

Reproducible and automated processing in high throughput NGS facilities

Gianmauro Cuccuru, Giorgio Fotia, Josh Moore, Luca Lianas, Luca Pireddu, Jason Swedlow, Gianluigi Zanetti

CRS4

Wellcome Trust Centre for Gene Regulation and Expression – University of Dundee

July 1st, 2013

Motivation

CRS4

- A public multi-disciplinary research center in Italy
- Focuses on applied computational sciences
- Within top 5 Italian computation facilities

CRS4 Sequencing and Genotyping Platform

- Currently the largest sequencing center in Italy
- Has enabled a number of studies on the Sardinian population

Sequencing Equipment: 3 Illumina HiSeq2000, plus older sequencers

Sequencing Capacity: about 5 Tbases/month

Since Sept. 2010 we've sequenced about. . .

- over 2000 whole-genome samples (low-pass, high-coverage)
 - some cancer genomes as well
- 800 total RNA samples
- 100 exomes
- a handful (30) of ChIP-Seq samples

- The number of samples and the amount of data to handle presented significant difficulties
 - difficulties scaling computational throughput
 - difficulties tracking the process
 - ample opportunities for inefficiencies

- The number of samples and the amount of data to handle presented significant difficulties
 - difficulties scaling computational throughput
 - difficulties tracking the process
 - ample opportunities for inefficiencies

Wishlist

We wanted to improve our process in several ways:

- automated processing
 - hands off from when the sequencer is started to deliverable data
- trace all data processing activities
- effectively manage file storage
- computational scalability

Our solution

Automated processing and tracking platform

To satisfy those requirements, we implemented a solution based on five core components:

- Galaxy
- the “Automator”
- OMERO.biobank
- iRODS
- Hadoop

Requirement

Automated processing

- Monitor sequencers, detect when new data is ready
- Automatically run data through standard pipelines
- Notifications to operator for data ready and for errors
- Automatically track these data sets and how they were generated

In part implemented with Galaxy; in part with custom software

Since you're at GCC, I can probably assume you know what Galaxy is!

Key features for our application

Workflows: give a way to define automated analysis “recipes”

Histories: saves sequence of tool invocations that produced a data set

- A convenient way to trace reproducible actions performed on data

REST API: provides some degree of programmatic access

- e.g., launch workflows, retrieve results

Familiarity: Galaxy was a desirable and familiar tool for our users

One tool for both automation and downstream analysis

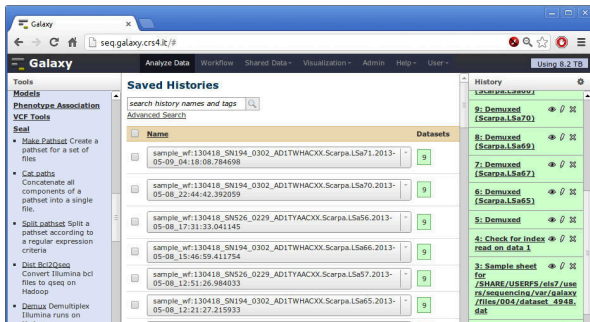
In our system, we use two Galaxy instances

Public instance

- Collect sample data and flowcell configuration from wet lab and/or client
- Return processed samples to users through data set libraries (WIP)
 - Possibly through integration with iRODS
- Currently runs a customized version of the nglims Galaxy fork by Brad Chapman (Harvard School of Public Health Bioinformatics)

Private instance

- Manages execution of our standard processing workflows
- Accumulates processing history for each flowcell and sample
 - We fetch this information through the Galaxy REST API (using the bioblend Python module by Afgan et al.)
- The private instance is a standard version of Galaxy



Galaxy

seq.galaxy.crs4.it/#

Galaxy Analyze Data Workflow Shared Data Visualization Admin Help User Using 8.2 TB

