**Basic Course on Bioinformatics tools for Next Generation Sequencing data mining**

# IT infrastructure and user interface: The Galaxy architecture and ARIES cluster

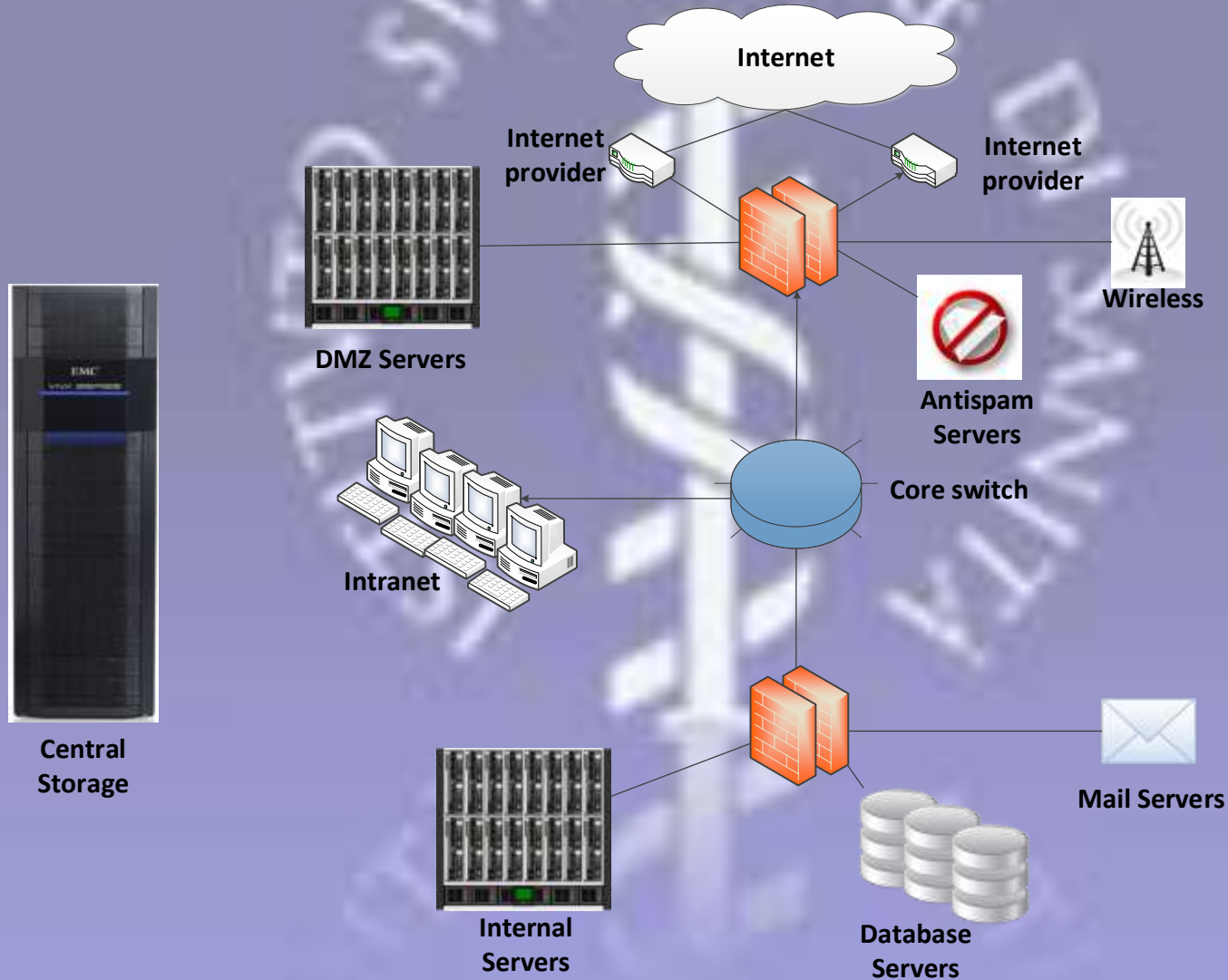Arnold Knijn

IT Sector - ISS

IT infrastructure and user interface: The Galaxy architecture and ARIES cluster

# Istituto Superiore di Sanità

- Personel: ~ 2400
- Wired end-points: ~ 3000
- LAN: 1
- Buildings: 37
- Core switches: 2
- Distribution switches: 87
- Appliances: ~ 20 (firewalls, IPS, etc.)
- 2 Blade systems: 16 hosts (16 logical CPUs, 32/36 GB RAM)
- Servers/virtual machines: > 130 (60% Windows, 40% Linux)
- Databases: > 100 (~ 450 GB)
- Mailboxes: > 3500 (> 6 TB)
- Central Storage: > 50 TB high-level, > 75 TB low-level

