# Galaxy as a Platform for High-throughput Genomics

Jeremy Goecks

EMORY UNIVERSITY

# Topics

**Galaxy**

Analyzing Cancer Genomes &
Transcriptomes

Web-based Visual Analysis

# Genomic Analyses are Difficult

Investigators unfamiliar with computation

Creating and reproducing workflows (pipelines) hindered by complexity: systems, scripts, tools, parameters

Collaboration and reuse difficult because current approaches do not support computational artifacts well

```
@HWI-ST565:241:D19R1ACXX:1:1101:5456:1997 1:N:0:CGATGT
NGCATAGGCAAGCACCGGAAGCACCCCGGCGGCCGCGGTAATGCTGGTGGTCTGCATCACCACCGGATCAACTTCGACAAATACCACCCAGGCTACTTTGG
+
#1=DDDDDHHHHHHIIIIIIIIIIIAHIIIIIIIIFECBBB;BBECCCCC?CBBCCCCCCCCCCCCBBBBBBBBBCCCCCBBBBBCCCCCCCBBBBBBBBCCCCDCC
@HWI-ST565:241:D19R1ACXX:1:1101:7520:1996 1:N:0:CGATGT
NCGCAACCTCAACACCACCTTCTTCGACCCCGCCGGAGGAGGAGACCCCATTCTATACCAACACCTATTCTGATTTTTCGGTCACCCTGAAGTTTATATTC
+
#4=DFFFFHHHHHJJJJJJJJJJJJJJJJJJJJJJIJHFHFFFDDDDDDDDDFEEEFDDDDDDDDDDDDEEEEEDEDDDDDDDDDDDDDCDDD@CDEDEEE>
@HWI-ST565:241:D19R1ACXX:1:1101:10117:1998 1:N:0:CGATGT
NATGTGCCCTCTGGCAGTCTGCTGCTGTGTCCAGAGTCCGACTCCAGCTGGGCTGTAACTGGGCTTGGCCCCCGCCTTAGGCCCCGCCAGCAGGCGAAGCA
+
#1=DDFFFHHHHHJIJJJJJJJJJJJJIJJJJJJJJJIJIJJJJJJJIJJJJJJJIIJIGIJJJJJJHHHFFFFDDDDB=BCDDDDDDDDDDDDCDDBB9>BDA
@HWI-ST565:241:D19R1ACXX:1:1101:10283:1992 1:N:0:CGATGT
NTTGTCACCAAGACCCTACTTCTAACCTCCCTGTTCTTATGAATTCGAACAGCATACCCCCGATTCCGCTACGACCAACTCATACACCTCCTATGAAAAAA
+
#4=DDFFFHHHHHJJJJJJJJJJJJJJJJJJJIJJJJJJJIJJJJJJJJJJJJJJJJGIHHHFFFFFDDDDDDDBDDDDDDDDDDDDDDDDDDDD@DDCD<
@HWI-ST565:241:D19R1ACXX:1:1101:10632:1993 1:N:0:CGATGT
NGTAAGCCTTCTATGCATCCACACCAAAATCCTGCAGAATGTAAGTAAGCTCTGCTTTATAAGATGGGTTCACCTTCATCGCAGACTGAAAGTTTCAGTTT
+
#1:ADDFFHHHHHGHGIGIGIJIIGGIGJIHGGIJGHIEHJDDHFBGIGGHJJJJJJJJJJIIGCFGII@CGCHGIHHEFGFDFDBEDDCCC@5>CDCA;>A
@HWI-ST565:241:D19R1ACXX:1:1101:10895:1991 1:N:0:CGATGT
NCTAGCACAGAGAGTTCTCCCAGTAGGTTAATAGTGGGGGGTAAGGCGAGGTTGGCGAGGCTTGCTAGAAGTCATCAAAAAGCTATTAGTGGGAGCAGAGT
+
#1=DFFFFHHHHHJGHIJJJJIJGHIJHIIJJJJFIIJJJD0:CDDDDDDD<BCDDDDDDDDDDDDDDADCCCDDCDDDDDDDDDDDEEDDBD?@@BCAAC
@HWI-ST565:241:D19R1ACXX:1:1101:10838:1994 1:N:0:CGATGT
NTCCCTGCTACTGCTGATGCACTGTCCTCTCCCTGCAGCCCCTGGCTTCCCAGCCTTCCTCCTGACCCCTTCCAACAGCCTTGGAACTCCAGCTGCCACCA
+
#1=DFFFFHHHHHJJJJJJJJJIJJJJJJJJJJJJJJJJJJJIJJJJJJJJJIJJJJJJJJJJJJJJJIHHHHFFFFFDEEDDCDDDDDDDDDDDDDDDDDD8
@HWI-ST565:241:D19R1ACXX:1:1101:11757:1991 1:N:0:CGATGT
NTTACCACTGGGAAAATCTTTTTAGTTAGATTGTAGGCTTCCTGGGGCCTCCCATGCCAGGACTGCAAAGTGATCCAGCCCTACCTGTCTTCCCACCTGTG
+
#4=DFFFFHHGHHJJJJIIJJJJJHIIJJJJIJIIIFJJJJIJGGHIJIJJJJJJJIHIJJJJIJJIHHHHHHFFFFFFFEDDEDDDDCDDDDDDDDDDBDC@
@HWI-ST565:241:D19R1ACXX:1:1101:11780:1992 1:N:0:CGATGT
NTAAATACTAAGCACAAGCTCACTTCCCTCTTGGTCAGGTGGTTGTTTTAGAGCTACTCGATATTTATAACTTTTTATAAGCACCGGTCATTTTTTGAGA
+
#1=DDFFFHHHHHJJJJJJJJIJJIJJJJJJJJJGIIJJJCGH?GHIIIJJJJJIJJJJJJJHIJJJIIHHEEHHHFFFFFEEEEDDD@DDDEFDDDDBDD
@HWI-ST565:241:D19R1ACXX:1:1101:12154:1991 1:N:0:CGATGT
NCGCTTTGGGAGGCTGAGGTGGGAAGATGACTTGAGCCCAGGAGTTCGAGACCAGCCTGGGCAACATGGTAAAACCCTTTCTCTGACCCCCACAAAAATAA
+
#1:DBDDDDDDDDID;CBB2ACA?EDDDI@DED?BDDDDDDID;B8.<@ADD@C@;5?A???@?@AA@=>5>>AA=<<AAAAAAAA>>>????>>>>?8A>
@HWI-ST565:241:D19R1ACXX:1:1101:12096:1998 1:N:0:CGATGT
NTCTGATGTTGCTGATCTCCGTGGCTGTGACCATCATGGCTGGTGACCACACTCCTTCTGCCCAGTTCGGCTGGAAAACTCTGGGAACTGCAGCACGAGAT
+
#1=DDFFFHHHHHJJJJJJJJJJJJJJJJIJJJJJJJJJIIJJIJJ?DA@FGIGGGGHGIGGGGIHGGDHEFFDDCD2>CC;;>;5>C5?<>ACCA>A?<@938
```