Tools

Models

Phenotype Association

VCF Tools

Seal

- **Make Pathset** Create a pathset for a set of files
- **Cat paths** Concatenate all components of a pathset into a single file.
- **Split pathset** Split a pathset according to a regular expression criteria
- **Dist Bcl2Seq** Convert Illumina bcl files to seq on Hadoop
- **Demux** Demultiplex Illumina runs on Hadoop

Saved Histories

search history names and tags

Advanced Search

Name	Datasets
sample_wf:130418_S1194_0302_AD1TWHACXX.Scarpa.L5e71.2013-05-09_04:18:08.784698	9
sample_wf:130418_S1194_0302_AD1TWHACXX.Scarpa.L5e70.2013-05-08_22:44:42.392059	9
sample_wf:130418_S1194_0229_AD1TYAACXX.Scarpa.L5a56.2013-05-08_17:31:33.041145	9
sample_wf:130418_S1194_0302_AD1TWHACXX.Scarpa.L5a66.2013-05-08_15:46:59.411754	9
sample_wf:130418_S1194_0229_AD1TYAACXX.Scarpa.L5a57.2013-05-08_12:51:26.984033	9
sample_wf:130418_S1194_0302_AD1TWHACXX.Scarpa.L5a65.2013-05-08_12:21:27.215933	9

History

- 9: Demuxed (Scarpa.L5a70)
- 8: Demuxed (Scarpa.L5a69)
- 7: Demuxed (Scarpa.L5a67)
- 6: Demuxed (Scarpa.L5a65)
- 5: Demuxed
- 4: Check for index read on data 1
- 3: Sample sheet for /SHARE/USERS/els7/user/sequencing/var/galaxy/files/004/dataset_4948.dat

Requirement

Automated processing

- We found Galaxy alone to be insufficient for full automation
- Clumsy for housekeeping tasks
 - e.g., move/rename files, interact with other services
- Workflows sometimes aren't sufficiently expressive
 - e.g., linking to variable number of outputs
 - no "if" operator

Our custom Automation package

- Part interfaces to programmatically control other components
 - i.e., Galaxy, iRODS, OMERO
- Part distributed event-dispatching daemon based on RabbitMQ
- Part custom-made event handlers
 - These use the aforementioned interfaces and anything else they need to implement actions, that may emit new events

Task division

- Galaxy handles all operations that transform or create datasets
 - allows us to easily create a history
- The Automator does all other operations, including driving Galaxy

Requirement

Trace all data processing activities

- Essential for reproducibility
- For any data set generated, track how it was created
 - Actions on data sets
- Track relations between data sets
- Database should support appropriate queries; e.g.,
 - From what flowcell was the dataset derived?
 - Through which operations/parameters?
 - With which other samples was it normalized?
 - What other data sets came from the same batch?

- OMERO is a “model-driven data management platform for experimental biology” (Allan, et al.; Nature Methods, 2012)
- Stores a graph structure where data set nodes are connected by actions
- Nodes and actions are tagged with model-dependant information
- We extended OMERO to handle data types produced in sequencing and microarray experiments
 - OMERO.biobank – will be included in official OMERO releases
 - In our model, we store information about samples, data paths, data format, and the data set's *entire processing history*

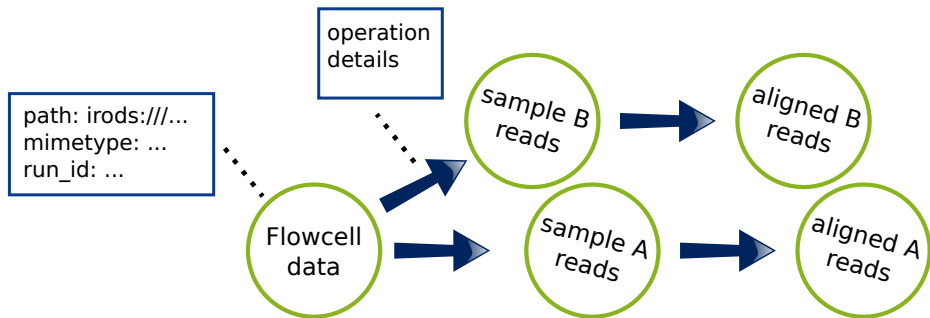


Fig: An example of an OMERO graph

Requirement

Effectively manage file storage

- We generate lots of files
- Incremental growth → multiple file systems
- Geographically dispersed collaborations

