# Galaxy

### Internals, organization, control flow

www.getgalaxy.org

James, Dan, Greg



@jxtx / #usegalaxy

#### The Plan

1. What's in the galaxy-central repository? 2. Galaxy web application architecture 3. Control flow in the Galaxy web application 4. Tools in the age of the toolshed 5. Galaxy Workflows 6. Galaxy data organization

#### 0. The right ways to be involved with Galaxy

#### IRC: <u>irc.freenode.net</u> #galaxyproject

Trello: <u>https://trello.com/b/75c1kASa/galaxy-</u> <u>development</u>

#### 1. Getting around the Galaxy repository

### (lets look at the code)

#### 2. Galaxy application architecture







The old way	
User stuff (prefs, etc)	
Tool forms	
Reports	
Tool shed	
*Many of these have an API but it is not yet used by the UI	

The new way Visualizations History Tool menu Most grids

In between

Workflows

Data Libraries

#### Example of the new way: Tool menu generation

#### Example of the new way: Tool form generation

#### So many languages!



### 3: Control flow: running jobs







#### Processes



#### 4. Tools in the age of the toolshed

#### The old way

1. Each tool specified by a tool.xml somewhere on the local filesystem (but typically under tools)

2. Tools to be loaded specified in tool\_conf.xml, loaded by Galaxy at startup — *no representation in database beyond tool ids* 

No way to access old tool configurations after updates

#### In the ToolShed

ToolShed Repository stored as mercurial repo on disk in ToolShed several types: unrestricted, suite, tool dependency unrestricted can have multiple installable revisions

#### lib.galaxy.webapps.tool\_shed / lib.tool\_shed

#### In the ToolShed



### Installed in Galaxy

#### app.install\_model



### 5. Galaxy workflows

#### lib/galaxy/workflow/modules.py



Workflow modules have:

# Config time state — in the workflow editor used to generate the form associated with a given step and update it

Runtime state — similar but used for parameters set at workflow runtime

As well as conversion from JSON <-> Workflow Module instance <-> workflow\_step encoded in database Workflow scheduling:

Currently workflows are scheduled like any other job

All intermediate datasets and connections are created and each step is sent as a job to the JobManager

Pausing: when intermediate steps fail the workflow is paused. Although, this actually applies to any dependent jobs

#### 6. Galaxy data organization

Where does data in Galaxy go?

 "Metadata" is stored in a SQL database (preferable Postgres): Users, workflows, histories, dataset metadata... everything a user creates interacting with Galaxy except the raw contents of datasets

2. Dataset contents is stored in file\_path, typically database/files

3. Data used by tools that is not user specific is stored in



https://wiki.galaxyproject.org/Admin/Internals/DataModel



Galaxy data model is not database entity driven

Entities are defined in galaxy.model as objects

SQLAIchemy is used for object relation mapping

Mappings are defined in galaxy.model.mapping in two parts — a table definition and a mapping between objects and tables including relationships

Migrations allow the schema to be migrated forward automatically

It *rarely* makes sense to access the Galaxy database directly

Where does data in Galaxy go?

- "Metadata" is stored in a SQL database (preferable Postgres): Users, workflows, histories, dataset metadata... everything a user creates interacting with Galaxy except the raw contents of datasets
  - 2. Dataset contents is stored in file\_path, typically database/files objectstore

3. Data used by tools that is not user specific is stored in



>> fh = open( dataset.file\_path, 'w' )
>>> fh.write('foo')
>>> fh.close()
>>> fh = open( dataset.file\_path, 'r')
>>> fh.read()

>>> update\_from\_file( dataset, file\_name='foo.txt')
>>> get\_data( dataset )
>>> get\_data( dataset, start=42, count=4096 )

**Distributed Object Store** 



#### Benefits

- Grow beyond original capacity
- Avoid migrating data offline
- Tier storage
- Let your users bring their own storage
- Use resources w/o a shared filesystem (with iRODS)
- Remove IO bottlenecks

#### Data tables and location files

Data tables provide an abstraction which tools use to access indexes of data which can be accessed on the local filesystem Tool config reference data table by name with abstract columns

<param name="mafType" type="select" label="Choose alignments">
 <options from\_data\_table='indexed\_maf\_files'>
 <filter type="data\_meta" ref="input1" key="dbkey" column="dbkey"
 <validator type="no\_options" message="No alignments are available
 </options>
</param>