3

# Galaxy Project: Fundamental Questions

When genomics (or any other biomedical science) becomes dependent on computational methods, how to:

- make tools and workflows **accessible** to scientists?
- ensure that analyses are **reproducible**?
- enable transparent **communication and reuse** of analyses?

# Vision

Galaxy is an **open, Web-based platform** for accessible, reproducible, and collaborative computational genomics

Goecks et al. (2010) *Genome Biology*

# Galaxy Demo

# What is Galaxy?

**Platform for high-throughput genomics**
1. get and integrate public, private data
2. analyze data and create workflows
3. visualization, sharing, publication

**Customizable open-source software on various HPC resources**
- public website — http://usegalaxy.org
- local instance
- on the cloud

Galaxy platform
- ✦ run tools, workflows on HPC resources
- ✦ create workflows, visualizations, pages
- ✦ share *everything*

8

# Cloud Launch

11

# Cloud Features

Resource configuration

Autoscaling

Snapshotting

# Galaxy is Very Popular

Public Website (http://usegalaxy.org), anybody can use:

- ✦ ~500 new users per month, ~200 TB of user data, ~130,000 analysis jobs per month

Used and cited in more than 1000 publications

# Galaxy is Very Popular

## Local installations all over the world

http://bioteam.net/slipstream/galaxy-edition/

# Topics

Galaxy

**Analyzing Cancer Genomes & Transcriptomes**

Web-based Visual Analysis

# Preliminary Data

6 patients, whole transcriptome sequencing (RNA-seq) of primary tumor
- mixed populations!
- 3 +ERCC, 3 -ERCC (via IHC)

