# Accessible, transparent, reproducible analysis with **Galaxy**

Try it now:
http://usegalaxy.org

Develop and deploy:
http://getgalaxy.org

Try it now:
http://usegalaxy.org

Develop and deploy:
http://getgalaxy.org

EMORY

PENNSTATE. 1855

Enis Afgan

Dannon Baker

Jeremy Goecks

Dan Blankenberg

Nate Coraor

Jennifer Jackson

Dave Clements

Kanwei Li

James Taylor

Greg von Kuster

Kelly Vincent

Anton Nekrutenko

Biology *has been* rapidly transformed into a data intensive science

**Illumina Hi-Seq:** ~25-50 GB per day, $16k-$20k per run Greater than 1Mb per dollar With multiplexing, as little as $100 per sample.
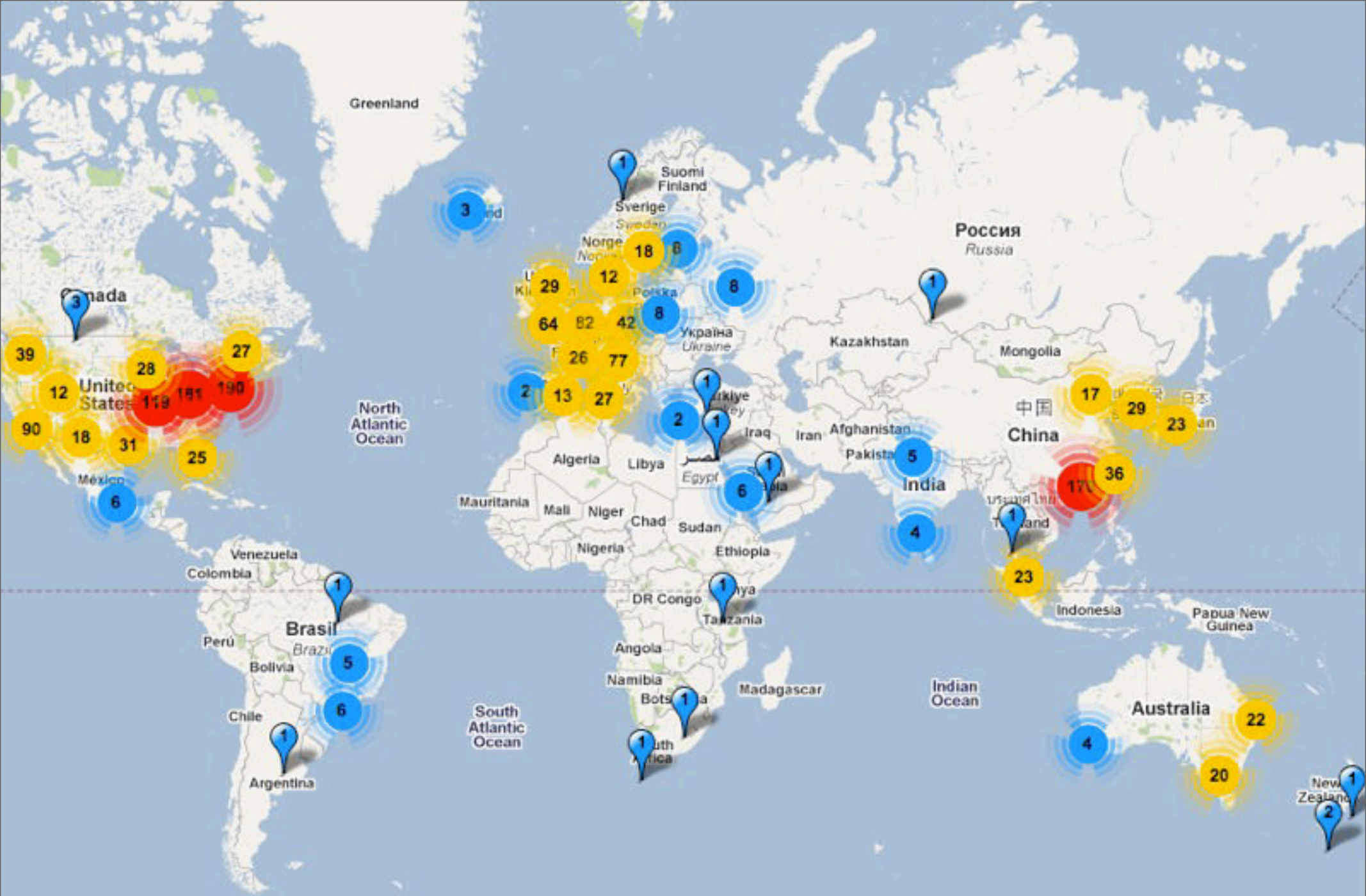
**454 GS / Junior:** 40-400Mb runs, but read lengths pushing 1kb

**Ion Torrent PGM:** 10Mb-1Gb runs, 200-400bp reads, 2 hour runtime, $500!

**PacBio RS:** Direct single molecule sequencing, only 35k reads, but long read lengths, 30 minute runs!
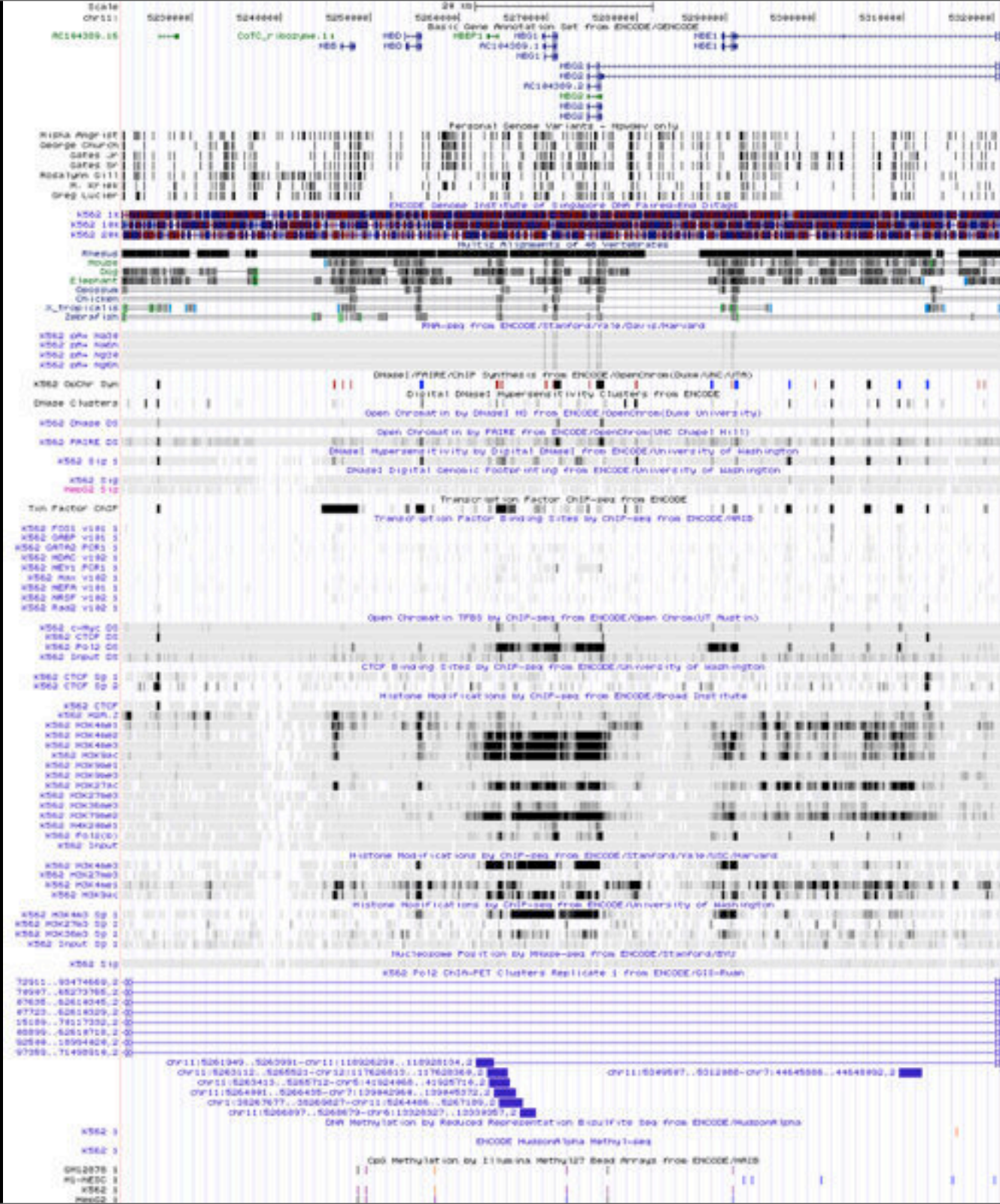
(http://pathogenomics.bham.ac.uk/hts/)