IT infrastructure and user interface: The Galaxy architecture and ARIES cluster

# ISS infrastructure

IT infrastructure and user interface: The Galaxy architecture and ARIES cluster

# Data Center Paradigm Evolution

- Mainframe

- One-application, one-server model

IT infrastructure and user interface: The
Galaxy architecture and ARIES cluster

# Standardisation/consolidation

Simplify, through reduction of system number and types.

IT infrastructure and user interface: The Galaxy architecture and ARIES cluster

# Simplification



Storage

CPU/RAM

Networking

IT infrastructure and user interface: The
Galaxy architecture and ARIES cluster

# Server virtualisation

IT infrastructure and user interface: The Galaxy architecture and ARIES cluster

# Virtualisation characteristics

✓ Physical server reduction (1:10 – 1:4)
✓ Decoupling from hardware
✓ Simpler management (installation, backup)
✓ Resource optimisation
✓ Adaptability/ Scalability
✓ Availability
✓ Test environment

- Some hypervisor overhead
- Non-virtualisable hardware (server non x86, etc.)
- More restrictive hardware requirements
- The infrastructure has to be solid
- Virtual machine proliferation

IT infrastructure and user interface: The
Galaxy architecture and ARIES cluster

# Centralised management

IT infrastructure and user interface: The
Galaxy architecture and ARIES cluster

# Galaxy architecture



standalone

web
server

File
server

Galaxy
cluster

Database
server

IT infrastructure and user interface: The
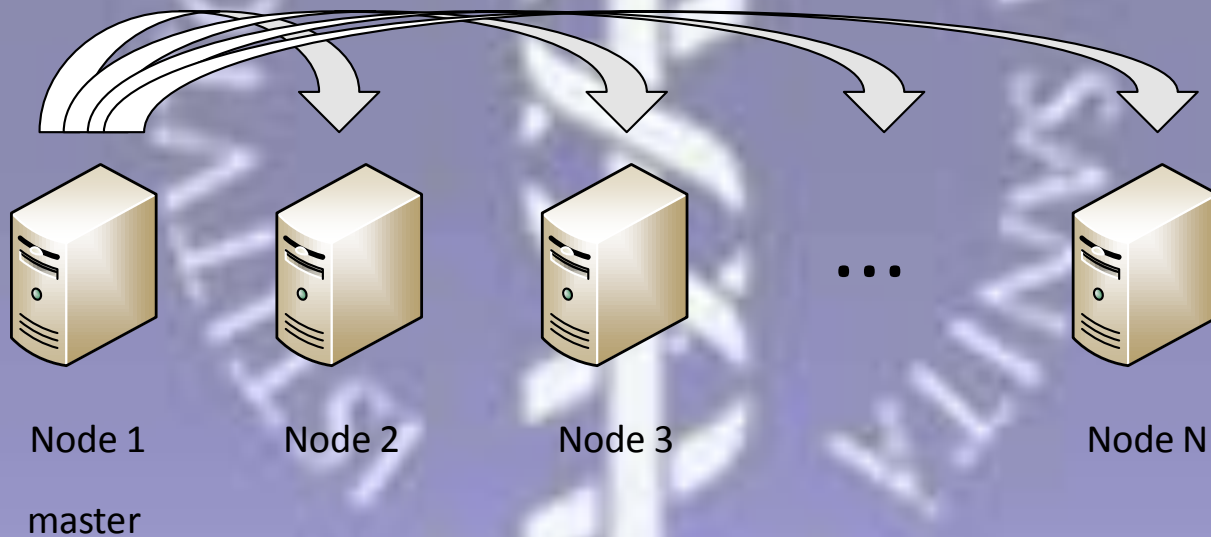Galaxy architecture and ARIES cluster