MiaPaCa2 cell line
- whole transcriptome
- targeted exome

Total sequencing data: ~70 GB



http://en.wikipedia.org/wiki/RNA-Seq

日本語要約

# The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity

Jordi Barretina, Giordano Caponigro, Nicolas Stransky, Kavitha Venkatesan, Adam A. Margolin, Sungjoon Kim, Christopher J. Wilson, Joseph Lehár, Gregory V. Kryukov, Dmitriy Sonkin, Anupama Reddy, Manway Liu, Lauren Murray, Michael F. Berger, John E. Monahan, Paula Morais, Jodi Meltzer, Adam Korejwa, Judit Jané-Valbuena, Felipa A. Mapa, Joseph Thibault, Eva Bric-Furlong, Pichai Raman, Aaron Shipway, Ingo H. Engels ⊞ *et al.*

Affiliations | Contributions | Corresponding authors

# Big Questions

How closely does MiaPaCa2 match to primary pancreatic tumors?

How to match patient to "best" CCLE cell line(s)?

# Using Galaxy for Analysis of Cancer Transcriptomes

New tools
- ✦ complement existing transcriptome analysis tools

New workflows
- ✦ workflows are understandable, extendable, sharable

New visual analysis applications
- ✦ visualize and call variants in a Web browser

# Single Sample Transcriptome Analysis

# Comparing Called Variants with Public Datasets

# Patient Mutations vs.

| | P1 | P2 | P3 | P4 | P5 | P6 | CL |
|---|---|---|---|---|---|---|---|
| OM MIA (4) | 0 | 1 | 1 | 0 | 0 | 0 | 4 |
| OM PC (11) | 0 | 1 | 1 | 0 | 0 | 0 | 4 |
| OM ALL (114) | 0 | 2 | 1 | 1 | 1 | 1 | 4 |
| HP MIA (84) | 3 | 6 | 4 | 5 | 4 | 3 | 15 |
| HP PC (1769) | 16 | 23 | 19 | 11 | 23 | 8 | 39 |
| HP ALL (64,669) | 110 | 180 | 143 | 97 | 136 | 65 | 87 |

OM = OncoMap, HP = hybrid capture with probes

# Patient Mutations vs.

**CCLE** Cancer Cell Line Encyclopedia

http://www.broadinstitute.org/ccle/home

| | P1 | P2 | P3 | P4 | P5 | P6 | CL |
|---|---|---|---|---|---|---|---|
| OM MIA (4) | 0 | 1 | 1 | 0 | 0 | 0 | 4 |
| OM PC (11) | 0 | 1 | 1 | 0 | 0 | 0 | 4 |
| OM ALL (114) | 0 | 2 | 1 | 1 | 1 | 1 | 4 |
| HP MIA (84) | 3 | 6 | 4 | 5 | 4 | 3 | 15 |
| HP PC (1769) | 16 | 23 | 19 | 11 | 23 | 8 | 39 |
| HP ALL (64,669) | 110 | 180 | 143 | 97 | 136 | 65 | 87 |