We can turn many **functional annotation** problems into **sequencing** problems
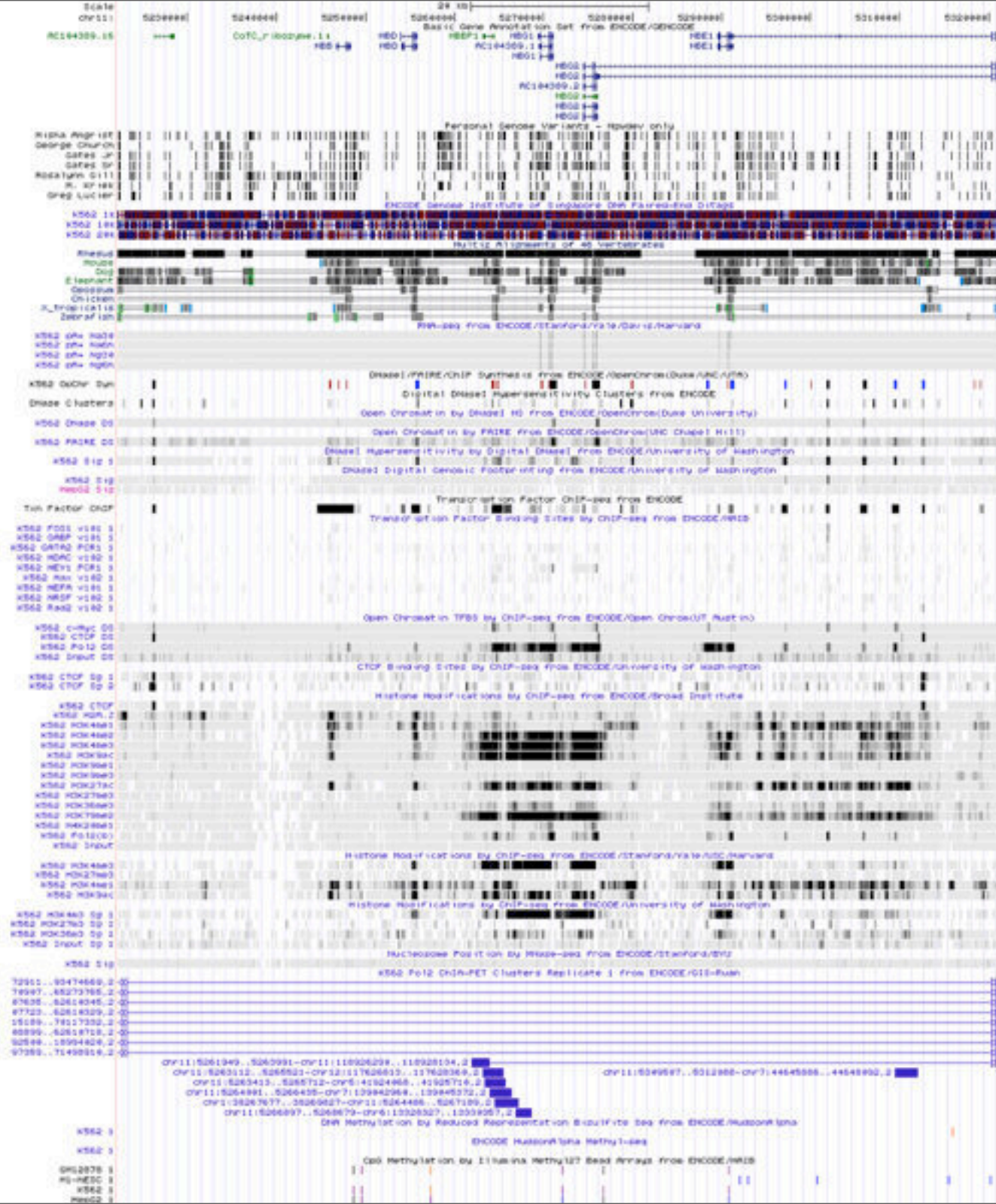
An individual genome is relatively static

