Workflows (and more!) with Galaxy

Expression Dynamics of Human XBP1 Revealed using an Integrative Approach

Jeremy Goecks The Galaxy Team & Emory University

Approach

Describe how Galaxy was used to complete an NGS biological analysis, and highlight Galaxy features along the way

The Investigation

Understanding human XBP1 expression dynamics

Unfolded Protein Response



Walter & Ron, Science 2011

IRE1 Pathway



Walter & Ron, Science 2011

Translational Pausing



Ron & Ito, Science 2011

XBP1 Isoforms



very similar isoforms

frameshift between isoforms

Chung et al., PLoS CB 2007

Questions

Can current RNA-seq tools be used to identify and quantify *XBP1*'s highly similar isoforms?

Can XBP1's translational pause point be identified?

Galaxy 101 (the *really* short version)

https://main.g2.bx.psu.edu/galaxy101

00	Galaxy
+ http://main.g2.bx	.psu.edu/ C Q* Google O
🗧 Galaxy	Analyze Data Workflow Shared Data Lab Visualization Admin Help User
Tools Options	Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ
Get Data Upload File from your computer UCSC Main table browser UCSC Archaea table browser BX main browser BX main browser Get Microbial Data BioMart Central server GrameneMart Central server Flymine server modENCODE fly server Mormbase server Wormbase server EuPathDB server EncodeDB at NHGRI EpiGRAPH server Send Data ENCODE Tools Lift-Over Text Manipulation Flater and Sort Join, Subtract and Group Extract Features Fetch Sequences Fetch Alignments	Help Table Browser Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see Using the Table Browser for a description of the controls in this form, the User's Guide for general information and sample queries, and the OpenHelix Table Browser tutorial for a narrated presentation of the software features and usage. For more complex queries, you may want to use Galaxy or our public MySQL server. To examine the biological function of your set through annotation enrichments, send the data to GREAT. Refer to the Credits page for the list of contributors and usage restrictions associated with these data. clade: Mammal • genome: Human • assembly: Feb.2009(GRCh37/hg19) • group: Genes and Gene Prediction Tracks • track: UCSC Genes • add custom tracks table: knownGene • position chr22 • lookup define regions identifiers (names/accessions): gaste list • upload list filter: (reate intersection: (reate) output format: BED - browser extensible data • @ Send output to • Galaxy • GREAT output file: (leave blank to keep output in browser) file type returned: • plain text • gzip compressed get output • summary/statistics To reset all user cart settings (including custom tracks), click here. there.
Get Genomic Scores Operate on Genomic Intervals	Using the Table Browser

00			1	Galaxy				
🔺 🕨 🕂 🚱 javascript:window.r	esizeTo(1024,7	768)					¢	Q+ Google
🗧 Galaxy	Analyze Data	Workflow	Shared Data	Lab	Visualization	Admin	Help	User
ToolsOptionsGet DataUpload File from your computerUCSC Main table browserUCSC Archaea table browserBX main browserGet Microbial DataBioMart Central serverGrameneMart Central serverFlymine servermodENCODE fly serverRatmine serverMormbase serverEuPathDB serverEncodeDB at NHGRIEpiGRAPH serverSend DataENCODE ToolsLift-OverText ManipulationConvert FormatsFASTA manipulationFilter and SortJoin, Subtract and GroupExtract FeaturesFetch AlignmentsGet Genomic ScoresOperate on Genomic Intervals		The following 1: UCSC Ma You can chec by refreshing will change fi 'error' if prob	I job has been s in k the status of a the History par rom 'running' to lems were encou	uccessfu queued ji ne. When 'finished untered.	lly added to the obs and view th a the job has ber I' if completed s	queue: e resulting en run the uccessfully	data status or	History Options - Unnamed history 1: UCSC Main on Human: @ 0 % knownGene (chr22:1-51304566)



00			(Galaxy				
+ http://main.g2.b	.psu.edu/						¢	Q+ Google
- Galaxy	Analyze Data	Workflow	Shared Data	Lab	Visualization	Admin	Help	User
A main provide the server A mode NCODE fly server Mormbase server Mormbase server EuPathDB server EuPathDB server EncodeDB at NHGRI EpiGRAPH server Send Data EncodeDB at NHGRI EpiGRAPH server Send Data ENCODE Tools Lift-Over Text Manipulation Convert Formats FASTA manipulation Filter and Sort Join, Subtract and Group Extract Features	Analyze Data	Workflow The following 2: UCSC Ma You can chec by refreshing will change fr 'error' if prob	Shared Data Job has been su in k the status of o the History par rom 'running' to lems were encou	Lab uccessful queued jo e. When 'finished' untered.	Visualization ly added to the obs and view the the job has bee if completed si	Admin queue: e resulting o en run the s uccessfully o	C Help	Qr Google User History Options - Image: Cooperative state
Fetch Sequences Fetch Alignments Get Genomic Scores Operate on Genomic Intervals) 4 4							