**Cell line does not appear very similar to tumors**
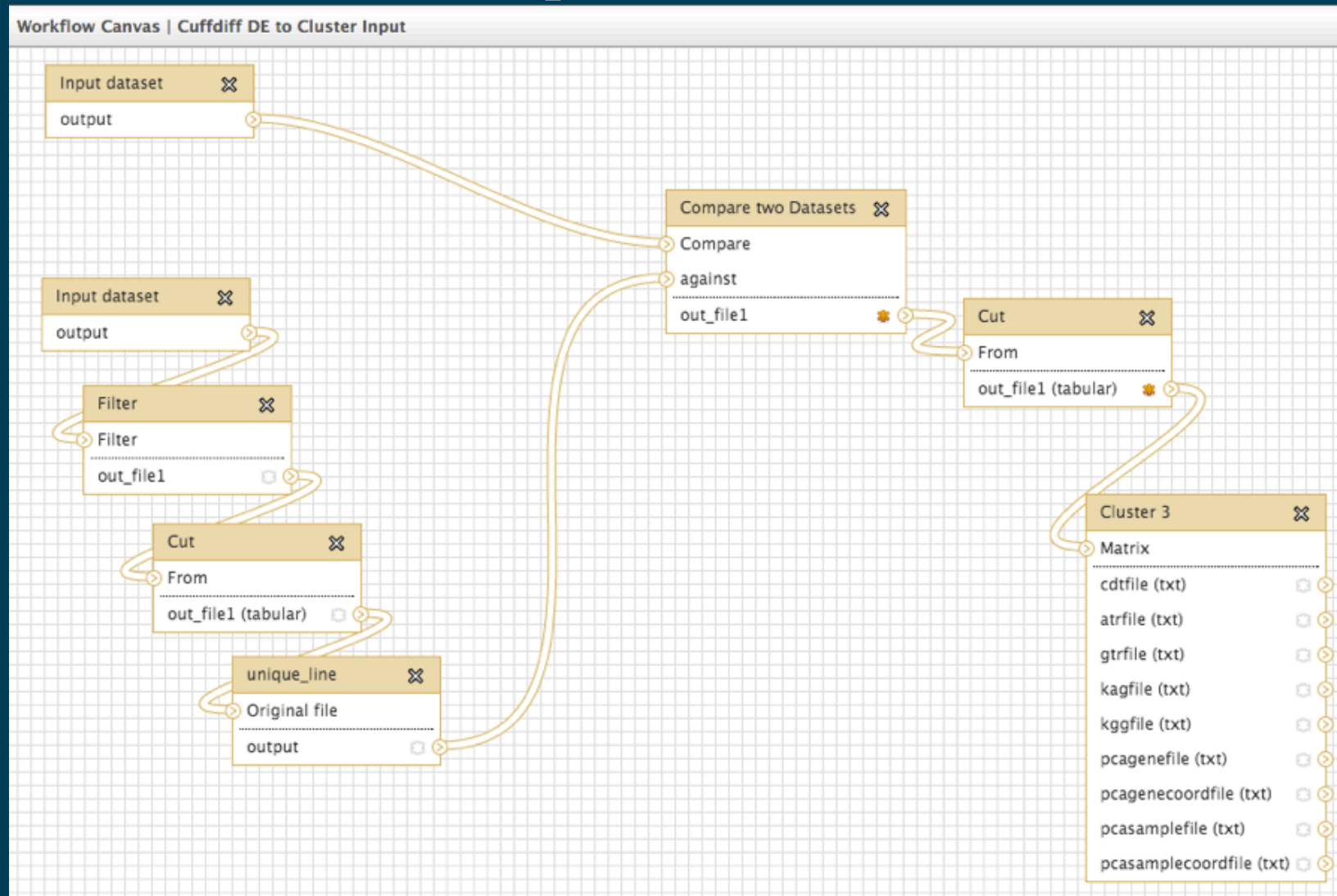
OM = OncoMap, HP = hybrid capture with probes

24

# Patient Mutations to Predict Tumor Attributes

|  | P1 | P2 | P3 | P4 | P5 | P6 |
|---|---|---|---|---|---|---|
| OM MIA (4) | 0 | 1 | 1 | 0 | 0 | 0 |
| OM PC (11) | 0 | 1 | 1 | 0 | 0 | 0 |
| OM ALL (114) | 0 | 2 | 1 | 1 | 1 | 1 |
| HP MIA (84) | 3 | 6 | 4 | 5 | 4 | 3 |
| HP PC (1769) | 16 | 23 | 19 | 11 | 23 | 8 |
| HP ALL (64,669) | 110 | 180 | 143 | 97 | 136 | 65 |
| Tumor % | 90% | 90% | 100% | 0%? | 60% | 40% |