Transcript levels, epigenomic modifications, and chromatin structure vary based on cell type, time, condition, ...
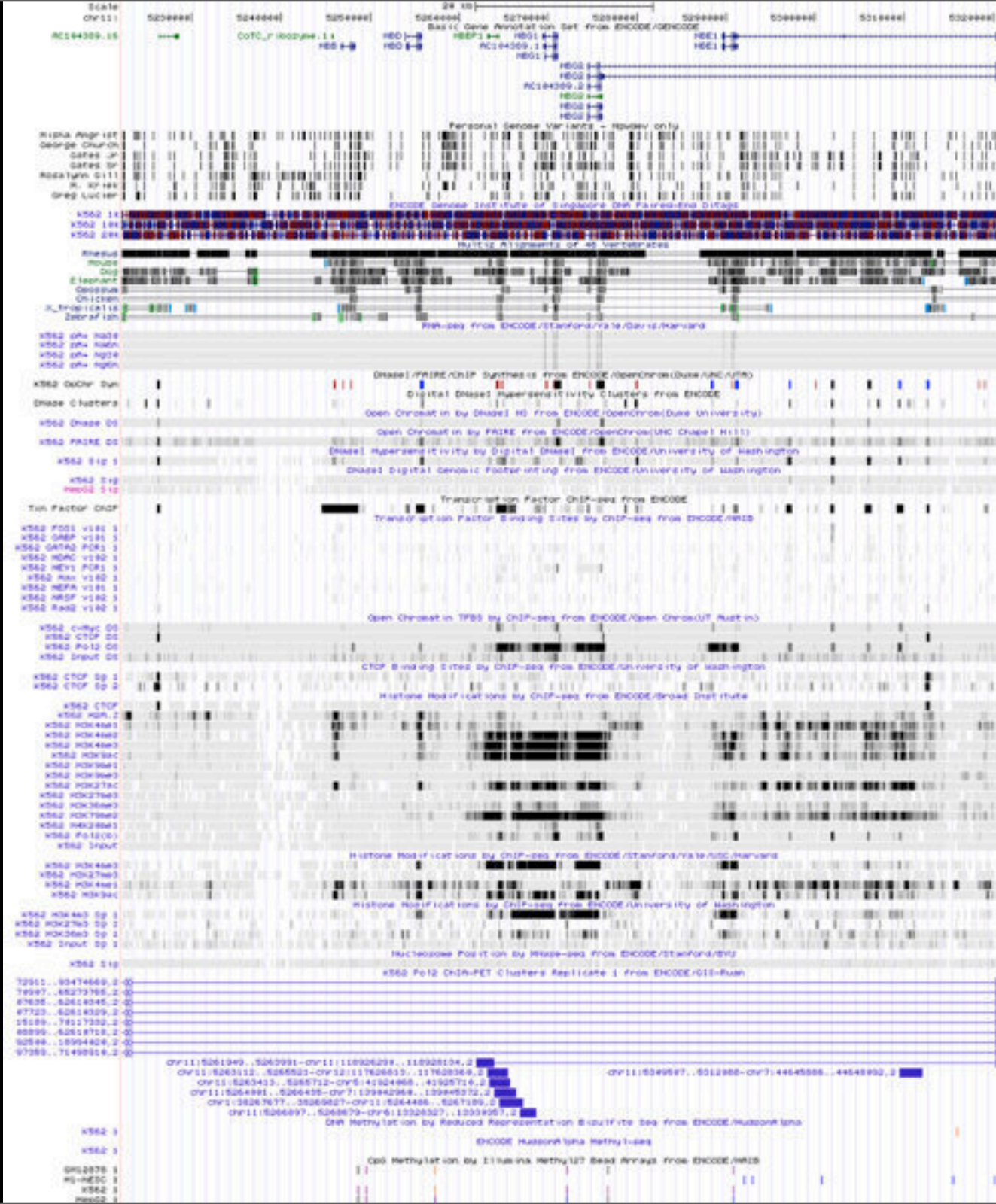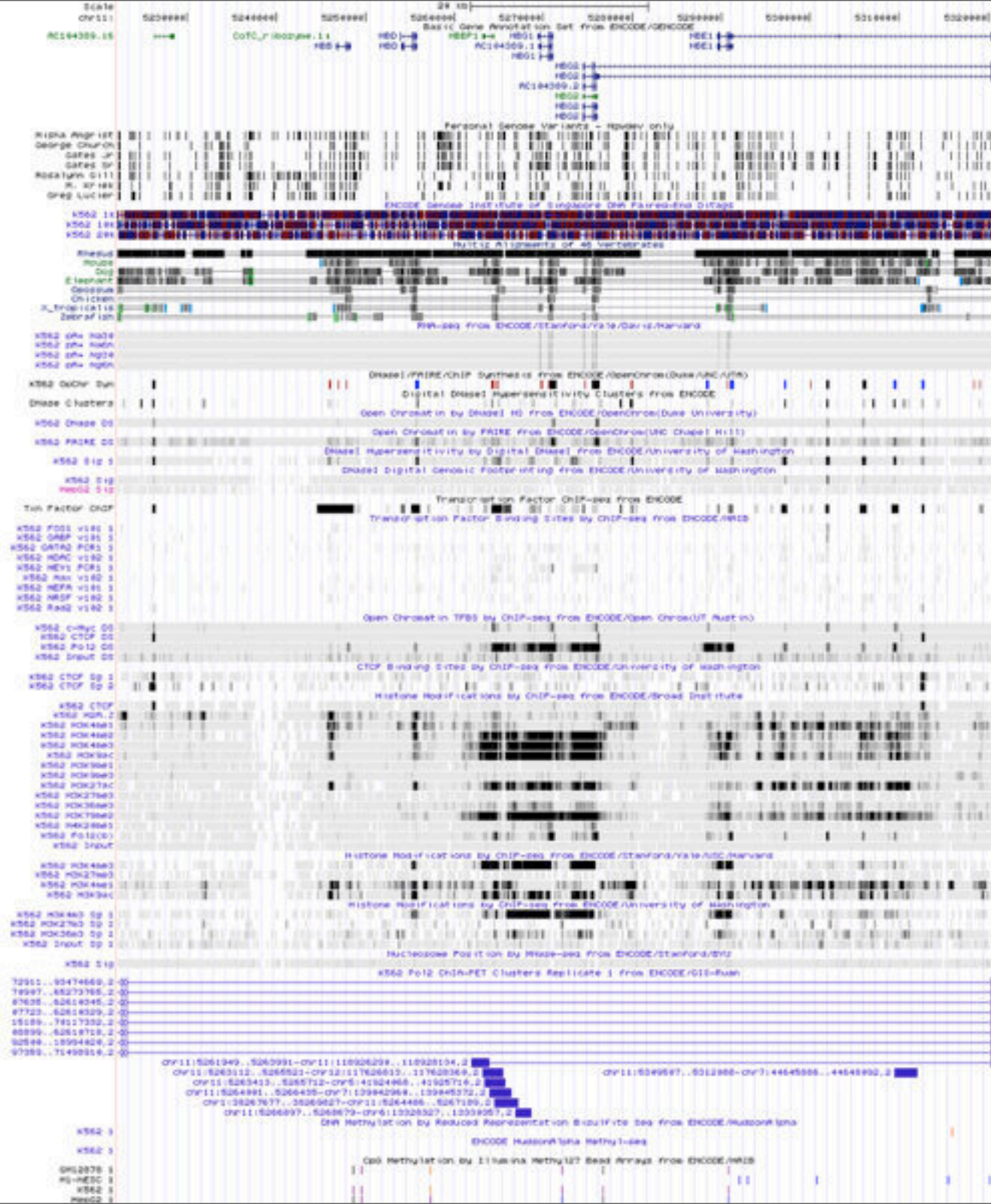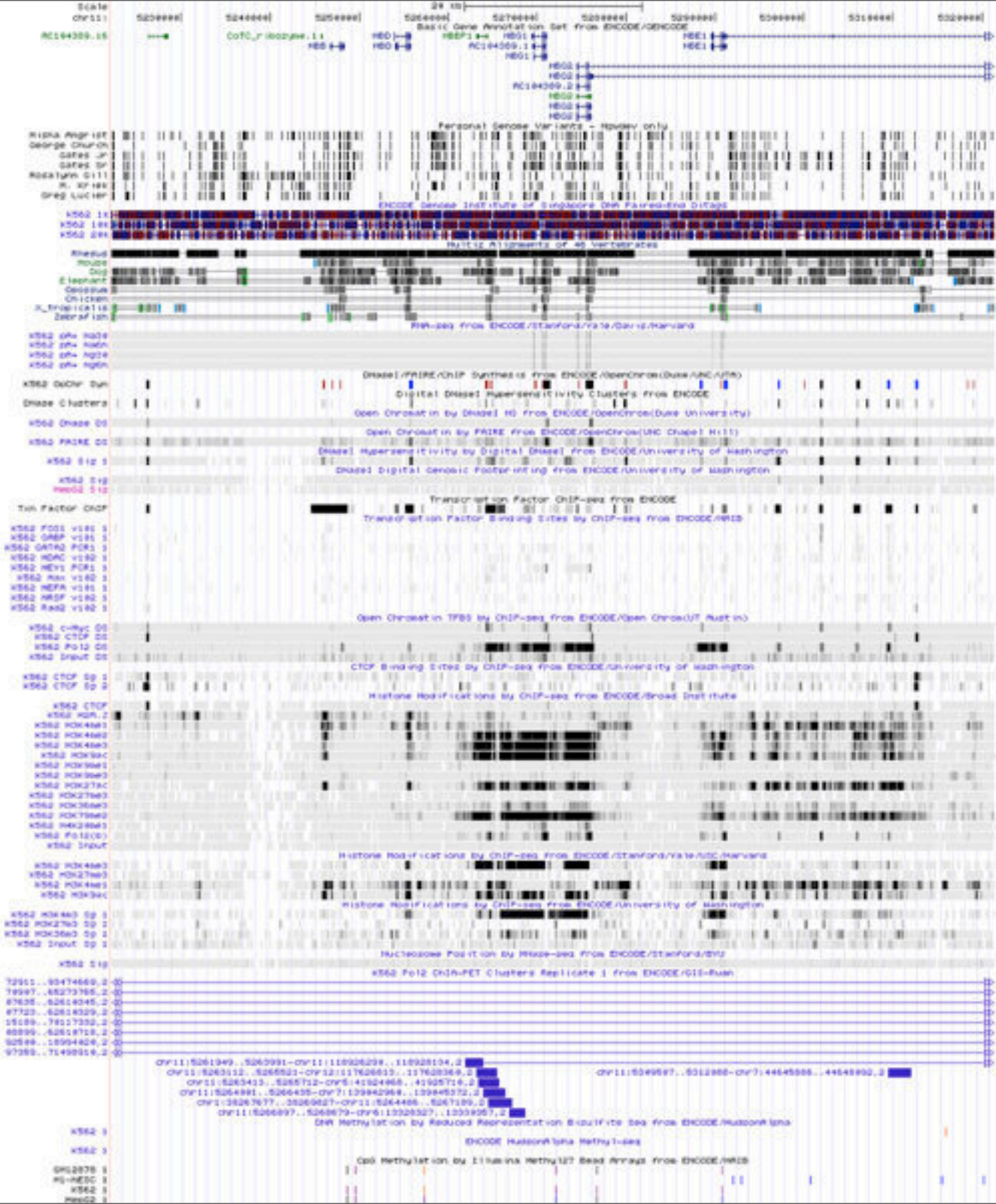
Enormous potential for data generation

