# The Galaxy service pilot in CSIRO

## A collaboration between science and IT

Steve McMahon, Philippe Moncuquet, Sean Li, Ondrej Hlinka, Josh Bowden, Sean McWilliam and Annette McGrath

## CSIRO BIOINFORMATICS CORE, INFORMATION MANAGEMENT & TECHNOLOGY
### www.csiro.au

Galaxy (1) (2) (3) is a web based bioinformatics toolkit with the ability to perform, reproduce and share complete analyses.

A Galaxy service pilot was set up in CSIRO (Australia's Commonwealth Scientific and Industrial Research Organisation) for the benefit of biologists and bioinformaticians. The service pilot was implemented as a collaboration between CSIRO's Information Management and Technology staff (IM&T) and the CSIRO Bioinformatics Core (CBC). This made the best use of the IT infrastructure and service delivery expertise (through IM&T) and the bioinformatics domain expertise (through CBC).
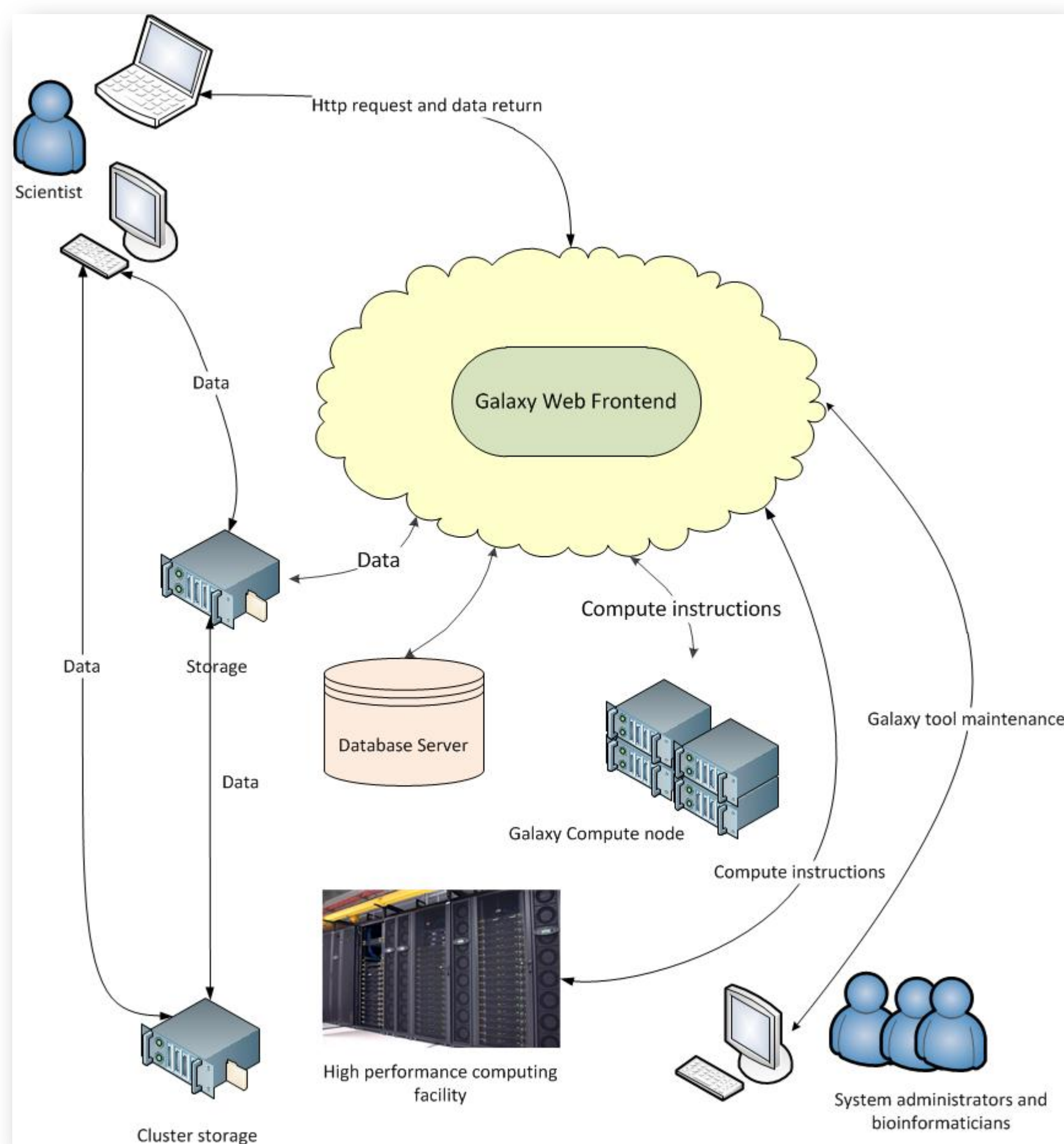
This presentation outlines Galaxy, the way it was implemented in CSIRO as a service pilot and some of the outcomes and related experiences.
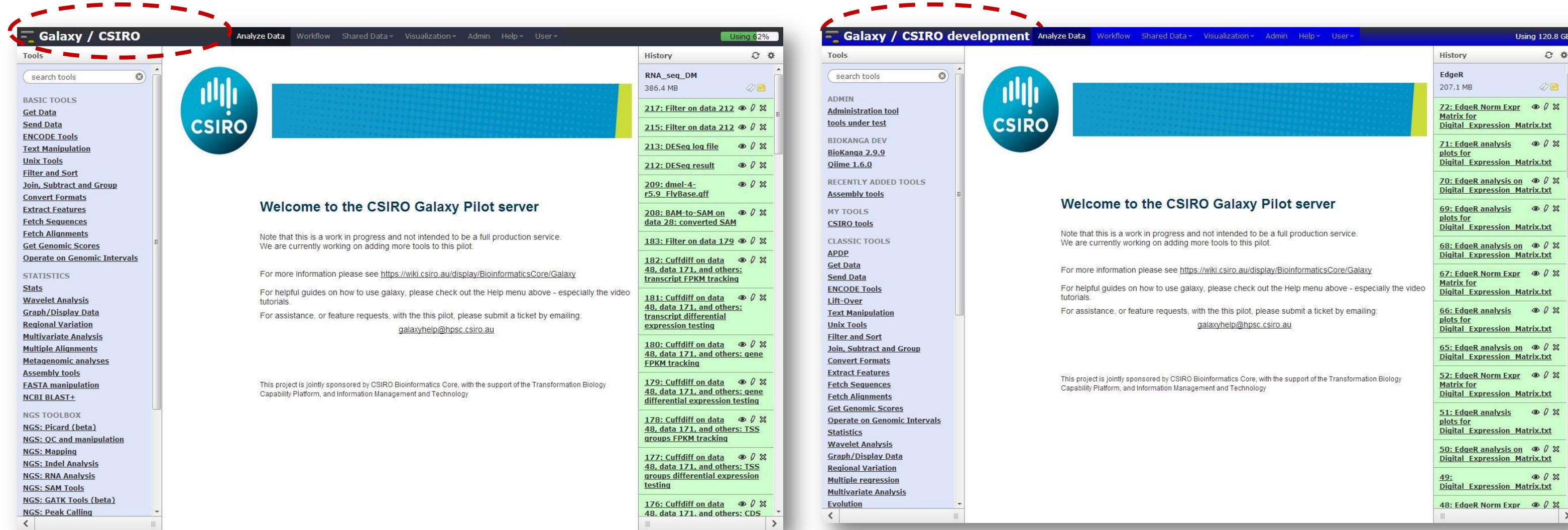
## Problem