OM = OncoMap, HP = hybrid capture with probes

25

# Clustering via Differential Expression

# Gene Expression Clustering



-0.14

0.31

0.46

0.48

0.62

0.77

P1    P2    P3    P5    P4    P6    CL
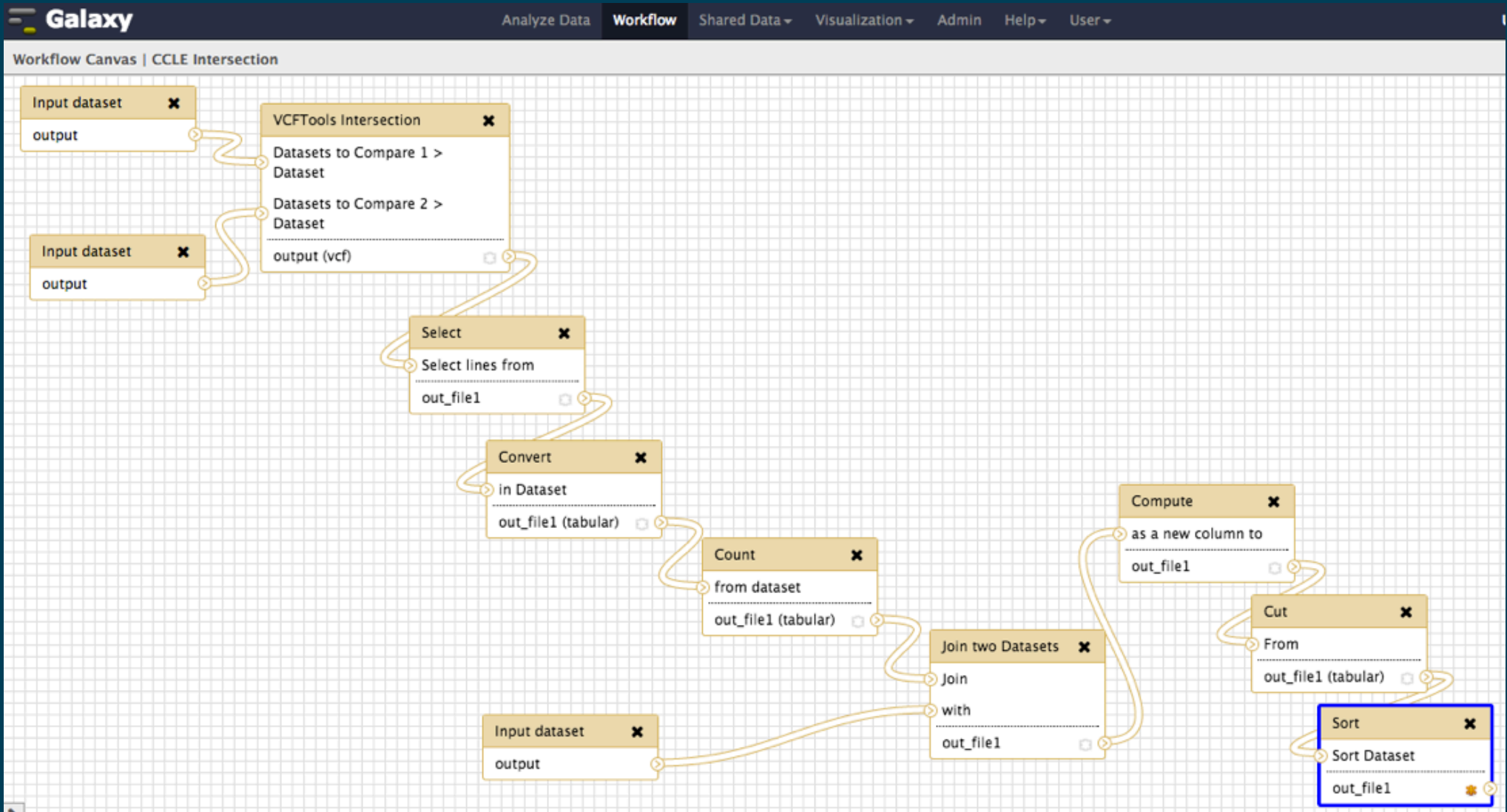
Spearman Correlation

