# High level distributed processing pipelines with Galaxy
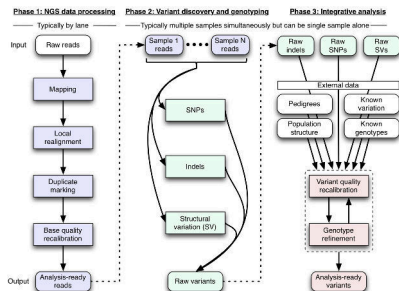
Brad Chapman, Bioinformatics Core at Harvard School of Public Health

Galaxy Community Conference, 27 July 2012

# Complex pipelines

Goal: simple interface to complex analyses

- Steps, lots of them
- Intermediate files
- Branching logic: if/else
- Idempotent
- Transactional
- Parallel: by record, by region
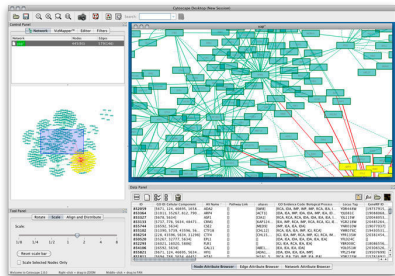- Associated resources
- Experimental metadata



http://www.broadinstitute.org/gsa/wiki/index.php/Best_Practice_
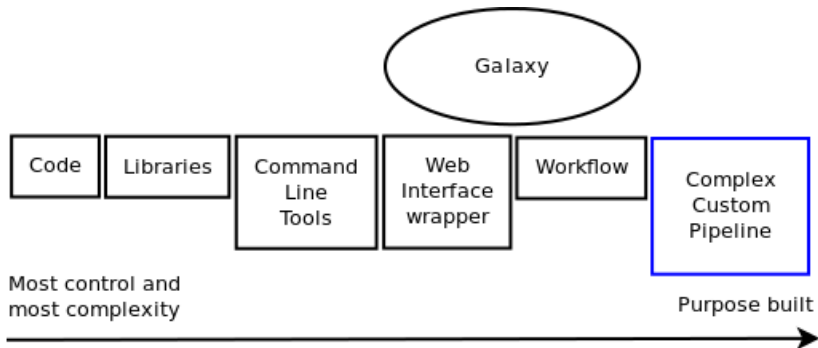Variant_Detection_with_the_GATK_v3

# (At least) two types of users

Galaxy

| Code | Libraries | Command Line Tools | Web Interface wrapper | Workflow | Complex Custom Pipeline |

Most control and
most complexity

Purpose built

**Levels of abstraction**

# Approaches

- Custom Galaxy interfaces
    - Variant calling pipeline
    - Galaxy views and controllers
- Web tools
    - BioCloudCentral
    - CloudMan
- External data upload
    - ISA-Tab experiments
    - Data Libraries and API
- Interoperable tools
    - Variant analysis platform
    - GenomeSpace

# Variant pipeline interface



https://bitbucket.org/hbc/galaxy-central-hbc

# Implementation approach

- Integrated into Galaxy fork
- Custom view (`http://your.galaxy.org/pipeline`)
- Custom controller
- Calls external server for processing
- Results upload to Galaxy Data Libraries

Demo with movies: `http://j.mp/uNXZY6`

# Backend variant processing

# BioCloudCentral

Easily launch CloudMan, CloudBioLinux and Galaxy platforms on Cloud Computing resources (including Amazon Web Services).

| Cluster name | | Name of your cluster used for identification. This can be any name you choose. |
| --- | --- | --- |
| Password | | Your choice of password, for the CloudMan web interface and accessing the instance via ssh or FreeNX. |
| Cloud | Amazon (AWS EC2) | Choose from the available clouds. The credentials you provide below must match (ie, exist on) the chosen cloud. |
| Access key | | Your Access Key ID. For the Amazon cloud, available from the security credentials page. |
| Secret key | | Your Secret Access Key. For the Amazon cloud, also available from the security credentials page. |
| Instance type | Large (4 ECUs / 7.5GB RAM) | Type (ie, virtual hardware configuration) of the instance to start. |

Show advanced startup options

Start an instance

http://biocloudcentral.org

# CloudMan: automate setup



http://usecloudman.org/

# Run analysis

# Experimental metadata: ISA-Tab



http://isatab.sourceforge.net/

# Data Library integration



- Stem Cell Discovery Engine:
  http://discovery.hsci.harvard.edu/
- https://github.com/hbc/projects/blob/master/scde_deploy/scripts/bii_datasets_to_galaxy.py

# Interoperable tools



- [https://github.com/chapmanb/bcbio.variation](https://github.com/chapmanb/bcbio.variation)
- [http://validationprotocol.org](http://validationprotocol.org)
- [http://genomespace.org](http://genomespace.org)

## Answers and Questions

- Build off core Galaxy features
    - Data Libraries
    - API
    - CloudMan
- Maximize interoperability
- Automate aggressively

Would love to hear your experiences.