# Galaxy

## Data intensive biology *for everyone.*

www.galaxyproject.org
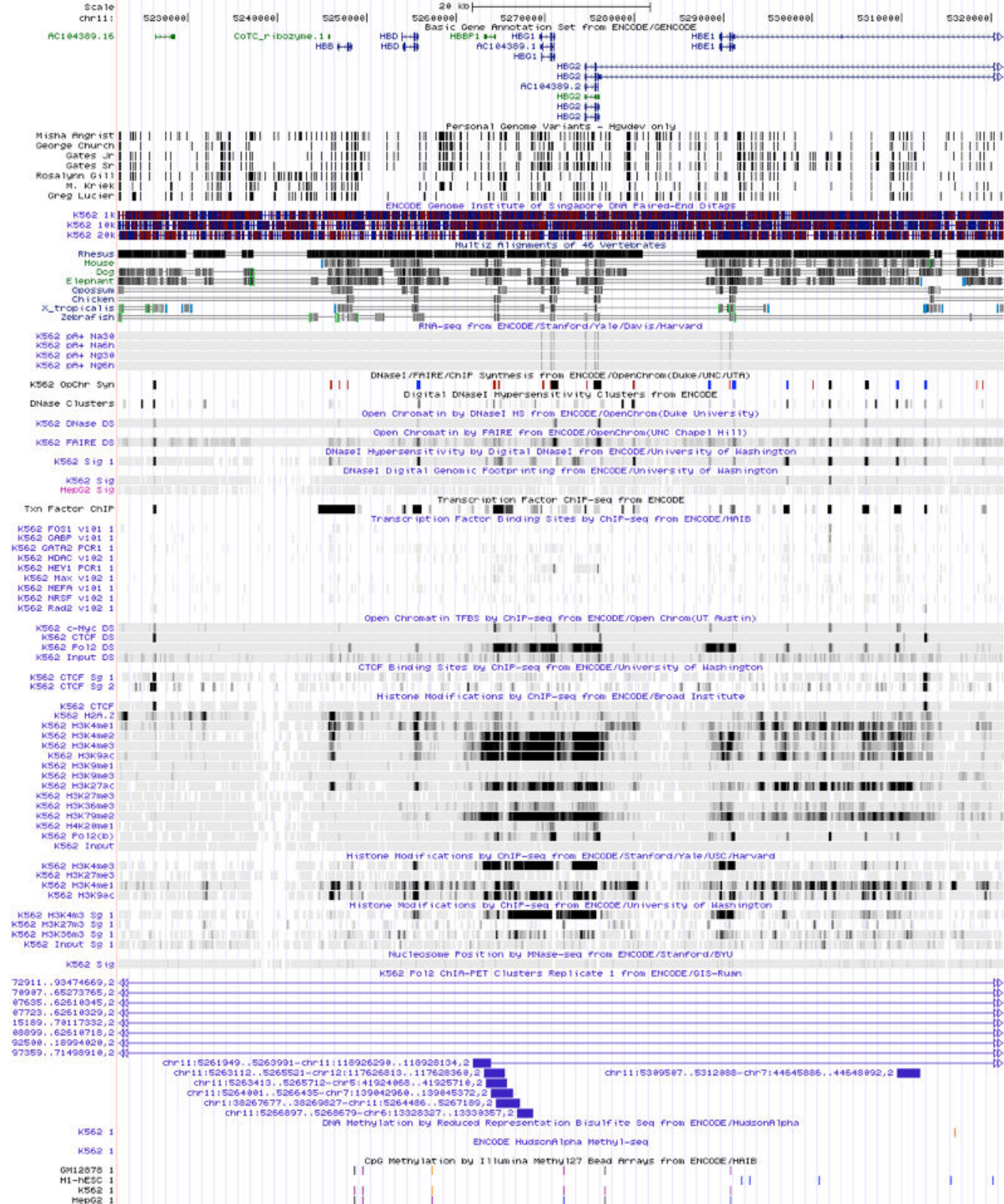
@jxtx / #usegalaxy

High-Throughput
∨
I ❤ SEQUENCING!