# Gene Expression Clustering

Spearman Correlation

# Matching Patients to Cell Lines

# Matching Patients to Cell Lines

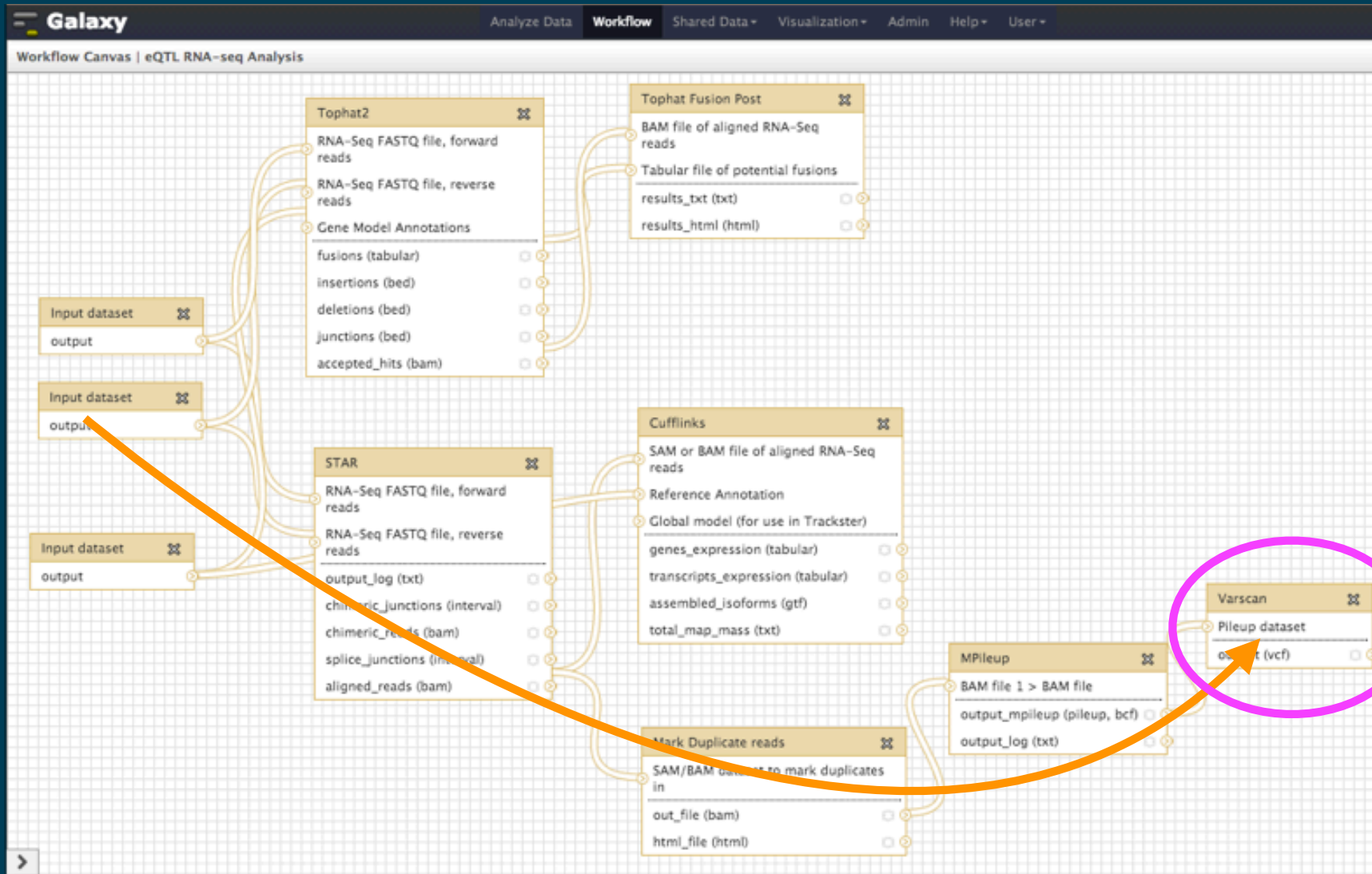| | Best Match (# muts) | Best Match (#, % muts) |
|---|---|---|
| P1 | KP3, KP2, KP4, PANC0327, PANC1005, QGP1 (4) | KP3 (4, 4.8%) |
| P2 | PANC0327 (8) | CAPAN2 (7, 9.1%) |
| P3 | SNU410, QGP1 (6) | KP3 (5, 6.0%) |
| P4 | CAPAN2, PANC0403, MIAPACA2, PANC0327 (5) | CAPAN (5, 6.5%) |
| P5 | T3M4 (8) | T3M4 (8, 8.7%) |
| P6 | CAPAN2, MIAPACA2 (3) | CAPAN2 (3, 4.8%) |