Resequencing

De-novo genome sequening

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)

Investigators across nearly all areas of biology can take advantage of these techniques

Investigator driven data production replacing large community data production projects

This "**democratization of sequencing**" has not yet been matched by democratization of analysis infrastructure, burden is largely on the investigator

However, making sense of this data *requires* sophisticated methods

How can these methods be made **accessible** to scientists?

How do we facilitate **transparent** communication of analyses?

How do we ensure that analyses are **reproducible**?

# A crisis in genomics research:
# reproducibility

# Microarray Experiment Reproducibility

- 18 Nat. Genetics microarray gene expression experiments

- **Less than 50% reproducible**

- Problems

  - missing data (38%)

  - missing software, hardware details (50%)

  - missing method, processing details (66%)

*Ioannidis, J.P.A. et al. Repeatability of published microarray gene expression analyses. Nat Genet 41, 149-155 (2009)*

# NGS Re-sequencing Experiment Reproducibility

- 14 re-sequencing experiments in Nat. Genetics, Nature, and Science (2010)

- **0% reproducible?**

- Problems

  - limited access to primary data (50%)

  - some or all tools unavailable (50%)

  - settings & versions not provided (100%)

# Galaxy: accessible analysis system

# What is Galaxy?

- **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

# Integrating existing tools into a uniform framework



- Defined in terms of an abstract interface (inputs and outputs)

  - In practice, mostly command line tools, a declarative XML description of the interface, how to generate a command line

- Designed to be as easy as possible for tool authors, while still allowing rigorous reasoning

## Left panel (Galaxy tool form)

**Cluster**

Cluster intervals of: `6: UCSC Main on Human: knownGene ▾`

max distance between intervals: `1`
(bp)

min number of intervals per cluster: `2`

Return type: `Merge clusters into single intervals ▾`

[ Execute ]

ℹ **TIP:** If your query does not appear in the pulldown menu -> it is not in interval format. Use "edit attributes" to set chromosome, start, end, and strand columns.
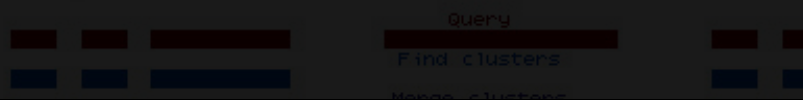
---

**Screencasts!**

See Galaxy Interval Operation Screencasts (right click to open this link in another window).

---

**Syntax**

- **Maximum distance** is greatest distance in base pairs allowed between intervals that will be considered "clustered". **Negative** values for distance are allowed, and are useful for clustering intervals that overlap.
- **Minimum intervals per cluster** allow a threshold to be set on the minimum number of intervals to be considered a cluster. Any area with less than this minimum will not be included in the ouput.
- **Merge clusters into single intervals** outputs intervals that span the entire cluster.
- **Find cluster intervals; preserve comments and order** filters out non-cluster intervals while maintaining the original ordering and comments in the file.
- **Find cluster intervals; output grouped by clusters** filters out non-cluster intervals, but outputs the cluster intervals so that they are grouped together. Comments and original ordering in the file are lost.

---

**Example**



## Right panel (cluster.xml)

```
1  <tool id="gops_cluster_1" name="Cluster">
2    <description>[[Cluster]] the intervals of a query</description>
3    <command interpreter="python2.4">
4      gops_cluster.py $input1 $output -1 $input1_chromCol,$input1_startC
5                      -d $distance -m $minregions -o $returntype
6    </command>
7    <inputs>
8      <param format="interval" name="input1" type="data">
9        <label>Cluster intervals of</label>
10     </param>
11     <param name="distance" size="5" type="integer" value="1" help="(bp
12       <label>max distance between intervals</label>
13     </param>
14     <param name="minregions" size="5" type="integer" value="2">
15       <label>min number of intervals per cluster</label>
16     </param>
17     <param name="returntype" type="select" label="Return type">
18       <option value="1">Merge clusters into single intervals</option>
19       <option value="2">Find cluster intervals; preserve comments and
20       <option value="3">Find cluster intervals; output grouped by clus
21       <option value="4">Find the smallest interval in each cluster</op
22       <option value="5">Find the largest interval in each cluster</opt
23     </param>
24   </inputs>
25   <help>
26
27 .. class:: infomark
28
29 **TIP:** If your query does not appear in the pulldown menu -> it is n
30
31 -----
32
33 **Screencasts!**
34
35 See Galaxy Interval Operation Screencasts  (right click to open this l
36
37 .. _Screencasts: http://www.bx.psu.edu/cgi-bin/trac.cgi/wiki/GopsDesc
38
39 -----
40
41 **Syntax**
42
43 - **Maximum distance** is greatest distance in base pairs allowed betw
44 - **Minimum intervals per cluster** allow a threshold to be set on the
45 - **Merge clusters into single intervals** outputs intervals that span
46 - **Find cluster intervals; preserve comments and order** filters out
47 - **Find cluster intervals; output grouped by clusters** filters out n
```

