



Galaxy Community Conference 2014: Visualization Workshop

0.75

Jeremy Goecks

Assistant Prof. of Comp. Biology

THE GEORGE WASHINGTON UNIVERSITY

WASHINGTON, DC

Sam Guerler

Research Software Engineer



Topics

Visualization history and introduction Biological Visualizations Numerical Visualizations Adding your own visualizations

Why Visualize?

Why Visualize?

Quick check: did it work?

Exploration and hypothesis generation

Sharing/publishing

Anscombe's Quartet



http://en.wikipedia.org/wiki/Anscombe's_quartet

Timeline of Visualization in Galaxy



Timeline of Visualization in Galaxy



- 1. visualization in Galaxy is nascent
- 2. you will be working with awesome new features
- 3. there may be bugs help us fix them!

Workshop Goals

Participants: learn about how to visualize your data in Galaxy

- biological visualizations
- numerical visualizations
- what Galaxy is doing underneath the covers

Instructors: feedback from you about what you like, don't like, and where to go next

Galaxy Visualizations

Visualizations are first-class objects in Galaxy, just like tools

A visualization can be added to Galaxy via a configuration file that specifies:

- datasets that can be used
- Iocation of visualization code (client-side or on server)

Galaxy handles visualization integration and data management, so users can focus on analyzing data (and developers can focus on creating visualizations)

Visualizations are 1st class Galaxy objects

Can be saved and versioned for reproducibility

Have a human-readable URL for sharing a fully interactive visualization: http://usegalaxy.org/u/jgoecks/v/tumor-mutations

Can embed interactive visualizations in online supplementary materials via Galaxy Pages

Visualization

Visual Analysis

Visualization Architecture

Client-server architecture



Lots of moving pieces

- prepare/process data on server
- send to client
- render on client

Topics

Visualization history and introduction Biological Visualizations Numerical Visualizations Adding your own visualizations **Analysis goal**: what similarities and differences can be found in cancer cell lines using exome and transcriptome sequencing?

Sequencing and Analysis

Profiled 3 pancreatic cancer cell lines using a 26gene panel targeting known oncogenic driver mutations

MiaPaCa2, HPAC, and PANC-1

Data available:

- exome: mapped reads, removed dups, called variants
- transcriptome: mapped reads, assembled transcripts, computed expression

Display Applications

11: MiaPaCa2: Varscan v ariants
~150,000 lines format: vcf , database: hg19
Log: tool progress Log: tool progress Picked up _JAVA_OPTIONS: - Djava.io.tmpdir=/tmp Got the following sample list: MiaPaCa2-exome Only variants will be reported Min coverage: 8 Min reads2: 2 Min var freq: 0.01 Min avg qual: 15 P-value thresh: 0.99
Read
display at UCSC <u>main</u> display with IGV <u>web current local</u> display at RViewer <u>main</u>

Display Applications

<u>11: MiaPaCa2: Varscan v</u> <u>ariants</u>
~150,000 lines
format: vcf, database: hg19
Log: tool progress
Log: tool progress
Picked up _JAVA_OPTIONS: -
Djava.io.tmpdir=/tmp
Got the following sample list:
MiaPaCa2-exome
Only variants will be reported
Min coverage: 8
Min reads2: 2
Min var freq: 0.01
Min avg qual: 15
P-value thresh: 0.99
Read
B 6 2 III S •
display at UCSC main
aispia, web current local
display at RViewer main

Display Applications

Advantages

- use tool that is familiar to you
- easy to view your data alongside public datasets

Disadvantages

- cannot save/share/version visualization
- many more visualizations than display applications in Galaxy
- no data processing or visual analysis, only visualization

Trackster—Galaxy's Genome Browser



Trackster—Galaxy's Genome Browser

Genome browsers remain a (the most?) powerful genome visualization

foundational tool

Trackster is for the high-throughput sequencing era

- very large datasets, numerous simultaneous tracks
- maximum flexibility for customization (e.g. rainbow tracks)
- 2-3 indices per dataset for fast visualization