*38 pancreatic cell lines in CCLE

# Topics

Galaxy

Analyzing Cancer Genomes  &
Transcriptomes

**Web-based Visual Analysis**

# Mutation Calling from RNA-seq



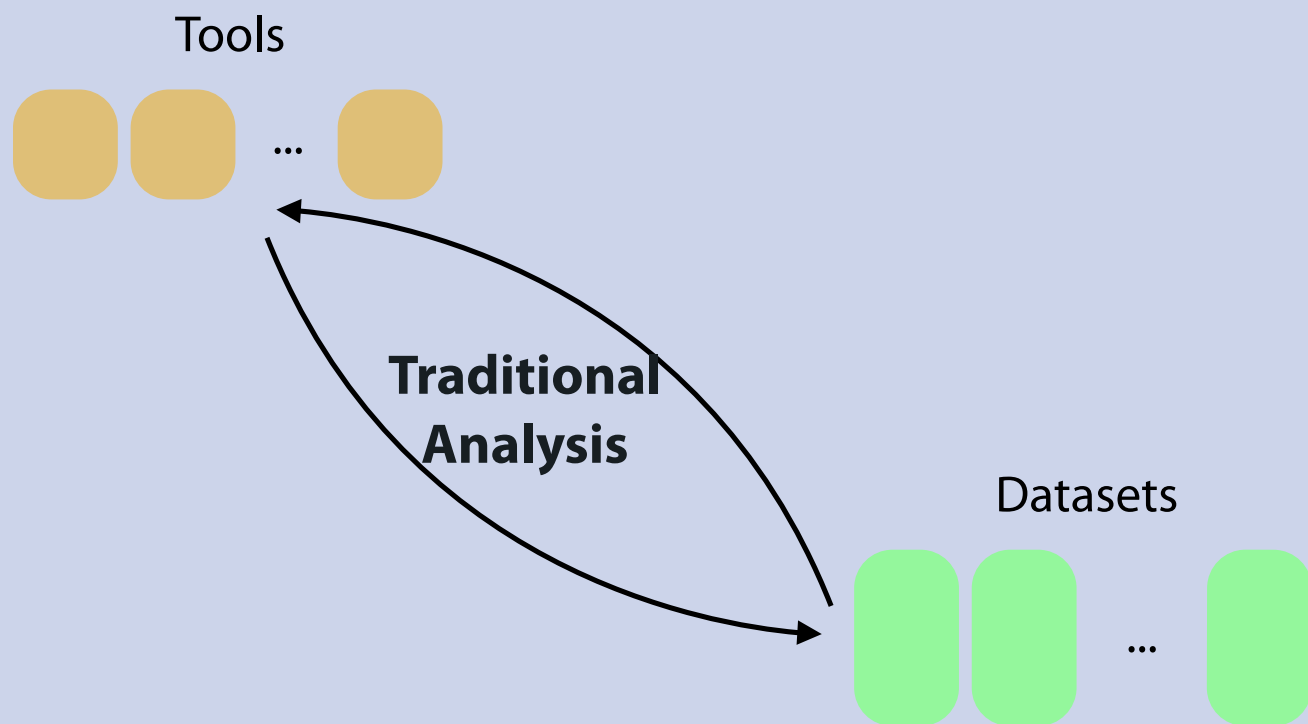**Variant calling from 6 patient, 700GB pileup file requires 48 hours to complete**

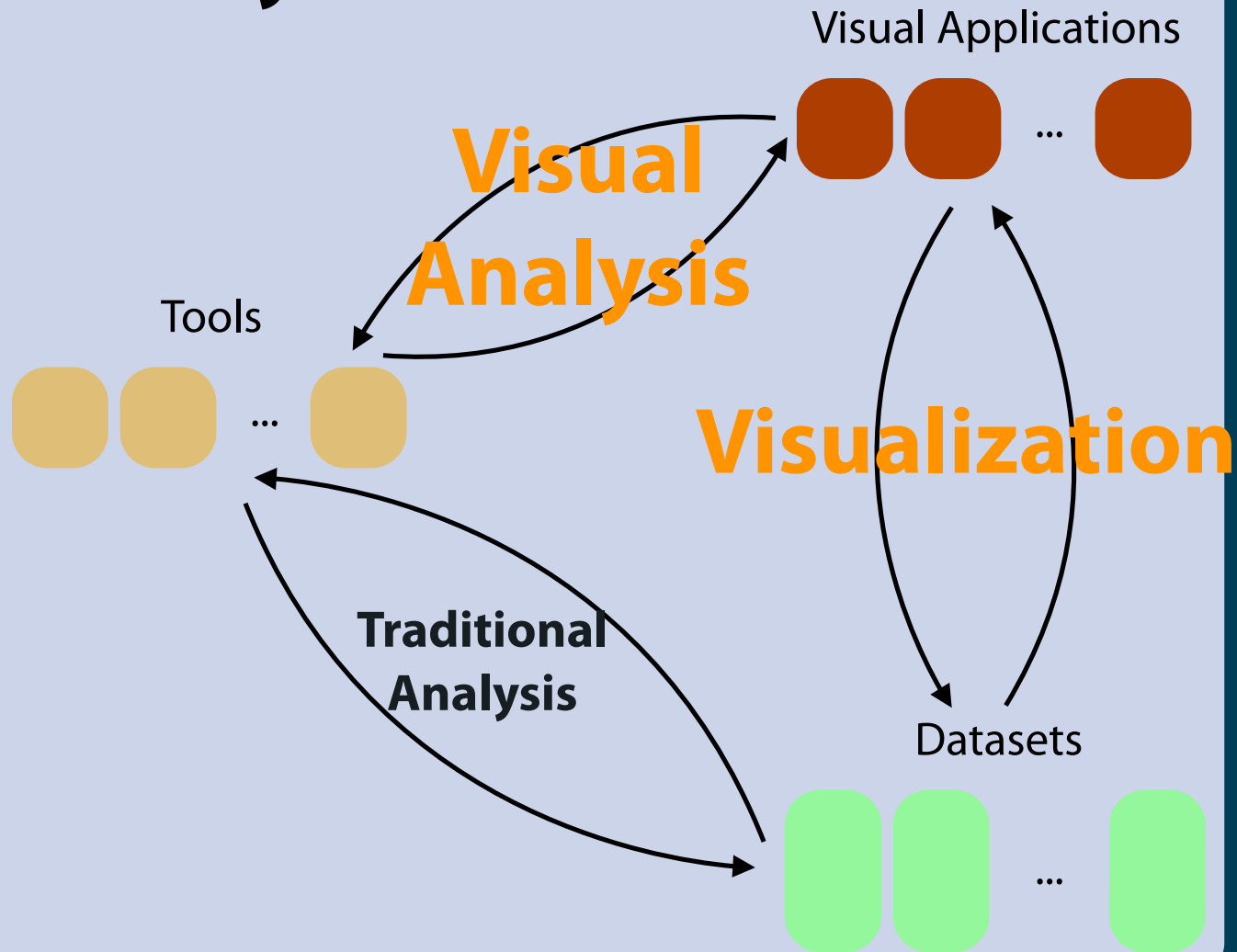# Why Visual Analysis?