Line: 87   Column: 8   XML   Soft Tabs: 2

HTML inputs generated from abstract parameter description

Template for generating command line from parameter values

```xml
<tool id="gops_cluster_1" name="Cluster">
  <description>[[Cluster]] the intervals of a query</description>
  <command interpreter="python2.4">
    gops_cluster.py $input1 $output -1 $input1_chromCol,$input1_startCol,$input1_endCol
                    -d $distance -m $minregions -o $returntype
  </command>
  <inputs>
    <param format="interval" name="input1" type="data">
      <label>Cluster intervals of</label>
    </param>
    <param name="distance" size="5" type="integer" value="1" help="(bp)">
      <label>max distance between intervals</label>
    </param>
    <param name="minregions" size="5" type="integer" value="2">
      <label>min number of intervals per cluster</label>
    </param>
    <param name="returntype" type="select" label="Return type">
      <option value="1">Merge clusters into single intervals</option>
      <option value="2">Find cluster intervals; preserve comments and order</option>
      <option value="3">Find cluster intervals; output grouped by clusters</option>
      <option value="4">Find the smallest interval in each cluster</option>
      <option value="5">Find the largest interval in each cluster</option>
    </param>
  </inputs>
  <help>

.. class:: infomark

**TIP:** If your query does not appear in the pulldown menu -> it is not in interval fo

.....

**Screencasts!**

See Galaxy Interval Operation Screencasts_ (right click to open this link in another wi

.. _Screencasts: http://www.bx.psu.edu/cgi-bin/trac.cgi/wiki/GopsDesc

.....

**Syntax**

- **Maximum distance** is greatest distance in base pairs allowed between intervals tha
- **Minimum intervals per cluster** allow a threshold to be set on the minimum number o
- **Merge clusters into single intervals** outputs intervals that span the entire clust
- **Find cluster intervals; preserve comments and order** filters out non-cluster inter
- **Find cluster intervals; output grouped by clusters** filters out non-cluster interv
```

Line: 61   Column: 45   XML   Soft Tabs: 2

```
41 **Syntax**
42
43 - **Maximum distance** is greatest distance in base pairs allowed between intervals tha
44 - **Minimum intervals per cluster** allow a threshold to be set on the minimum number o
45 - **Merge clusters into single intervals** outputs intervals that span the entire clust
46 - **Find cluster intervals; preserve comments and order** filters out non-cluster inter
47 - **Find cluster intervals; output grouped by clusters** filters out non-cluster interv
48
49 -----
50
51 **Example**
52
53 .. image:: ../static/operation_icons/gops_cluster.gif
54
55 </help>
56
57   <outputs>
58     <data format="input" name="output" metadata_source="input1" />
59   </outputs>
60   <code file="operation_filter.py">
61     <hook exec_after_process="exec_after_cluster" />
62   </code>
63   <tests>
64     <test>
65       <param name="input1" value="1.bed" />
66       <param name="distance" value="1" />
67       <param name="minregions" value="2" />
68       <param name="returntype" value="1" />
69       <output name="output" file="gops-cluster-1.dat" />
70     </test>
71     <test>
72       <param name="input1" value="1.bed" />
73       <param name="distance" value="1" />
74       <param name="minregions" value="2" />
75       <param name="returntype" value="2" />
76       <output name="output" file="gops-cluster-2.dat" />
77     </test>
78     <test>
79       <param name="input1" value="1.bed" />
80       <param name="distance" value="1" />
81       <param name="minregions" value="2" />
82       <param name="returntype" value="3" />
83       <output name="output" file="gops-cluster-3.dat" />
84     </test>
85   </tests>
86
87 </tool>
```

cluster.xml

Line: 61   Column: 45   XML   Soft Tabs: 2   —

Functional tests to be run with the "full stack" in place

Much more complex interfaces can be defined

Repeating groups of parameters

Conditional groups, grouping constructs can be nested

```
build_ucsc_custom_track.xml

1   <tool id="build_ucsc_custom_track_1" name="Build custom track">
2     <description>for UCSC genome browser</description>
3     <command interpreter="python2.4">
4       build_ucsc_custom_track.py
5         "$out_file1"
6         #for $t in $tracks
7           "${t.input.file_name}"
8           "${t.input.ext}"
9           #if $t.input.ext == "interval"
10            ${t.input.metadata.chromCol},${t.input.metadata.startCol},${t.input.metadata.endCol},${t.input.metadata.strandCol}
11          #else
12            "NA"
13          #end if
14          "${t.name}"
15          "${t.description}"
16          "${t.color}"
17          "${t.visibility}"
18        #end for
19    </command>
20    <inputs>
21      <repeat name="tracks" title="Track">
22        <param name="input" type="data" format="interval,wig" label="Dataset"/>
23        <param name="name" type="text" size="15" value="User Track">
24          <validator type="length" max="15"/>
25        </param>
26        <param name="description" type="text" value="User Supplied Track (from Galaxy)">
27          <validator type="length" max="60"/>
28        </param>
29        <param label="Color" name="color" type="select">
30          <option selected="yes" value="0-0-0">Black</option>
31          <option value="255-0-0">Red</option>
32          <option value="0-255-0">Green</option>
33          <option value="0-0-255">Blue</option>
34          <option value="255-0-255">Magenta</option>
35          <option value="0-255-255">Cyan</option>
36          <option value="255-215-0">Gold</option>
37          <option value="160-32-240">Purple</option>
38          <option value="255-140-0">Orange</option>
39          <option value="255-20-147">Pink</option>
40          <option value="92-51-23">Dark Chocolate</option>
41          <option value="85-107-47">Olive green</option>
42        </param>
```

Template language for building complex command lines

```
 70      </conditional>
 71     </repeat>
 72   </inputs>
 73
 74   <configfiles>
 75     <configfile name="script_file">
 76       ## Setup R error handling to go to stderr
 77       options( show.error.messages=F,
 78               error = function () { cat( geterrmessage(), file=stderr() ); q( "no", 1, F ) } )
 79       ## Determine range of all series in the plot
 80       xrange = c( NULL, NULL )
 81       yrange = c( NULL, NULL )
 82       #for $i, $s in enumerate( $series )
 83         s${i} = read.table( "${s.input.file_name}" )
 84         x${i} = s${i}[,${s.xcol}]
 85         y${i} = s${i}[,${s.ycol}]
 86         xrange = range( x${i}, xrange )
 87         yrange = range( y${i}, yrange )
 88       #end for
 89       ## Open output PDF file
 90       pdf( "${out_file1}" )
 91       ## Dummy plot for axis / labels
 92       plot( NULL, type="n", xlim=xrange, ylim=yrange, main="${main}", xlab="${xlab}", ylab="${ylab}" )
 93       ## Plot each series
 94       #for $i, $s in enumerate( $series )
 95         #if $s.series_type['type'] == "line"
 96           lines( x${i}, y${i}, lty=${s.series_type.lty}, lwd=${s.series_type.lwd}, col=${s.series_type.col} )
 97         #elif $s.series_type.type == "points"
 98           points( x${i}, y${i}, pch=${s.series_type.pch}, cex=${s.series_type.cex}, col=${s.series_type.col} )
 99         #end if
100       #end for
101       ## Close the PDF file
102       devname = dev.off()
103     </configfile>
104   </configfiles>
105
106   <outputs>
107     <data format="pdf" name="out_file1" />
108   </outputs>
109
110   <help>
111   .. class:: infomark
```

Or additional configuration files, scripts, …

As data sizes grow, increasingly important
to be able to express within tool parallelism

Naturally parallel (split/join) constructs can be
specified in configuration

Parallel environments (MPI) can be used, but
management delegated to underlying resources

Ongoing work to support more
complex scenarios

# Customization extends beyond tools

- Everything in the Galaxy framework is either configuration driven or pluggable (or both)

- Tools conventionally extended through configuration, but new tool types can be added