SAM/BAM, BED, GFF/GTF, VCF, Wiggle, BigWig, BigBed, BedGraph

- 1. Create visualization
- 2. Add gene annotation (RefSeq)
- 3. Save visualization
- 4. Exit
- 5. Reopen visualization

1. Create visualization

- 1. Create visualization
- 2. Add gene annotation (RefSeq)

- 1. Create visualization
- 2. Add gene annotation (RefSeq)
- 3. Save visualization

- 1. Create visualization
- 2. Add gene annotation (RefSeq)
- 3. Save visualization

4. Exit

- 1. Create visualization
- 2. Add gene annotation (RefSeq)
- 3. Save visualization
- 4. Exit
- 5. Reopen visualization

Behind the Scenes

Galaxy is indexing datasets for

- viewing large genomic regions (coverage plots)
- viewing small genomic regions (getting individual data points)
- feature names and locations

Indexes is the primary way that big datasets are visualized quickly

Modes and Searching

Tracks can be displayed differently

- coverage
- individual features

Let's try different modes

 this is fast because data is sent from Galaxy server and rendered in your Web browser

Let's try searching for a gene: ERBB2

Let's Call Variants

VarScan

- Sample names: MiaPaCa2, PANC1, HPAC
- + Run

Rename output: "Cell line variants"

Let's Assemble Transcripts

Cufflinks

- input dataset is #7
- + run

Rename output: "MiaPaCa2 Assembled Transcripts"

Let's add data to Trackster

Add exome data for all cell lines... ...but where is our data?

Circster

Interactive Circos plot

Whole genome view with structural variation



Let's view our data in Circster

Double-click or use trackpad to zoom in

change track min/max

what do we see?

Let's add data to Circster and adjust options

1. Add transcriptome coverage data

Let's add data to Circster and adjust options

1. Add transcriptome coverage data

2. Change arc dataset height

Let's add data to Circster and adjust options

1. Add transcriptome coverage data

2. Change arc dataset height

3. Change max for tracks
Let's add data to Circster and adjust options

1. Add transcriptome coverage data

2. Change arc dataset height

3. Change max for tracks

4. Save visualization

Back to Trackster: Rainbow Track for Coverage

- 1. Navigate to ERBB2 gene
- 2. Create group
- 3. Add transcriptome coverage tracks to group
- 4. Create composite track
- 5. Adjust max
- 6. what do we see?

Add More Data!

Add RNA-seq mapped reads, variants, and assembled transcripts

Look at ERBB2

* bookmark!

Look at STK11

• bookmark!

Look at KRAS —> LYRM5

• bookmark!

Visual Analysis

KRAS and Variants

Sweepster

Topics

Visualization history and introduction Biological Visualizations Numerical Visualizations Adding your own visualizations

What is Galaxy Charts?



Import data files

10.0	- Galaxy Analyze Data W	Shared Data - Visualizatio	on v Help	+ User+	Using 2.1
		Data Libraries	2		
	Data Library "Charts"	Data Libraries Beta			
	🗋 Name	Published Histories	Data type	Date uploaded	File size
	amino_acid_features.txt -	Published Workflows	tabular	Mon Jun 30 04:13:31 2014 (UTC)	974 bytes
	http://www.compsysbio.org/bacteriome/dataset/functional_in	Published Visualizations Published Pages	tabular	Mon Jun 30 16:33:33 2014 (UTC)	81.6 KB
	For selected datasets: Import to current histo				

1 TIP: You can download individual library datasets by selecting "Download this dataset" from the context menu (triangle) next to each dataset's name.

1 TIP: Several compression options are available for downloading multiple library datasets simultaneously:

- · gzip: Recommended for fast network connections
- · bzip2: Recommended for slower network connections (smaller size but takes longer to compress)
- · zip: Not recommended but is provided as an option for those who cannot open the above formats

Click on **Shared Data** and select **Data Libraries**. Navigate to the **Chart** library and import it into your history (*data reference: http://dna.cs.byu.edu/treesaap and bacteriome.org*).

