (QCIF)
- [Life Sciences Computation Centre](#) (LSCC) [at the VLSCI](#)
- [Victorian eResearch Strategic Initiative](#) (VeRSI)

National Agencies

- [NeCTAR, DIIS RTE](#)
- [CSIRO](#)
- [EMBL Australia](#)
- [Bioplatforms Australia](#) (BPA)
- [Australian Genome Research Facility](#) (AGRF)
- [Australian National Data Service](#) (ANDS)