As next generation sequencing becomes more affordable, more experiments requiring bioinformatic analysis are designed by biologists. In CSIRO they had been relying on a limited number of skilled bioinformaticians to carry out this analysis. Alternatively they may have chosen to educate themselves in these same skills which is time consuming. Analysis was not always performed on the best computing resources and not always in the most optimal way. It was proposed that a Galaxy instance would **empower biologists** to perform their own data analysis by providing easy access to bioinformatics analysis tools. The success of this project relied on both scalable hardware solutions and bioinformatics knowledge and support.

## Architecture

An effective bioinformatics platform needs to be scalable as computational power and storage needs are variable. The CSIRO Galaxy pilot project was designed to service approximately **100 users**. As a pilot it had to be capable enough to be useful whilst not being too expensive.
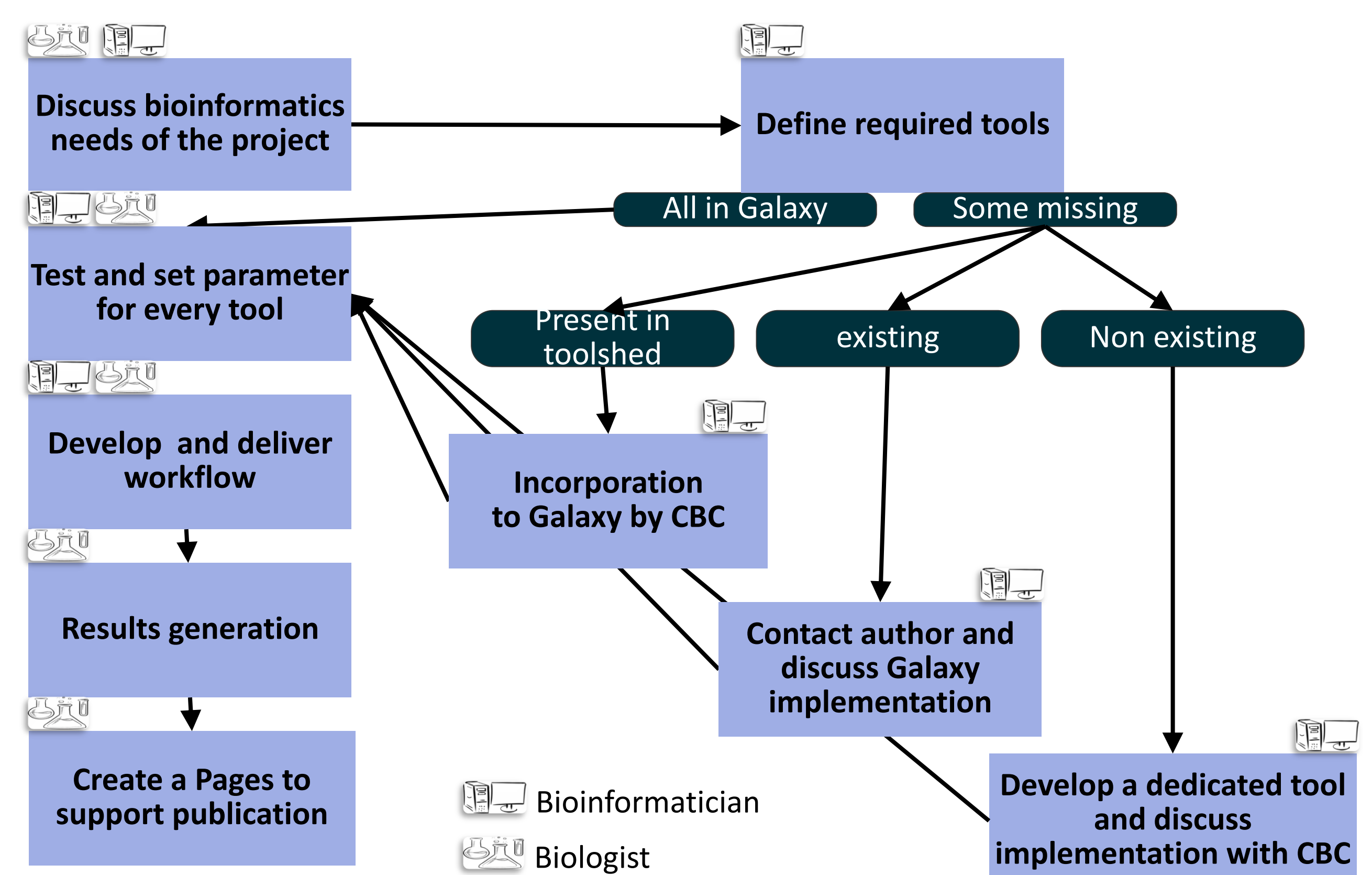


The infrastructure consisted of 2 physical machines – one to host the Galaxy database and the other to run Galaxy and be the general compute engine. There were also **2 virtual machines** running the web front ends for a production instance and a development instance of Galaxy. The compute machine had **32 cores** and **192 GB RAM** and ran the Torque batch system.