Make a new chart (1 of 4)

1	51: http://www.compsy 💿 🖋 🗙
	sbio.org/bacteriome/dataset/functio
	nal_interactions.txt
	3,989 lines
	format: tabular, database: ?
	uploaded tabular file
3	Charts
	l Scatterplot
	Trackster
	LIGIL DUDO 0.000 100
	B4200 B4202 0.933934
	B0779 B4058 0.933934
	B0032 B0033 0.933183

Wait for the upload to complete. Select your **Dataset** and click on the **Visualization Icon** then select **Charts**.

Give your chart a name

III Unclustered Heatmap	2
Start Configuration O Add Data	
Provide a chart title:	
Chart title	
How many data points would you like to analyze?	
Few (<500) Some (<10k) Many (>10k)	
• Bar diagrams	
	() () () () () () () () () ()
Pequilar (NVD3) Stacked (NVD3) Herizontal S	tacked
(NVD3) ho	rizontal
	NVD3)
• Others	

Name your chart Unclustered Heatmap.

Select a chart type



Double click on the **Heatmap** icon.

Select data columns

•
•
•

At first click on **Row labels** and select **Column 2**. Then, click on **Draw**.

Unclustered Heatmap



Make a new chart (2 of 4)

1	51: http://www.compsy 💿 💉 🗙
	sbio.org/bacteriome/dataset/functio
	nal_interactions.txt
	3,989 lines
	format: tabular, database: ?
	uploaded tabular file
3	Charts
	l Scatterplot
	l Trackster
	LIGIL DUDGE CLUBSCIDE
	B4200 B4202 0.933934
	B0779 B4058 0.933934
	B0032 B0033 0.933183

Select your **Dataset** and click on the **Visualization Icon** then select **Charts**.

Give your chart a name

🔟 Unclustered Heatma	р			
Start Configurati	on o Add Data			
Provide a chart title:				
Chart title				
How many data poin	ts would you like to	o analyze?		
Few (<500) Some	e (<10k) Many (>	10k)		
• Bar diagrams				
Regular (NVD3)	Stacked (NVD3)	Horizontal (NVD3)	Stacked horizontal (NVD3)	
• Others				

Name your chart **Clustered Heatmap**.

Select a new chart type

• Area charts



⊕Regular (NVD3)





@Stream (NVD3)

Pie chart (NVD3)

• Data processing (requires 'charts' tool from Toolshed)



Histogram (NVD3)



⊕Discrete Histogram (jqPlot)



⊕Box plot (jqPlot)



Double click on the **Clustered Heatmap** icon.

Select data columns

•
•
•

At first click on **Row labels** and select **Column 2**. Then, click on **Draw**.

Clustered Heatmap



Use the mouse wheel or your touch pad to zoom into the highlighted area.

Enlarged view



Tooltips popup if you move the mouse pointer over a box. Here the interaction between **B4143** and **B3295** is highlighted. Click on **Editor** again to further customize this chart.

Chart configuration

Stort Con	figuration 1	· Data label	• Add Date		
Start Con	nguration 1	<u>: Data label</u> 🗢	O Add Data	1	
Provide a cha	rt title:				
New Chart					
New Chart					
New Chart			170-8- 5		
New Chart Iow many da	ta points woul	d you like to ana	lyze?		
New Chart Iow many da Few (<500)	ta points woul Some (<10k)	d you like to ana Many (>10k)	lyze?		
New Chart Iow many da Few (<500) Bar diagrams	ta points woul Some (<10k)	d you like to ana Many (>10k)	lyze?		
New Chart Iow many da Few (<500) Bar diagrams	ta points woul Some (<10k)	d you like to ana Many (>10k)	lyze?		
New Chart Iow many da Few (<500) Bar diagrams	ta points woul Some (<10k)	d you like to ana Many (>10k)	lyze?		

Go to the **Configuration** tab.