# High-throughput sequencing is
## **transformative**

Resequencing

De novo genome sequencing

Direct RNA sequencing

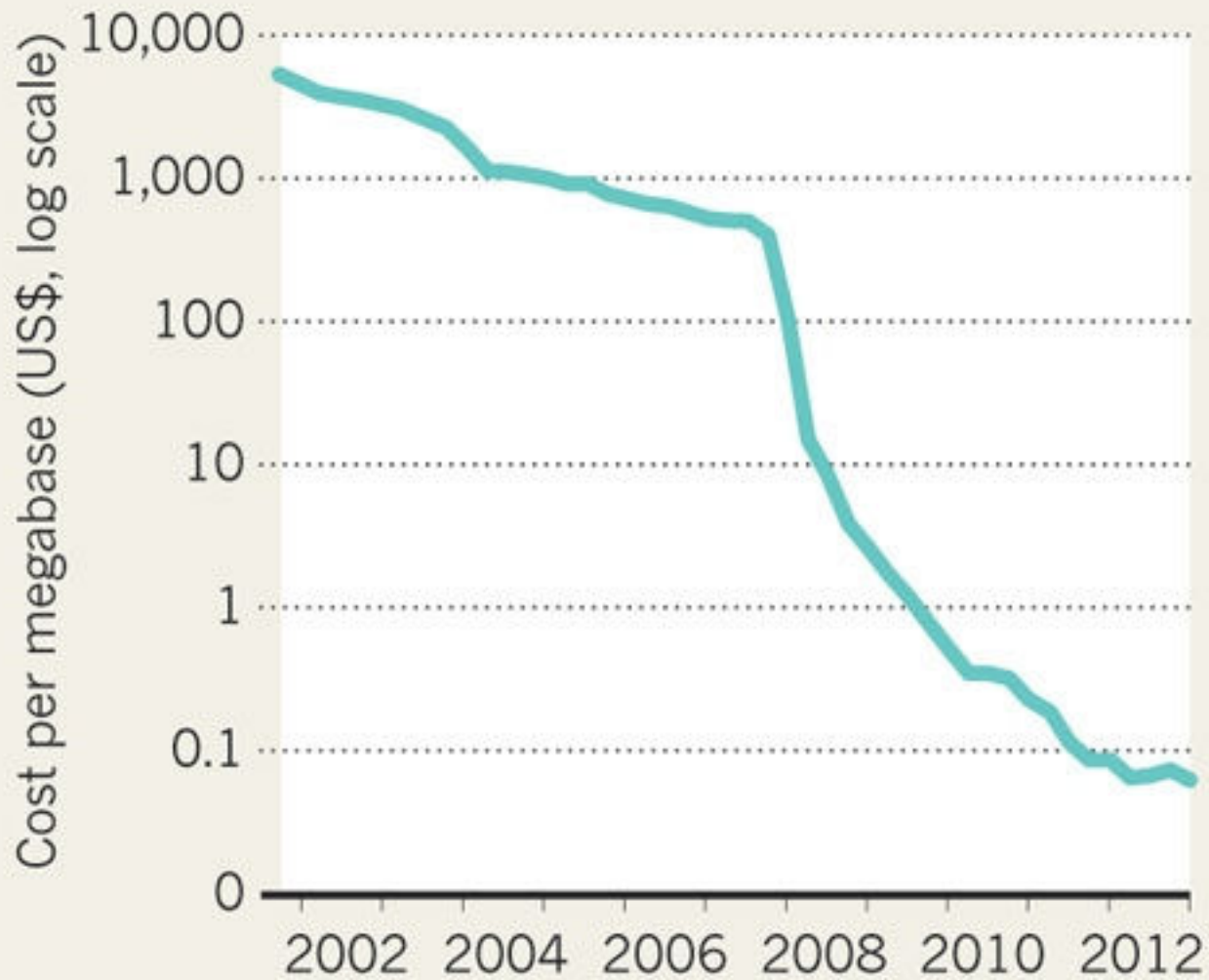Open Chromatin assays (DNase, FAIRE)

Transcription factors (ChIP-seq)

Histones variants (ChIP-seq, MNase-seq)

Long range interactions (5C, Hi-C, ChIA-PET

Methylation (Bisulfite-seq)

# High-throughput sequencing is
# democratizing

# It is widely available…



(http://omicsmaps.com/)

# ...and practically free!



(NHGRI / *Nature* 497:546–547)

**Making sense of this data requires
sophisticated methods**

How can we ensure that these methods are
**accessible** to researchers?

...while also ensuring that scientific results
remain **reproducible**?

