# GCC2015 Galaxy Hackathons

GCC Hackathon (Code and Data!) Organizers
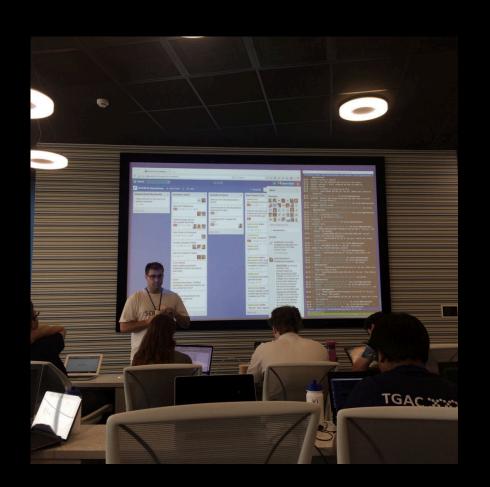
# Making Connections, Providing Footholds

- Community Building
- Bandwidth is important
- Great opportunity for 'fun' projects
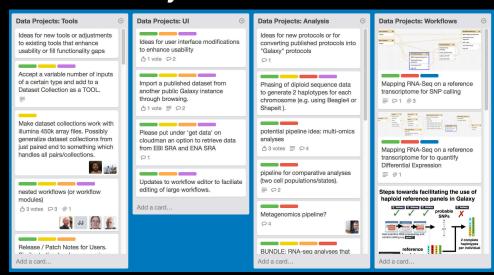
# By the Numbers

- **>50** people
- **2** days
- **1** trello board
- **6** major projects
- **20** pizzas
- **163** cups of coffee

# Main Themes

bit.ly/hackboard

- Workflows
- Tools
- API
- Code Cleanup
- Smashing Eggs

# Nested Workflows

- Goal: Enable nesting of workflows inside of other workflows.
- Steps
  a. construct tests of workflows as we want them to appear (both from the API and as unit tests)
  b. refactor tool logic to a new AbstractExecutable object and reimplement tool or workflow-specific logic

Work is happening at: https://github.com/nebiolabs/galaxy

(Peter van Huesden, Marius van den Beek, Brad Langhorst)

# Tools

- BLAST+ XML2 Datatype
  a. NCBI Released a new output format: 14
  b. New challenges - produces multiple files
  c. Working on a Composite datatype to capture the relationships between these files

(Peter Cock, Carrie Ganote, Dave Bouvier; See https://github.com/peterjc/galaxy_blast)

# [WIP] Mega flake8 linting #433

**Open** **nsoranzo** wants to merge 51 commits into `galaxyproject:dev` from `bgruening:lint2`

💬 Conversation 3    ◦ Commits 51    ⊞ Files changed 379

**nsoranzo** commented 3 days ago | Owner | ✎

This is the joint work of **@bgruening** and **@remimarenco** with help and fixes from me, **@dannon** and other GCC2015 hackathon participants.
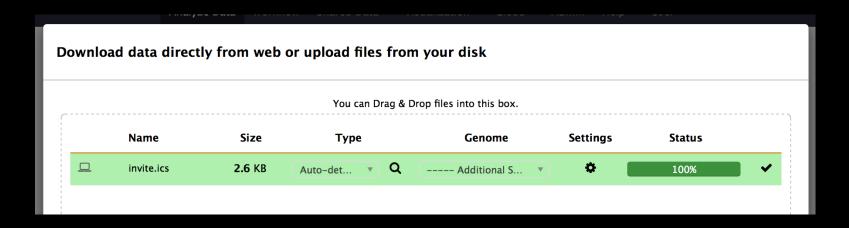
# Data Libraries API and Beta UI Tweaks

- Add the ability to import whole folders to a history
- Add ability to recursively download whole folders
- Fix Galaxy internal StreamBall module so it will create file archives with duplicated filenames properly

https://github.com/galaxyproject/galaxy/pull/426

Made by: Eric Enns - @EJEnns, Martin Čech

# Composite Uploads

- Goal: Enable the upload of composite datasets.
- Steps
  a. construct UI to select composite dataset components
  b. submit selected files to API and execute the upload tool

# Smashing Eggs

Eggs provide Galaxy framework dependencies

Why Smash?

- Outdated format
- Replace archaic, fragile download and build code with standard *pip*-based installation

+136 −5,471

Still WIP but we're close:

*github.com/galaxyproject/galaxy/pull/428*

# Data Hack



*Purpose*

- Bring together end-users & coders
- Identify issues with Galaxy usability and other bottlenecks from an **end-user perspective**
- Generation of "Best Practices" analysis workflows
- Wrap tools and pipelines to simplify and address common data manipulations
- Plan to publish (*GigaScience*) our initial workflow solutions along with a cloudshare with these workflows and example data

# Data Hack

- *Results*

- "Best Practices" analysis workflows for RNA-seq for multiple use cases, variant calling mini-pipeline

- New tools wrapped: Beagle, Impute and StepIt and pipelines for variant calling

- Identified & fixed: broken tools, new capabilities, wrapping existing tools

This will continue in a more formalized form with the formation of a "GalaxyScientists" group to provide a synthesized view from the end-user community for new functionalities, tools and capabilities.

**Johns Hopkins University Data Science Specialization Program**

# Thanks!

- Johns Hopkins Data Science

- TGAC - The Genome Analysis Centre

# The End