Chart settings

X axis:		
Axis label	X-axis	
	Provide a label for the axis.	
Axis value type	Auto	•
had the type	Select the value type of the axis.	
Y axis:		
Axis label	Y-axis	
	Provide a label for the axis.	
Axis value type	Auto	-
Axis value type	Select the value type of the axis.	
Others:		
Show legend	Yes No Would you like to add a legend?	
Color scheme	Jet	•
color scheme	Select a color scheme for your heatmap	
Liri template	http://someurl.com?id=LABEL	
ontemplate	Enter a url to link the labels with external sources. Use LABEL as placeholder.	

Heatmap specific options are **highlighted**. Feel free to set **axis labels** or other options.

Define a URL template

Paste a **database URL** into the template URL field and add the **__LABEL__** tag. You may use **http://www.ncbi.nlm.nih.gov** or any other database. Click on **Draw** to redraw the chart.

Data points linked to web sources



Double click on a **box** and the browser will open two new tabs using the previously defined **URL template**.

Profilee			
S NCBI Resources	⊙ How To ⊙	<u>Sign in</u>	to N
GEO Profiles	GEO Profiles + b4143 Save search Advanced	Search) H
Show additional ters	<u>Display Settings:</u> ⊘ Summary, 20 per page, Sorted by Subgroup effect <u>Send to:</u> ⊘	Filters: Manage Filters	
Gene symbol	Results: 1 to 20 of 47 << First < Prev Page 1 of 3 Next> Last> groEL - Stress factor RooS regulation in exponential-phase	Profile data Download profile data	
Gene keyword Select	1. <u>bacteria</u> Annotation: groEL, molecular chaperone GroEL (multiple annotations exist) Organism: Escherichia coli K-12	Profile pathways	
Organism Select	Reporter: GPL199, mopA_b4143_at (ID_REF), GDS3123, 1037522 (Gene ID), 913705 (Gene ID), 948665 (Gene ID), 959980 (Gene ID), b4143 (ORF) DataSet type: Expression profiling by array, transformed count, 6 samples ID: 49311248	Find pathways	
Gene ontology Select	groL - Indole-3-acetic acid effect on Escherichia coli Anostation and Confit Changmain GmEL large subunit of	Database: Select +	
Differential expression Up/down genes	GroESL (multiple annotations exist) Organism: Escherichia coli, Escherichia coli K-12 Reporter: GPL189, 1240 (ID_REF), GDS2181, b4143 (ORF)		
DataSet keyword	DataSet type: Expression profiling by array, count, 12 samples ID: 27346540 GEO DataSets Gene Profile neighbors Chromosome neighbors	b4143[All Fields]	
Select			

Cluster selection and analysis



Select one element from each **highlighted row**. What are the corresponding **protein functions**?

Identified protein categories



Please return to the **Editor**.

Make a new chart (3 of 4)

1	51: http://www.compsy 💿 💉 🗙
	sbio.org/bacteriome/dataset/functio nal interactions.txt
	3,989 lines format: tabular , database: <u>?</u>
	uploaded tabular file
3	Charts
	l Scatterplot
	Trackster
	B4200 B4202 0.933934
	B0779 B4058 0.933934
	B0032 B0033 0.933183

Select your **Dataset** and click on the **Visualization Icon** then select **Charts**.

Give your chart a name

III Unclustered Heatmap	p			
Start Configuration	on o Add Data			
Provide a chart title:				
Chart title				
How many data point	ts would you like to	analyze?		
Few (<500) Some	(<10k) Many (>	10k)		
• Bar diagrams				
Regular (NVD3)	Stacked (NVD3)	Horizontal	Stacked	
		(1103)	(NVD3)	
• Others		(1122)	(NVD3)	

Name your chart Score Histogram.

Analyze the score distribution



Double click on the **Histogram** icon and click on **Draw**.

Give your chart a name

II Unclustered Heatmap		🖺 Draw
Please select data colu	Imns before drawing the chart.	
Start Configuration	<u>1: Data label</u> • Add Data	
Provide a label:		
Data label		
Select columns:		
Observations	Column: 3 [float]	-

Click on **Draw**.

Export as **PNG**



Click on **Screenshot** and select **Save as PNG**. Finally, return to the **Editor** again.