00	Galaxy
+ http://main.g2.bx.	psu.edu/ C Q Google O
- Galaxy	Analyze Data Workflow Shared Data Lab Visualization Admin Help User
Tools Options -	Group Options -
Convert Formats FASTA manipulation	Select data: 3: Join on data 2 and data 1 Galaxy 101
Filter and Sort Join, Subtract and Group	Query missing? See TIP below. Group by column: 3: Join on data 2 and data 1 @ 0 % 16,190 regions, format: interval,
 Join two Queries side by side on a specified field Compare two Origins to find 	Ignore case while grouping?:
Compare two oueres to find common or distinct rows Subtract Whole Query from	Operations
another query <u>Group</u> data by a column and perform aggregate operation	Add new Operation Add new Opera
on other columns. Column Join	Execute chr22 16266928 16267095 uc002z1h. 1_cc Chr22 16266928 16266928 16267095 uc002z1h. 1_cc TIP: If your data is not TAB delimited, use Text Manipulation->Convert 0hr22 16266928 16267095 uc002z1h. 1_cc
Extract Features Fetch Sequences Fetch Allegements	Syntax
Get Genomic Scores Operate on Genomic Intervals	This tool allows you to group the input dataset by a particular column and perform aggregate functions like Mean, Median, Mode, Sum, Max, Min, Count, Random draw and Concatenate on other columns.
Statistics Graph/Display Data	All invalid, blank and comment lines are skipped when performing the aggregate functions. The number of skipped lines is displayed in the
Regional Variation Multiple regression Multivariate Analysis	resulting history item. If multiple modes are present, all are reported.
Evolution Metagenomic analyses	Example For the following input:
Human Genome Variation EMBOSS	chr22 1000 1003 TTT chr22 2000 2003 ana chr10 2200 2203 TTT chr22 2000 2203 TTT
NGS TOOLBOX BETA	Grouping on column 4 while ignoring case, and performing operation
NGS: Mapping NGS: SAM Tools	Count on column 1 will return: AAA 2



Concert Formats Convert forma	00	Galaxy
Callaxy Analyze Data Workflow Shared Data Lab Visualization Admin Help User Tools Options ▼ Lift-Over Text Manipulation • Add column to an existing query Select first • Concatenate queries tail-to-head Sistent first • Concatenate queries tail-to-head Sistent of data 4 • Concatenate queries tail-to-head Execute • Marge Columns together • What it does • Concatenate single interval as a new query This tool outputs specified number of lines from the beginning of a dataset • Cance Columns from a table • Charge Case of selected columns • Charge Case of selected • Charge Case of selected columns • Select first • Select first • Charge Case of selected • Charge Select Birstow Carefield • Charge Select Birstow Carefield • Select first times from a Query • Select first • Select Birstow Carefield • O X • Select first • Charge Calumn Indig form a file • Charge Calumn Indig form a file • Charge Case of selected • Charge Case of selected • Charge Select Birstow Carefield • Charge Select Birstow Carefield • O X • Select first lines from a	+ http://main.g2.bx.	psu.edu/ C Qr Google 🕥
Tools Options ▼ LIH-Over Select first Select first Select first: Select first: Se	💳 Galaxy	Analyze Data Workflow Shared Data Lab Visualization Admin Help User
FASTA manipulation Filter and Sort Join, Subtract and Group Extract Features Fetch Sequences Fetch Alignments	Tools Options Litt-Over Text Manipulation • Add column to an existing query • Compute an expression on every row • Concatenate queries tail-tohead • Concatenate queries tail-tohead • Condense consecutive characters • Convert delimiters to TAB • Merge Columns together • Create single interval as a new query • Cut columns from a table • Change Case of selected columns • Paste two files side by side • Remove beginning of a file • Select first lines from a Query • Select last lines from a Query • Select last lines from a Query • Trim leading or trailing characters Convert Formats FASTA manipulation Filter and Sort Join, Subtract and Group Extract Features Fetch Alignments Get Genomic Scores • Cores	Select first Select first Select first Since From: Siston a data 4 Creace What it does This tool outputs specified number of lines from the beginning of a dataset Example Selecting 2 lines from this: chr 2 \$6652 \$17003 CTCF 26 \$310 * chr 3 \$6756 \$17003 CTCF 26 \$310 * chr 3 \$6756 \$17003 CTCF 26 \$310 * chr 3 \$6756 \$17003 CTCF 26 \$310 *
	Get Genomic Scores	

Galaxy 101 Highlights

Requires only a Web browser

Integrated data sources, including UCSC

Collection of many tools that can be combined into complex analyses

Datasets grouped together into a history



OPINION

The real cost of sequencing: higher than you think!