- Datatypes added through configuration, or plugin classes for advanced functionality

- *Nothing inherently specific to genomics!*

**NGS: QC and manipulation**

ILLUMINA DATA

- FASTQ Groomer convert between various FASTQ quality formats
- FASTQ splitter on joined paired end reads
- FASTQ joiner on paired end reads
- FASTQ Summary Statistics by column

ROCHE-454 DATA

- Build base quality distribution
- Select high quality segments
- Combine FASTA and QUAL into FASTQ

AB-SOLID DATA

- Convert SOLiD output to fastq
- Compute quality statistics for SOLiD data
- Draw quality score boxplot for SOLiD data

GENERIC FASTQ MANIPULATION

- Filter FASTQ reads by quality score and length
- FASTQ Trimmer by column
- FASTQ Quality Trimmer by sliding window

---

Evolution
**Metagenomic analyses**
**Human Genome Variation**
**EMBOSS**

NGS TOOLBOX BETA

**NGS: QC and manipulation**
**NGS: Mapping**

ILLUMINA

- Map with Bowtie for Illumina
- Map with BWA for Illumina

ROCHE-454

- Lastz map short reads against reference sequence
- Megablast compare short reads against htgs, nt, and wgs databases
- Parse blast XML output

AB-SOLID

- Map with Bowtie for SOLiD

**NGS: SAM Tools**
**NGS: Indel Analysis**
**NGS: Peak Calling**
**NGS: RNA Analysis**

RGENETICS

**SNP/WGA: Data; Filters**
**SNP/WGA: QC; LD; Plots**
**SNP/WGA: Statistical Models**

Workflows

---

NGS TOOLBOX BETA

**NGS: QC and manipulation**
**NGS: Mapping**
**NGS: SAM Tools**

- Filter SAM on bitwise flag values
- Convert SAM to interval
- SAM-to-BAM converts SAM format to BAM format
- BAM-to-SAM converts BAM format to SAM format
- Merge BAM Files merges BAM files together
- Generate pileup from BAM dataset
- Filter pileup on coverage and SNPs
- Pileup-to-Interval condenses pileup format into ranges of bases
- flagstat provides simple stats on BAM files

**NGS: Indel Analysis**
**NGS: Peak Calling**
**NGS: RNA Analysis**

RGENETICS

**SNP/WGA: Data; Filters**
**SNP/WGA: QC; LD; Plots**
**SNP/WGA: Statistical Models**

Workflows

---

**NGS: SAM Tools**
**NGS: Indel Analysis**

- Filter Indels for SAM
- Extract indels from SAM
- Indel Analysis

**NGS: Peak Calling**

- MACS Model-based Analysis of ChIP-Seq
- GeneTrack indexer on a BED file
- Peak predictor on GeneTrack index

**NGS: RNA Analysis**

RNA-SEQ

- Tophat Find splice junctions using RNA-seq data
- Cufflinks transcript assembly and FPKM (RPKM) estimates for RNA-Seq data
- Cuffcompare compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments
- Cuffdiff find significant changes in transcript expression, splicing, and promoter use

FILTERING

- Filter Combined Transcripts using tracking file

Dozens of tools for different NGS applications packaged with Galaxy

# Analysis environment

# Galaxy analysis interface



- Consistent tool user interfaces automatically generated

- History system facilitates and tracks multistep analyses

# Automatically and transparently tracks
# every step of every analysis

# As well as user-generated metadata and annotation...

# Workflows

# Galaxy workflow system



- Workflows can be constructed from scratch *or* extracted from existing analysis histories

- Facilitate reuse, as well as providing precise reproducibility of a complex analysis

**Example**: Workflow for differential expression analysis of RNA-seq using Tophat/Cufflinks tools

**Example**: Diagnosing low-frequency heterosplasmic sites in two tissues from the same individual

# Galaxy deployment models

# Galaxy main site
## (http://usegalaxy.org)

- Public web site, anybody can use

- ~500 new users per month, ~100 TB of user data, ~130,000 analysis jobs per month, every month is our busiest month ever...

- Will continue to be maintained and enhanced, but with limits and quotas

- Centralized solution cannot scale to meet data analysis demands

# Local Galaxy instances
# (http://getgalaxy.org)

- Galaxy is designed for local installation and customization

  - Just download and run, completely self-contained

  - Easily integrate new tools

  - Easy to deploy and manage on nearly any (unix) system

  - Run jobs on existing compute clusters

# Scale up on existing resources

- Move intensive processing (tool execution) to other hosts

- Frees up the application server to serve requests and manage jobs

- Utilize existing resources

- Supports any batch scheduler that supports DRMAA (most of them)

- All levels of job running and scheduling are pluggable

# Galaxy Cloud
# (http://usegalaxy.org/cloud)

- On-demand resource acquisition fits well with the irregular resource needs of many labs working with sequence data

- Our goal is to approach the ease of use of a "software as a service" solution while maintaining the flexibility and control of an infrastructure based solution

# Using Amazon EC2: Startup in 3 steps

Can use like any other Galaxy instance, with additional compute nodes acquired and released (*automatically*) in response to usage

Share a snapshot of this instance

# Galaxy Cloudman

## Galaxy Cloudman Console

Welcome to Galaxy Cloudm... ...loudMan. Your previous
data store has been recon... ...d and remove 'worker'
nodes for running jobs.

Terminate clust...    ...Access Galaxy

## Status

**Cluster name:**    jame...

**Disk status:**    181...

**Worker status:**    **Idle**

**Service status:**    Appli...

**External Logs:**    Gala...

Cluster status log...

Autoscaling is off.
Turn on?