Make a new chart (4 of 4)

1	51: http://www.compsy 💿 💉 🗙
	sbio.org/bacteriome/dataset/functio
	nal_interactions.txt
	3,989 lines
	format: tabular, database: ?
	uploaded tabular file
3	Charts
	Scatterplot
	l Trackster
	LIGIE 0000 0.000100
	B4200 B4202 0.933934
	B0779 B4058 0.933934
	B0032 B0033 0.933183

Select your **Dataset** and click on the **Visualization Icon** then select **Charts**.

Give your chart a name

III Unclustered Heatmap			B
Start Configuration • Add Data			
Provide a chart title:			
Chart title			
How many data points would you like	to analyze?		
Few (<500) Some (<10k) Many (>10k)		
• Bar diagrams			
		in m	
Regular (NVD3) Stacked (NVD3)	Horizontal	Stacked	
	(NVD3)	(NVD3)	
• Others			
1 million	· .		
X / X X			

Name your chart **Discrete Histogram**.

Analyze the protein distribution



Double click on the **Discrete Histogram** icon.

Add more data

) Draw
Start Configuration 1: Data label • • Add Data	
Data label	
Select columns:	
Observations Column: 1 [str]	•

Click on Add Data.

Select a second data group

Start Configuration 1: Data label 2: Data label O Add Data Provide a label:	
Provide a label:	
Data label	
Select columns:	
Observations Column: 2 [str]	-

At first click on **Observations** and select **Column 2**. Then, click on **Draw**.

Which proteins have most interactions?



Done with Part I.
Scratchbook

Activate the Scratchbook

- Galaxy	Analyze Data Workflow Shared Data - Visualization - Help - User -		Using 2.1 MB
Tools		History	C 🕈
search tools		Unnamed history	
Get Data		2.1 MB	Q 🗹 📎 🗩
Lift-Over	Welcome to Galaxy on the Cloud	25: Charts on data 1	
Text Manipulation	managed by CloudMan	25. Charts on data 1	• / *
Filter and Sort		24: amino acid featur	<u>es.</u> 👁 🖋 🗙
Join, Subtract and Group		txt	
Convert Formats		1: http://www.compsy	<u>/sb</u> 👁 🖋 🗙
Extract Features		io.org/bacteriome/da	taset/functiona
Fetch Sequences		l interactions.txt	
Fetch Alignments		3,989 lines	
Get Genomic Scores		format: tabular, databa	ase: <u>?</u>
Operate on Genomic Intervals		uploaded tabular file	
Statistics			
Graph/Display Data			• •
Regional Variation		1 2 3	
Multivariate Analysis		B1882 B1888 1.000000	
Evolution		B0728 B0729 0.966967	
Motif Tools		B1812 B3360 0.956456	
Multiple Alignments		R0779 R4058 0 933934	
FASTA manipulation		D0027 D0022 0 022102	
<			>

Activate the **Scratchbook** by clicking on the above icon.

Activate the Scratchbook

🗧 Galaxy	Analyze Data Workflow Shared Data - Visualization - Help - User -	Using 2	2.1 MB
Tools	New Track Browser	History	C 🕈
search tools	Saved Visualizations	Unnamed history	
Get Data		2.1 MB Q 🗹	۲
Lift-Over	Welcome to Galaxy on the Cloud	25: Charts on data 1	1 X
Text Manipulation	managed by CloudMan		
Filter and Sort		24: amino acid features.	I ×
Join, Subtract and Group		txt	
Convert Formats		1: http://www.compsysb	€ ×
Extract Features		io.org/bacteriome/dataset/fu	nctiona
Fetch Sequences		L interactions.txt	
Fetch Alignments		3,989 lines	
Get Genomic Scores		format: tabular, database: ?	
Operate on Genomic Intervals		uploaded tabular file	
Statistics			
Graph/Display Data			• •
Regional Variation		1 2 3	
Multivariate Analysis		B1882 B1888 1.000000	
Evolution		B0728 B0729 0.966967	
Motif Tools		B1812 B3360 0.956456	
Multiple Alignments		B4200 B4202 0.933934	
FASTA manipulation		B0779 B4058 0.933934	
<			>

Click on Saved Visualizations.