# Differences

By default, Galaxy uses:

- SQLite
- Built-in HTTP server for all tasks
- Local job runner
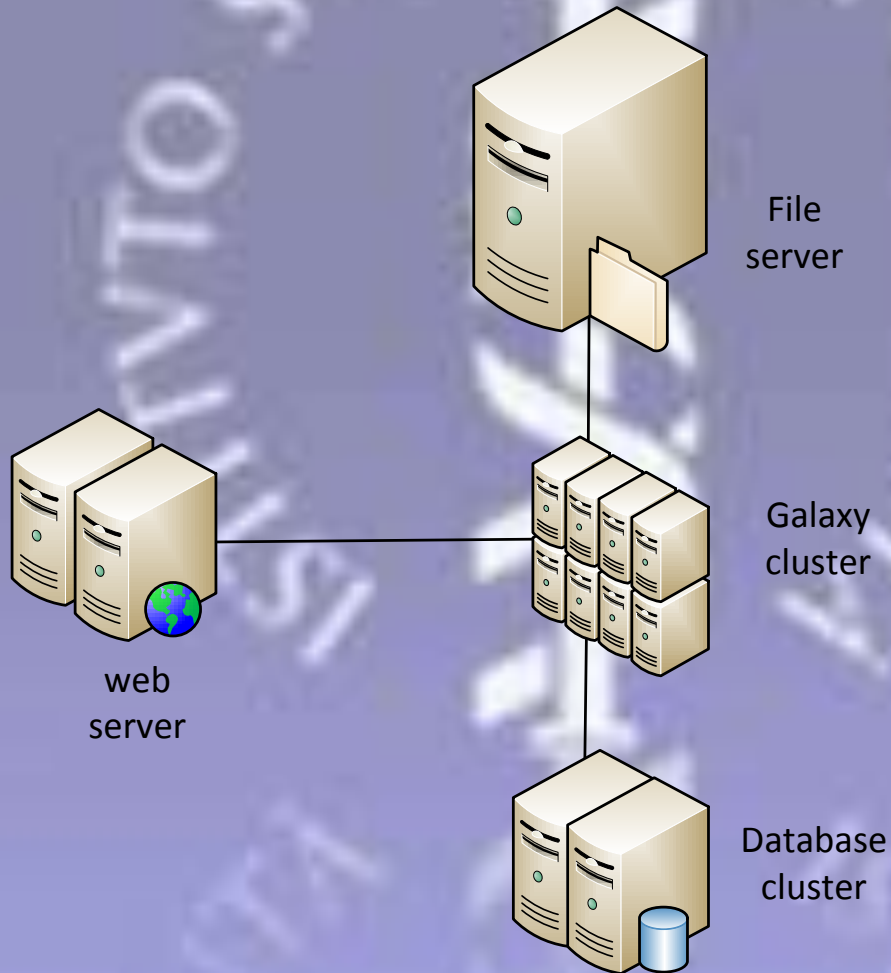- Single process
- Simplest error-proof configuration

In production:

- Real database
- Real HTTP server for many tasks
- Cluster job runner
- Multi process
- More complex configuration