### Currently shared instances

### Share-an-instance

**This form allows you to share this cluster
instance, at its current state, with others.** You can
make the instance public or share it with specific users
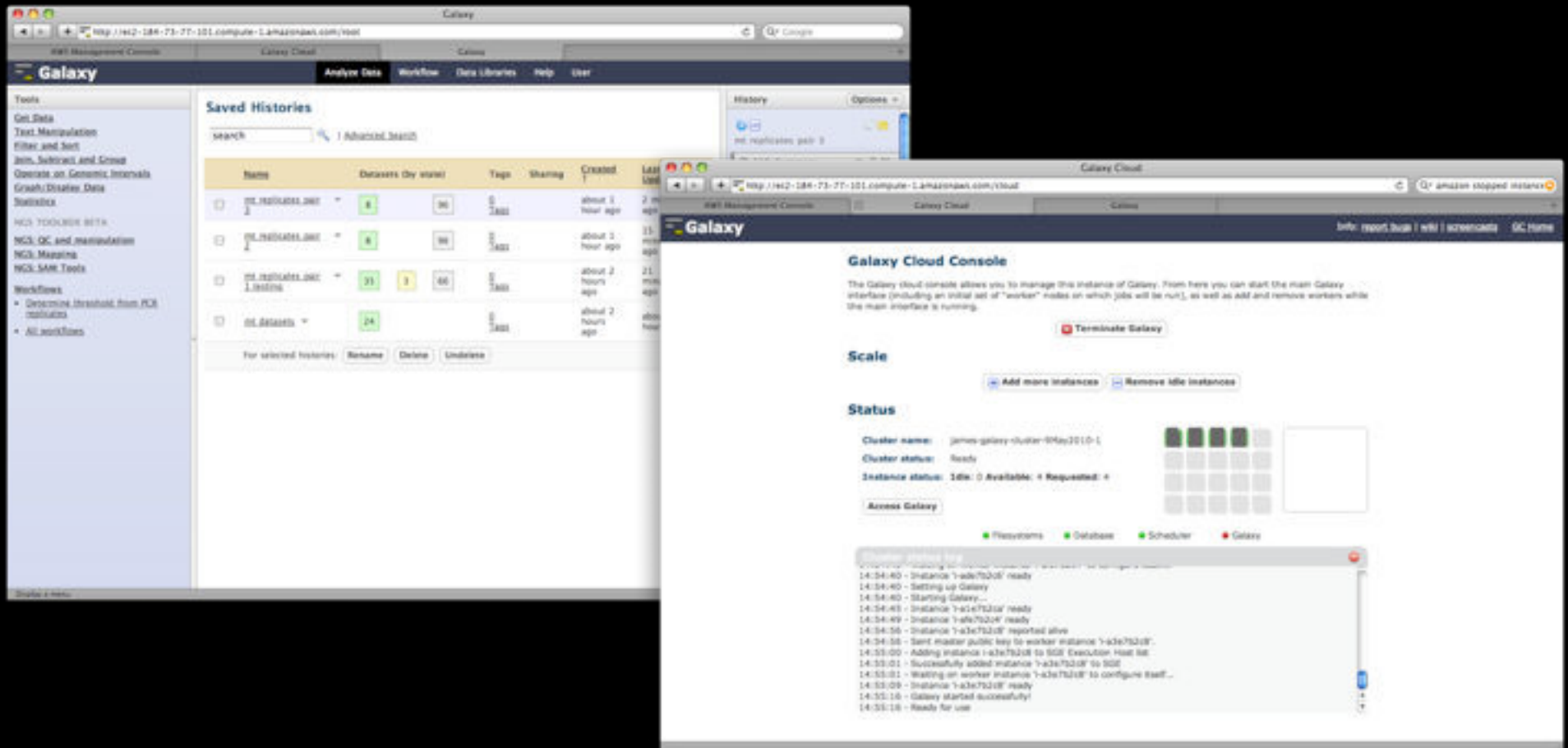by providing their account information below.
You may also share the instance with yourself by
specifying your own credentials, which will have the
effect of saving the instance at its current state.

**While setting up an instance to be shared, all
currently running cluster services will be stopped.**
Then, a snapshot of your data volume and a folder in
your cluster's bucket will be created (under
'shared/[current date and time]); this folder will contain
your cluster's current configuration. The created
snapshot and the folder will be given READ permissions
to the users you choose (or make it public). This will
enable those users to instantiate their own instances of
the given cluster instance. This implies that you will only
be paying for the created snapshot while users deriving a
cluster from yours will incur costs for running the actual
cluster. After the sharing process is complete, services
on your cluster will automatically resume.

◉ Public  ○ Shared

Share-an-instance

Display a menu

Tool installation and configuration, image creation, etc, all **completely automated and extensible**

Cloud instances include all tools available
in main Galaxy *and more*

Same automation approach can be used for
configuring tool dependencies for a local Galaxy

VM image with just tools available, currently at
http://s3.amazonaws.com/usegalaxy/UseGalaxy.ova

**Why we love clouds and cloud-like things:**

Reasonably cost effective and efficient
(elasticity + autoscaling definitely save money)

Analysis costs are more directly quantifiable

Infrastructure as an abstraction + standard APIs for
provisioning reduces risk of vendor lock-in

Virtualization makes *so* many things easier

# Publishing and sharing

# *Everything* can be shared



**Sharing and Publishing History 'Variant Analysis for Sample E18'**

**Making History Accessible via Link and Publishing It**

This history **accessible via link and published**.

Anyone can view and import this history by visiting the following URL:

http://main.q2.bx.psu.edu/u/jgoecks/h/variant-analysis-for-sample-e18 ✎

This history is publicly listed and searchable in Galaxy's Published Histories section.

You can:

Unpublish History

Removes history from Galaxy's Published Histories section so that it is not publicly listed or searchable.

Disable Access to History via Link and Unpublish

Disables history's link so that it is not accessible and removes history from Galaxy's Published Histories section so that it is not publicly listed or searchable.

**Sharing History with Specific Users**

You have not shared this history with any users.

Share with a user

Back to Histories List

Pervasive search allows others to find published items of interest

Galaxy Page for a recent study on mitochondrial heteroplasmy

Actual histories and datasets directly accessible from the text

Histories can be imported and the exact parameters inspected

Workflows and other entities can also be embedded

And imported for inspection, verification, and reuse

# The power of Galaxy publishing

- Galaxy's publishing features facilitate access and reproducibility without any extra leg work

- One click grants access to the *actual analysis* you performed to generate your original results

  - Not just data access: the full pipeline

  - Annotate each step

  - Anyone can import your work and immediately reproduce or build on it

HOME | ABOUT | ARCHIVE | SUBMIT | SUBSCRIBE | ADVERTISE | AUTHOR INFO | CONTACT | HELP

Institution: PENN STATE UNIV  Sign In via User Name/Password

Search for Keyword: [ ] Go
Advanced Search

# Windshield splatter analysis with the Galaxy metagenomic pipeline

Sergei Kosakovsky Pond[1,2,6,9], Samir Wadhawan[3,6,7], Francesca Chiaromonte[4], Guruprasad Ananda[1,3], Wen-Yu Chung[1,3,8], James Taylor[1,5,9], Anton Nekrutenko[1,3,9] and The Galaxy Team[1]

[+] Author Affiliations

## Abstract

How many species inhabit our immediate surroundings? A straightforward collection technique suitable for answering this question is known to anyone who has ever driven a car at highway speeds. The windshield of a moving vehicle is subjected to numerous insect strikes and can be used as a collection device for representative sampling. Unfortunately the analysis of biological material collected in that manner, as with most metagenomic studies, proves to be rather demanding due to the large number of required tools and considerable computational infrastructure. In this study, we use organic matter collected by a

## Footnotes

[Supplemental material is available online at http://www.genome.org. All data and tools described in this manuscript can be downloaded or used directly at http://galaxyproject.org. Exact analyses and workflows used in this paper are available at http://usegalaxy.org/u/aun1/p/windshield-splatter.]

Article published online before print. Article and publication date are at http://www.genome.org/cgi/doi/10.1101/gr.094508.109.

**OPEN ACCESS ARTICLE**

### This Article

- » Abstract *Free*
- Full Text (PDF) *Free*
- Supplemental Material

- All Versions of this Article:
  - gr.094508.109v1
  - 19/11/2144 *most recent*

[–] Article Category

Resource

[+] Services
[+] Citing Articles
[+] Google Scholar
[+] PubMed
[+] Social Bookmarking

[+] Recent Updates

Follow us on twitter

[+] Most Read Articles

View all ...