Data tables configuration maps abstract data table to a concrete file

## Which can provide any information, but typically locations of data at a given site

#This is a sample file distributed with Galaxy that is used by some #alignment tools. The maf\_index.loc file has this format (white space #characters are TAB characters): #<Display\_name UID> <indexed\_for:build1,build2,build3> <exists\_in\_maf:build1,build2,build3> <Comma\_Se # ENCODE TBA (hq17) ENCODE\_TBA\_hg17 armadillo,baboon,galGal2,panTro1,colobus\_monkey,cow,canFam1,dusky\_ti ENCODE MAVID (hg17) ENCODE\_MAVID\_hg17 armadillo,baboon,galGal2,panTro1,colobus\_monkey,cow,canFam1,dusky ENCODE TBA (hq16) ENCODE\_TBA\_hg16 armadillo,baboon,galGal2,panTro1,colobus\_monkey,cow,canFam1,dusky\_ti 8-way multiZ (hg17) 8\_WAY\_MULTIZ\_hg17 canFam1,danRer1,fr1,galGal2,hg17,mm5,panTro1,rn3 canFam1, danRe 17-way multiZ (hq18) 17\_WAY\_MULTIZ\_hg18 hg18, panTro1, bosTau2, rheMac2, mm8, rn4, canFam2, echTel1, loxAfr1 3-way multiZ (hg18,panTro2,rheMac2) 3\_WAY\_MULTIZ\_hg18 hg18,panTro2,rheMac2 hg18,panTro2,rheMac2 5-way multiZ (hg18,panTro2,rheMac2,mm8,canFam2) 5\_WAY\_MULTIZ\_hg18 hg18, panTro2, rheMac2, mm8, canFam2 28-way multiZ (hg18) 28\_WAY\_MULTIZ\_hg18 hg18 hg18,dasNov1,otoGar1,felCat3,galGal3,panTro2,bosTau3 15-way multiZ (dm2) 15\_WAY\_MULTIZ\_dm2 dm2,droSim1,droSec1,droYak2,droEre2,droAna3,dp4,droPer1,droWil1,d 17-way multiZ (mm8) 17\_WAY\_MULTIZ\_mm8 mm8 hg18,panTro1,bosTau2,rheMac2,mm8,rn4,canFam2,echTel1,loxAfr1 8-way multiZ (ponAbe2) 8\_WAY\_MULTIZ\_ponAbe2 ponAbe2,hg18,panTro2,rheMac2,calJac1,mm9,monDom4,ornAna1

#### Data Managers

Special class of Galaxy tool which allows for the download and/or creation of data that is stored within Data Tables and their location files.

These tools handle e.g. the creation of indexes and the addition of entries/lines to the data table / .loc file via the Galaxy admin interface.

Data Managers can be defined locally or installed through the Tool Shed.

Available in: Admin GUI, Workflows, API

#### Special class of Galaxy tool

<tool id="data\_manager\_fetch\_genome\_all\_fasta" name="Reference Genome" version="0.0.1" tool\_type="manage\_data">

<outputs> <data name="out\_file" format="data\_manager\_json"/> </outputs> Writes a JSON description of new data table entries as content of tool output file "data\_tables" [ "all\_fasta":[ "path" "sacCer2.fa", "dbkey" "sacCer2" "name": "S. cerevisiae June 2008 (SGD/sacCer2) (sacCer2)", "value": "sacCer2" }

#### This creates a new entry in the Tool Data Table:

#<unique\_build\_id> <dbkey> <display\_name> <file\_path>
sacCer2 sacCer2 S. cerevisiae June 2008 (SGD/sacCer2) (sacCer2) /Users/dan/galaxy-central/tool-data/sacCer2/seq/sacCer2.fa

# Where the sacCer2.fa file was placed by the tool in the output file's extra\_files\_path

#### data\_manager entry inside <data\_managers> tag in data\_manager\_conf.xml

```
<data_manager too]_file="data_manager/bwa_index_builder.xml" id="bwa_index_builder" version="0.0.1">
    <data_table name="bwa_indexes">
        <output>
            <column name="value" />
            <column name="dbkey" />
            <column name="name" />
            <column name="path" output_ref="out_file" >
                <move type="directory" relativize_symlinks="True">
                    <target base="${GALAXY_DATA_MANAGER_DATA_PATH}">${dbkey}/bwa_index/${value}</target>
                </move>
                <value_translation>${GALAXY_DATA_MANAGER_DATA_PATH}/${dbkey}/bwa_index/${value}/${path}</value_translation</pre>
                <value_translation type="function">abspath</value_translation>
            </column>
        </output>
    </data_table>
</data_manager>
```

### informs Galaxy about which data tables to expect for new entries special handling of provided JSON values and files

### Q&A