As the pilot was being developed the storage allocated to it went from about 1TB to 5TB and then later to 16TB. The available storage seemed to be a key factor for whether some users made use of the pilot.
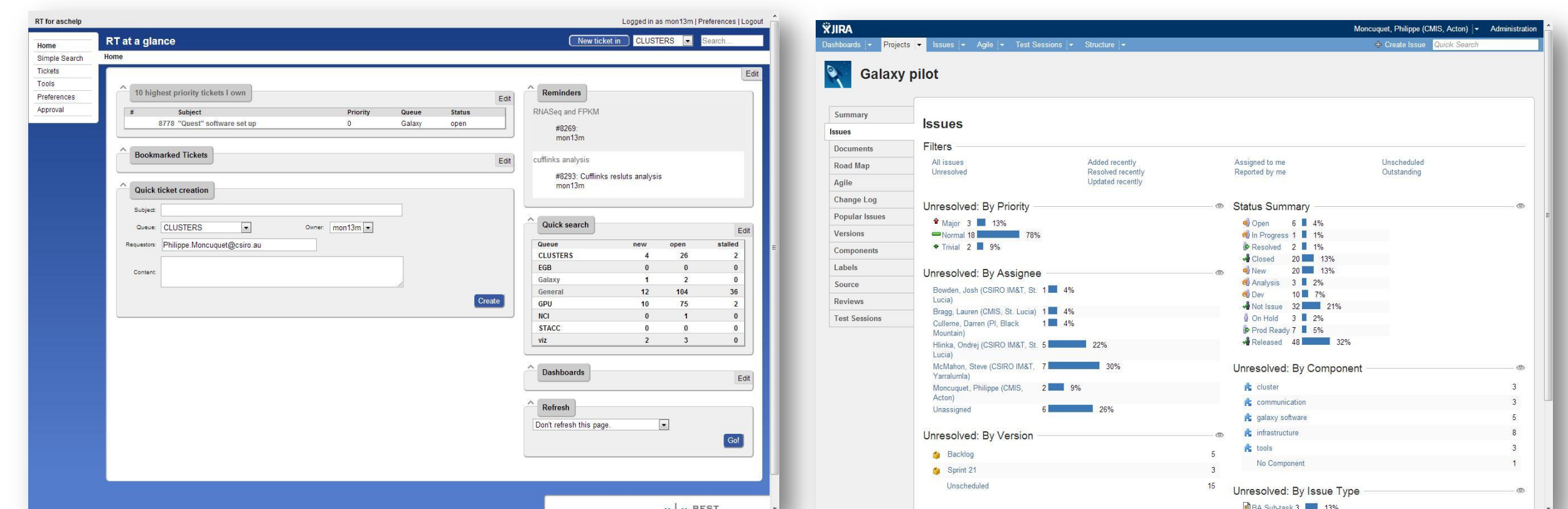
## A typical NGS project

A good collaboration between biologist and bioinformatician is essential for any successful NGS experiment. Galaxy offers the means to design the analysis of an experiment from beginning to end. Galaxy workflows are a great option to **tackle reproducibility issues** and enhance sharing of the scientific work.



## Support, development and user education are key

As for all web-based platforms, development, maintenance, **support and user training** play a central role in our Galaxy production instance. Maintenance and development were dealt with by both bioinformaticians and IT staff. Development followed an "agile" process with issues managed in Jira – web based software suited to this process. User initiated issues were submitted to RT (Request Tracker), another issue tracker. The vast majority of the user support queries were related to the correct choice and use of bioinformatics tools.



User training was carried out via live demos. There are also plans for hands-on tutorials.

## Pilot results

The pilot gathered **113 users**. More than 300 tools have been made available in the platform and over **1000 jobs** are launched each month. The pilot has been used in 1 submitted article so far. 19 workflows are shared with users as well as some useful datasets.

Since the announcement of the pilot last October, more than 60 support tickets have been created and resolved.

The pilot was considered a great success showing how CSIRO IT and science staff could work together to achieve project goals and a full production service is being planned.

**REFERENCES**

1. *Galaxy: a web-based genome analysis tool for experimentalists.* **Blankenberg D, Von Kuster G, Coraor N, Ananda G, Lazarus R, Mangan M, Nekrutenko A, Taylor J.** January 2010, Current Protocols in Molecular Biology, pp. 1-21.
2. *Galaxy: a platform for interactive large-scale genome analysis.* **Giardine B, Riemer C, Hardison RC, Burhans R, Elnitski L, Shah P, Zhang Y, Blankenberg D, Albert I, Taylor J, Miller W, Kent WJ, Nekrutenko A.** 10, October 2005, Genome Research, Vol. 15, pp. 1451-5.
3. *Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences.* **Goecks, J, Nekrutenko, A, Taylor, J and The Galaxy Team.** 8, August 2010, Genome Biol., Vol. 11.