### Current Issue

October 2010, 20 (10)

[+] From the Cover

Alert me to new issues of *Genome Research*

- **Advance Online Articles**
- **Submit a Manuscript**
- **GR in the News**
- **Editorial Board**
- **E-mail Alerts & RSS Feeds**
- **Recommend to Your Library**
- **Job Opportunities**

# Visualization

Integration with many existing browsers (extensible)

**Visualization and analytics:**
**Galaxy Track Browser**

Entirely web standards based to support
sharing, communicating, and collaborating
around visualizations

Dynamic and responsive

Open source and extremely extensible

http://main.g2.bx.psu.edu/u/jeremy/v/gcc2011-1-viewing-and-navigating

# Galaxy

Analyze Data    Workflow    **Shared Data**    Visualization    Help    User

Published Visualizations | jeremy | GCC2011-1: Viewing and    chr19    625,719 – 682,581

| 630,000 | 640,000 | 650,000 | 660,000 | 670,000 | 680,000 |

UCSC Main on Human: knownGene (chr19)                                                    Auto (Squish)

UCSC Main on Human: all_est (chr19)                                                      Dense

UCSC Main on Human: phyloP46wayPrimates (chr19)                                          Histogram

1

-1

h1-hESC Tophat Mapped Reads                                                             Auto (Squish)

| 630,000 | 640,000 | 650,000 | 660,000 | 670,000 | 680,000 |

Display a menu

With increasingly complex tools, more experimentation with parameters is necessary, visual feedback aids exploration

Galaxy already provides a very sound model for abstracting interfaces to analysis tools

Existing tool framework can be leveraged for **visual analytics**

Dynamic filtering on element properties (here, FPKM for putative transcripts)

Modifying Cufflinks parameters and locally reassembling

Arbitrary visualization types supported
(but not implemented)

Access to tools and visual analytics specific features
(e.g. local computation using global models)
can be used by new visualization types

# Scaling Galaxy: two distinct problems

- So much data, not enough infrastructure.

  - Solution, encourage local Galaxy instances, cloud Galaxy, support increasingly decentralized model, *improve access to exiting resources*

- So many tools and workflows, not enough manpower

  - Focus on building infrastructure to allow community to integrate and share tools, workflows, and best practices

# Galaxy toolshed vision

- Allow users to share "suites" containing tools, datatypes, workflows, sample data, and automated installation scripts for tool dependencies

- Version controlled

- Community annotation, rating, comments, review

- Dependency resolution

- Integration with Galaxy instances to automate tool installation and updates

# Galaxy Tool Shed

Galaxy Tool Shed

**Repositories**

- Browse by category
- Browse all repositories
- Login to create a repository

## Categories

search repository name, description

| Name | Description | Repositories |
|------|-------------|--------------|
| Assembly | Tools for working with assemblies | 10 |
| Computational chemistry | Tools for use in computational chemistry | 2 |
| Convert Formats | Tools for converting data formats | 7 |
| Data Source | Tools for retrieving data from external data sources | 2 |
| Fasta Manipulation | Tools for manipulating fasta data | 12 |
| Graphics | Tools producing images | 4 |
| Next Gen Mappers | Tools for the analysis and handling of Next Gen sequencing data | 12 |
| Ontology Manipulation | Tools for manipulating ontologies | 2 |
| SAM | Tools for manipulating alignments in the SAM format | 3 |
| Sequence Analysis | Tools for performing Protein and DNA/RNA analysis | 27 |
| SNP Analysis | Tools for single nucleotide polymorphism data such as WGA | 2 |
| Statistics | Tools for generating statistics | 4 |
| Text Manipulation | Tools for manipulating data | 9 |
| Visualization | Tools for visualizing data | 4 |

# Galaxy Tool Shed

**Repositories**     **Help**     **User**

**Repositories**

- Browse by category
- Browse all repositories
- Login to create a repository

## Repositories

search repository name, description

Advanced Search

| ☐ | Name ↓ | Synopsis | Revision | Category | Owner | Averag |
|---|--------|----------|----------|----------|-------|--------|
| ☐ | agile_wrapper | Quickly match reads to a reference genome or sequence file | 0:d6a426afaa46 | • Next Gen Mappers<br>• Sequence Analysis | simonl | ★★ |
| ☐ | assemblystats | Summarise an assembly (e.g. N50 metrics) | 0:6544228ea290 | • Next Gen Mappers<br>• Sequence Analysis | konradpaszkiewicz | ★★ |
| ☐ | blast2go | Maps BLAST results to GO annotation terms | 1:0f159cf346c8 | • Ontology Manipulation<br>• Sequence Analysis | peteric | ★★ |
| ☐ | clustalomega | multiple sequence alignment program for proteins | 0:ff1768533a07 | • Fasta Manipulation<br>• Sequence Analysis | clustalomega | ★★ |
| ☐ | contamrm | For fast contaminant filtering from nextgen reads. | 0:6e61b7ddb5f9 | • Sequence Analysis | edward-kirton | ★★ |
| ☐ | cpg_island | TODO | -1:000000000000 | • Sequence Analysis | tiemoon | ★★ |

# Galaxy Tool Shed

Repositories    Help    User

Repository Actions ▾

## clustalomega

**Clone this repository:**
hg clone http://toolshed.g2.bx.psu.edu/repos/clustalomega/clustalomega

**Name:**
clustalomega

**Synopsis:**
multiple sequence alignment program for proteins

**Detailed description:**
Clustal Omega is a general purpose multiple sequence alignment program for proteins. It produces high quality alignme

**Version:**
0:ff1768533a07

**Owner:**
clustalomega

**Times downloaded:**
7

## Categories

Fasta Manipulation

Sequence Analysis

## Repository metadata

**Tools:**

| name | description | version | requirements |
|------|-------------|---------|--------------|
| Clustal Omega | multiple sequence alignment program for proteins | version: 0.2 | none |

# Galaxy Tool Shed

Repositories    Help    User

Repository Actions ▾

## Clustal Omega

**Name for output files:**

co_alignment

**Output guide tree:**

☐ Yes

**Output distance matrix:**

☐ Yes

Clustal-Omega is a general purpose multiple sequence alignment (MSA) program for proteins. It produces high quality MSAs and is capable of handling data-sets of hundreds of thousands of sequences in reasonable time.

In default mode, users give a file of sequences to be aligned and these are clustered to produce a guide tree and this is used to guide a "progressive alignment" of the sequences. There are also facilities for aligning existing alignments to each other, aligning a sequence to an alignment and for using a hidden Markov model (HMM) to help guide an alignment of new sequences that are homologous to the sequences used to make the HMM. This latter procedure is referred to as "external profile alignment" or EPA.

Clustal-Omega uses HMMs for the alignment engine, based on the HHalign package from Johannes Soeding [1]. Guide trees are optionally made using mBed [2] which can cluster very large numbers of sequences in O(N*log(N)) time. Multiple alignment then proceeds by aligning larger and larger alignments using HHalign, following the clustering given by the guide tree.

In its current form Clustal-Omega can only align protein sequences but not DNA/RNA sequences. It is envisioned that DNA/RNA will become available in a future version.

A full version of these instructions is available at http://www.clustal.org/

This is a beta version of Clustal Omega. Bugs should be reported to clustalw@ucd.ie

A standalone version of Clustal Omega for Linux/Windows/Mac is available from http://www.clustal.org/

[1] Johannes Soding (2005) Protein homology detection by HMM-HMM
      comparison. Bioinformatics 21 (7): 951-960.
[2] Blackshields G, Sievers F, Shi W, Wilm A, Higgins DG. Sequence
      embedding for fast construction of guide trees for multiple sequence alignment. Algorithms Mol Biol. 2010
      May 14;5:21.

# Some future challenges

- Capturing and automatically deploying tool dependencies, automatic tool acquisition in Galaxy instances

- Better interfaces for highly parallel analysis (e.g. running the same workflow across 192 individuals)

- Various workflow engine improvements, partial data streaming, combined experimental/computational workflows

**Try it now:**

**http://usegalaxy.org**

**Develop and deploy:**

**http://getgalaxy.org**

**Join us, contact me at:**

**james@jamestaylor.org**

Opportunities for collaboration, positions for postdocs, researchers, software engineers