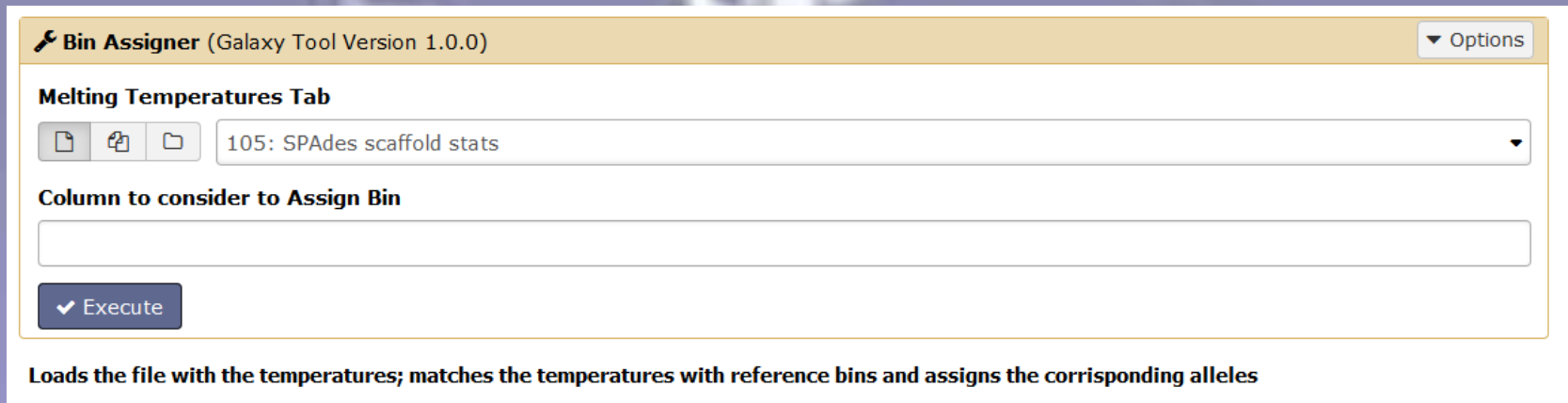IT infrastructure and user interface: The Galaxy architecture and ARIES cluster

# Galaxy cluster



Node 1

Node 2

Node 3

...

Node N

master

IT infrastructure and user interface: The
Galaxy architecture and ARIES cluster

# Architecture scalability



File server

web server

Galaxy cluster

Database cluster

IT infrastructure and user interface: The Galaxy architecture and ARIES cluster

# Galaxy user interface

IT infrastructure and user interface: The
Galaxy architecture and ARIES cluster

# Intuitive and self-documenting

# UI vs Command-Line



**Tool: Bin Assigner**

| | |
|---|---|
| Name: | BinAssigner Log File |
| Created: | Fri Feb 13 07:43:59 2015 (UTC) |
| Filesize: | 877 bytes |
| Dbkey: | ? |
| Format: | txt |
| Galaxy Tool ID: | binassigner |
| Galaxy Tool Version: | 1.0.0 |
| Tool Version: | |
| Tool Standard Output: | stdout |
| Tool Standard Error: | stderr |
| Tool Exit Code: | 0 |
| API ID: | e9fb797960230e8a |
| History ID: | f597429621d6eb2b |
| UUID: | dc1676ef-87b7-48bf-a24e-4359f57cf2fa |
| Full Path: | /home/galaxy/galaxy-dist/database/files/001/dataset_1439.dat |
| Job Command-Line | python /home/galaxy/galaxy-dist/tools/Hrevap/BinAssigner.py -t /home/galaxy/galaxy-dist/database/files/001/dataset_1433.dat -o /home/galaxy/galaxy-dist/database/files/001/dataset_1438.dat -c 7 > /home/galaxy/galaxy-dist/database/files/001/dataset_1439.dat |
| Job Runtime (Wall Clock) | 1 seconds |
| Cores Allocated | 1 |
| Job Start Time | 2015-02-13 08:44:00 |
| Job End Time | 2015-02-13 08:44:01 |

| Input Parameter | Value | Note for rerun |
|---|---|---|
| Melting Temperatures Tab | 176: TermoTyping Summary File | |
| Column to consider to Assign Bin | 7 | |

IT infrastructure and user interface: The
Galaxy architecture and ARIES cluster

# Home –made tools

```xml
<tool id="binassigner" name="Bin Assigner">
    <description>Bin Assigner tool</description>
    <command interpreter="python">
        BinAssigner.py -t $tmstab -o $output -c $columntab > $logfile
    </command>
    <inputs>
        <param name="tmstab" type="data" format="tabular" label="Melting Temperatures Tab"/>
        <param name="columntab" type="text" format="integer" label="Column to consider" />
    </inputs>
    <outputs>
        <data format="tabular" name="output" label="Allele Table"/>
        <data format="txt" name="logfile" label="BinAssigner Log File" />
    </outputs>
    <help>
        **Loads the file with the temperatures; matches the temperatures with reference bins
            and assigns the corrisponding alleles**
    </help>
</tool>
```

IT infrastructure and user interface: The
Galaxy architecture and ARIES cluster