# Galaxy: accessible analysis system

**A free (for everyone) web service** integrating a wealth of tools, compute resources, terabytes of reference data and permanent storage

**Open source software** that makes integrating your own tools and data and customizing for your own site simple

**An open extensible platform** for sharing tools, datatypes, workflows, ...

Describe analysis tool behavior abstractly

Analysis environment automatically and transparently tracks details

Workflow system for complex analysis, constructed explicitly or automatically

Pervasive sharing, and publication of documents with integrated analysis

Visualization and visual analytics

The free service is still the easiest way for users with no informatics infrastructure to analyze their data

How can we possibly sustain this?

Registered Users versus Jobs Submitted on Galaxy Main

# usegalaxy.org data growth

+128 cores for NGS/multicore jobs

Data quotas implemented...

New Data per Month (TB)

Nate Coraor

usegalaxy.org frustration growth

Nate Coraor

**How can this possibly scale?**

1. Leverage exisiting public cyber-infrastructure

2. Decentralize, provide many deployment models
(cloud and local — not talking about this today)

Best place to build this robust entry point is clearly a national supercomputing center

The Texas Advanced Computing Center (TACC) has already built substantial infrastructure in the context of the iPlant project

(Including multi petabyte online storage, cloud infrastructure, collocated with some of the worlds largest HPC machines)

However, the iPlant and TACC cyber-infrastructure was underused; thus we established a collaboration

Since October 2013 Galaxy Main has run from TACC

# Transparent Migrations using Galaxy's Hierarchical Object Store

Galaxy Server Processes

Read Data

Write Data

In Corral? — Yes → Corral

No

In Staging? — Yes → Corral Staging

No

In PSU? — Yes → Penn State

No → Object Not Found

Nate Coraor

# Expanding to more XSEDE resources

Galaxy can already run jobs on almost any batch system, but most XSEDE resources do not provide direct access for job submission…

# Pulsar

Galaxy job runner that can
run almost anywhere

No shared filesystem, stages all necessary
Galaxy components

John Chilton

Galaxy Server VMs (TACC)

Galaxy Server Processes

Messaging Server

Job control (AMQP)

Data transfer (HTTPS)

Data transfer (HTTPS)

Blacklight (PSC)

Pulsar

Stampede (TACC)

Pulsar

Nate Coraor

# Moving long running jobs out to XSEDE

- Problem:
  - Jobs wait in the queue for a long time
  - Jobs may fail immediately upon run due to bad parameters
  - Most jobs run quickly! Can we relocate the long ones?

- Goals:
  - Shorten wait from submission to start
  - Allow testing params without waiting

- Solutions:
  - Set a short walltime, resubmit jobs to bigger resources (new code)
  - User selection of resources (Stampede - longer wait to start, but more concurrent jobs allowed)
  - Create "development" queues w/ short walltime

Nate Coraor and John Chilton

# State of Affairs

- Today
  - Galaxy Test jobs to Stampede and Blacklight
  - Galaxy Main jobs to Stampede

- Up next
  - Galaxy Main jobs to Blacklight
  - Optimize Trinity tools for Blacklight
  - Linking XSEDE allocations to Galaxy accounts

# Credits

- **Texas Advanced Computing Center**
  - Dan Stanzione
  - Matt Vaughn
  - Chris Jordan
  - Mike Packard
  - Nathaniel Mendoza

- **iPlant Collaborative**
  - Stephen Goff

- **Pittsburgh Supercomputing Center**
  - Philip Blood
  - Kathy Benninger
  - Robert Budden
  - Jared Yanovich
  - Josephine Palencia
  - J. Ray Scott
  - Joe Lappa

## ... and the Galaxy Team and community
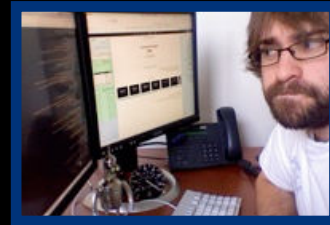
# Engineering



Enis Afgan

Dannon Baker

Dan Blankenberg

Dave Bouvier

Nate Coraor

Martin Čech

John Chilton

Carl Eberhard

Sam Guerler

Nick Stoler

# Support and outreach

Dave Clements

Jennifer Jackson

# Leadership

James Taylor

Anton Nekrutenko

Jeremy Goecks