Challenges for bioinformatics tools
- ✦ time and compute intensive
- ✦ many parameters
- ✦ parameter sensitive and unpredictable

Approach: repeatedly run tools and visualize outputs to compare and find best parameters/approach

# Galaxy

Tools

...

**Traditional Analysis**

Datasets

...

# Web-based Visualization for High-throughput Genomic Datasets

State-of-the-art data management
- automatic indexing for aggregate data and individual data points
- retrieve data on demand and cache

Render in browser for speed and flexibility

Can share and publish fully-functional visualizations

Goecks et al. (2011) *IEEE BioVis*

# Demo: Visual Analysis

# Real-time Visual Analysis

**Interactive use of production tool to call and visualize variants for multiple patients using parameter sweeps**

A general approach for interactive visual analysis on very large genomics datasets

- ✦ any Galaxy visual application, many tools (original application: transcript assembly)
- ✦ can decide what data to analyze on the fly

Exploration

Automation

Goecks et al. (2012) *Nature Biotechnology*

# Concluding Thoughts

Galaxy is a very useful platform for high-throughout genomics

- ✦ accessible, reproducible, collaborative
- ✦ public, local, cloud

New tools, workflows, and visual analysis tools for analyzing high-throughput cancer sequencing data

- ✦ match patients to drug-profiled cancer cell lines via variants
- ✦ and soon variants + gene expression

Visualization/visual analysis are first-class objects in Galaxy

- ✦ visual analysis affords rapid experimentation with tool parameters
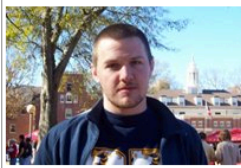
Galaxy

EMORY | WINSHIP CANCER INSTITUTE
A Cancer Center Designated by the National Cancer Institute
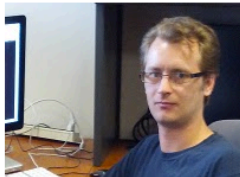
Mike Rossi

Enis Afgan
IRB

Guru Ananda
Penn State

Dannon Baker
Emory

Dan Blankenberg
Penn State
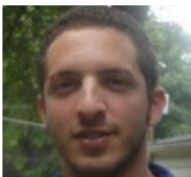
Dave Bouvier
Penn State

Dave Clements
Emory

Nate Coraor
Penn State

Carl Eberhard
Emory

Jeremy Goecks
Emory

Sam Guerler
Emory

Jennifer Hillman Jackson
Penn State

Greg von Kuster
Penn State

Ross Lazarus
BakerIDI

Anton Nekrutenko
Penn State

James Taylor
Emory

EMORY UNIVERSITY

PENNSTATE 1855

genome.gov
National Human Genome Research Institute
National Institutes of Health

NSF National Science Foundation
WHERE DISCOVERIES BEGIN

# Thanks!
# Questions?

http://galaxyproject.org

http://jeremygoecks.com
jeremy.goecks@emory.edu