Activate the Scratchbook

💳 Galaxy		Analyze I	Data Workflow Share	ed Data + Visua	lization - Help -	User -		Using 2.1 MB
(Ŧ	S	aved Visu	alizations		Ē	Ð	History	20
Saved Visualization	s			Create	new visualization	oud ^{budMan}	Unnamed history 2.1 MB 25: Charts on data 1	Q 🕑 🗞 🗩
<u>Title</u>	Type <u>Dbkey</u>	Tags	Sharing	<u>Created</u>	Last Updated †		24: amino acid featur txt	r <u>es.</u> 🕑 🖋 🗙
Unclustered Heatmap 🚽	Charts	<u>0 Tags</u>	Accessible, Published	~4 hours ago	~3 hours ago		1: http://www.comps	vsb 💿 🖋 🗙
Clustered Heatmap -	Charts	<u>0 Tags</u>	Accessible, Published	~21 hours ago	~5 hours ago		io.org/bacteriome/da	taset/functiona
For 0 selected items: Delete]						3,989 lines format: tabular , datab uploaded tabular file	ase: <u>7</u>
Visualizations share	ed with yo	u by	others			2	 C II 2 3 B1882 B1888 1.000000 	•
Evolution							B0728 B0729 0.966967 B1812 B3360 0.956456	
Motif 1001s Multiple Alignments FASTA manipulation							B4200 B4202 0.933934 B0779 B4058 0.933934 B0022 B0022 0.032192	
<								>

Select a Visualization and repeat the process by selecting **Saved Visualizations** again.

Scratchbook for multiple charts



Resize all visualizations so they fit into the screen.

More Examples

Create a pie chart



Select the imported datasets, create a new chart and select **Pie chart**. Then, click on **Add data**.

Add first data group

					- Curreer	Diam
Start	<u>Configuration</u>	1: Helix frequency •	2: Beta frequency •	• Add Data		5
Provide a l	label:					
Helix frequ	uency					
Select colu	imns:					
Labels		Column: 1 [str]			•
Values		Column: 7 [flo	at]			•
		1				

Configure the Helix frequency column.

Add second data group

III New C	hart				🖺 Draw
<u>Start</u>	<u>Configuration</u>	1: Helix frequency •	2: Beta frequency •	• Add Data	
Provide	a label:				
Beta fre	quency				
Select co	olumns:				
Labels		Column: 1 [str	1		•
Values		Column: 8 [flo	at]		•
		L			

Configure the **Beta frequency** column.

Configure the pie chart

🖻 New Chart		🛾 Cancel 🖺 Draw
Start Configuration 1:	Helix frequency O 2: Beta frequency O O Add Data	
Pie chart settings:		
Donut ratio	50% Determine how large the donut hole will be.	•
Show legend	Yes No Would you like to add a legend?	
Label settings:		
Donut label	Label column What would you like to show for each slice?	•
Show outside	Yes No Would you like to show labels outside the donut?	

Configure the **Pie chart** as shown above. Then, click on **Draw**.

Configure the pie chart



Glutamic acids seem to fit much better into **helices** than **beta sheets**. In other words, "Aspartic and Glutamic Acids are Important for Alpha-helix Folding", JBSD 2007.

Create a bar diagram



Create data groups for the following features: Hydrophobicity, Membrane frequency, Flexibility, Helix frequency and Beta frequency.

Bar diagram of amino acid features



Use the **tooltips** to identify the amino acids which are likely to be found within membrane proteins.

Topics

Visualization history and introduction Biological Visualizations Numerical Visualizations Adding your own visualizations

Adding your own Visualizations

Go to config/plugins/visualizations/charts

charts/others/YOURVIZNAME

Add three files to this directory:

Logo (logo.png) Configuration (config.js) Wrapper (wrapper.js)

charts/types.js

Rebuild by typing 'npm install' and 'grunt'