iRODS

integrated Rule-Oriented Data-management System

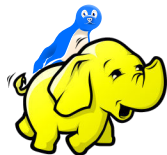
- A file cataloguing system
- We have used it to create a single go-to place to find data files
 - Simplifies accessing data on complex storage architecture
- Allows us to tag files with attributes (e.g., run id, sample id, etc.) and use the tags in queries
- Optimized file transfers

Requirement

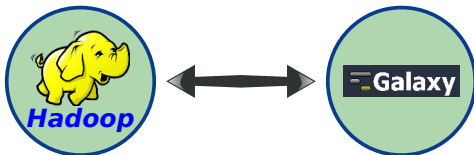
Computational scalability

- Large data generation rate from sequencers
- Interest in minimizing turn-around time
- Need to scale out computation over many nodes

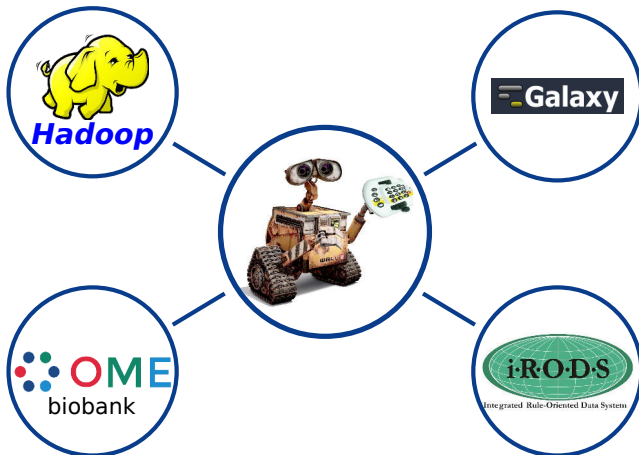
- The system that enables many data-centered companies
 - e.g., Twitter, Facebook, Yahoo, ...
- Automatically handles:
 - distribution of computation and data
 - node and task failures
- To leverage Hadoop in this context, we adopted:
 - Seal: toolkit for Hadoop-based sequencing data processing
 - demultiplexing, alignment (based on BWA, sorting, etc.)
 - Pydoop: Python API for Hadoop
 - A dependency for Seal, but also used for custom tools and scripts
 - SeqPig: SQL-like scripting for Hadoop with sequencing-specific functionality



- We've implemented a thin integration between Galaxy and Hadoop
- Can run Hadoop-based programs via Galaxy
- Big Hadoop datasets referenced by Galaxy through a “pointer” type, *pathset*
 - a file containing a list of paths
- Hadoop jobs executed through a wrapper that knows how to interpret pathset files and passes the correct arguments to the Hadoop job
- Galaxy dataset clean-up program also has to learn about pathsets (WIP)



Bringing it all together...



Conclusion

Currently. . .

- The system is currently being used to process the sequences produced by our center
 - is able to process flowcells in complete automation
- Through the stored histories, we ensure reproducibility
- Development is ongoing to improve it
- Is proving to be a good solution to our original problem

- Better management and monitoring
- Better error handling and restart

- Better management and monitoring
- Better error handling and restart

Open source?

- Some parts already are: <https://github.com/crs4>
 - Omero.biobank
 - Seal, Pydoop
- Hopefully by the fall we'll also release:
 - the automator
 - Hadoop-Galaxy integration

A shameless plug...

- Public Galaxy instance for NGS microbiology data by CRS4
- If you're interested, go see the poster!

<http://orione.crs4.it>

A shameless plug...

- Public Galaxy instance for NGS microbiology data by CRS4
- If you're interested, go see the poster!

<http://orione.crs4.it>

Thanks for your attention!