Andrea Sboner^{1,2}, Xinmeng Jasmine Mu¹, Dov Greenbaum^{1,2,3,4,5}, Raymond K Auerbach¹ and Mark B Gerstein^{#1,2,6}



Computation in Science?

Scientists unfamiliar with computation

Reproducibility hindered by complexity: systems, scripts, tools, parameters

Collaboration and publishing difficult because current media do not support computational artifacts well

Galaxy Project: Fundamental Questions

When Biology (or any science) becomes dependent on computational methods, how to:

- make tools and methods 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 transparent computational biomedical research

What is Galaxy?

GUI for genomics

+ for complete analyses: analyze, visualize, share, publish

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

Amplification

Many tools available in a single place means that tools can be combined in novel ways

Everything is open source: framework, protocols, and libraries

reused and extended by anyone

Users, developers, community benefit

Questions

Can current RNA-seq tools be used to identify and quantify *XBP1*'s highly similar isoforms?

Can XBP1's translational pause point be identified?



http://www.fml.tuebingen.mpg.de/raetsch/members/research/transcriptomics

Illumina BodyMap 2.0

RNA-seq (Hi-Seq) data from 16 human tissues

Two datasets for each tissue

- 50bp paired-end
- 75bp single-end

~80 million reads per dataset, ~160 million reads per tissue

In a public data library on usegalaxy.org

http://www.ensembl.info/blog/2011/05/24/human-bodymap-2-0-data-from-illumina/

Differential Expression Analysis in Galaxy



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

DE NOVO ASSEMBLY

 <u>Trinity</u> De novo assembly of RNA-Seq data

FILTERING

 <u>Filter Combined Transcripts</u> using tracking file

Trapnell, C., Pachter, L. and Salzberg, S.L. TopHat: discovering splice junctions with RNA-Seq. Bioinformatics 25, 1105-1111 (2009).
 Trapnell et al. Transcript assembly and abundance estimation from RNA-Seq reveals thousands of new transcripts and switching among isoforms. Nature Biotechnology doi:10.1038/nbt.1621

Differential Expression Analysis in

NATURE PROTOCOLS | PROTOCOL

Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks

Cole Trapnell, Adam Roberts, Loyal Goff, Geo Pertea, Daehwan Kim, David R Kelley, Harold Pimentel, Steven L Salzberg, John L Rinn & Lior Pachter

Affiliations | Contributions | Corresponding author

Nature Protocols 7, 562–578 (2012) | doi:10.1038/nprot.2012.016 Published online 01 March 2012

Abstract

Abstract · Accession Codes · Author Information

Recent advances in high-throughput cDNA sequencing (RNA-seq) can reveal new genes and splice variants and quantify expression genome-wide in a single assay. The volume and complexity of data from RNA-seq experiments necessitate scalable, fast and mathematically principled analysis software. TopHat and Cufflinks are free, open-source software tools for gene discovery and comprehensive expression analysis of high-throughput mRNA sequencing (RNA-seq) data. Together, they allow biologists to identify new genes and new splice variants of known ones, as well as compare gene and transcript expression under two or more conditions. This protocol describes in detail how to use TopHat and Cufflinks to perform such analyses. It also covers several

-	print
4	email
J)	download citation
9	order reprints
9	rights and permissions
1	share/bookmark

pression Analysis

-SEQ

tat for Illumina Find splice tions using RNA-seq data

<u>inks</u> transcript assembly FPKM (RPKM) estimates for -Seq data

<u>compare</u> compare mbled transcripts to a ence annotation and track inks transcripts across iple experiments

<u>merge</u> merge together ral Cufflinks assemblies

<u>diff</u> find significant changes anscript expression, ing, and promoter use

IOVO ASSEMBLY

<u>ty</u> De novo assembly of -Seq data

ERING

<u>Combined Transcripts</u> tracking file

ics 25, 1105-1111 (2009). ranscripts and switching

 Trapnell, C., Pachter,
 Trapnell et al. Transc among isoforms. Na accessory tools and utilities that aid in managing data, including CummeRbund, a tool for visualizing RNAseq analysis results. Although the procedure assumes basic informatics skills, these tools assume little to no background with RNA-seq analysis and are meant for novices and experts alike. The protocol begins with raw sequencing reads and produces a transcriptome assembly, lists of differentially expressed and regulated genes and transcripts, and publication-quality visualizations of analysis results. The protocol's execution time depends on the volume of transcriptome sequencing data and available computing resources but takes less than 1 d of computer time for typical experiments and ~1 h of hands-on time.





Working with NGS Tools

Often challenging

- many parameters
- time intensive
- evaluating results difficult

Three ways Galaxy can help:

- experimentation: can rerun tools, workflows
- visualization: Trackster, display applications
- reproducibility: parameters tracked, workflows available

🖶 🧶 🌒 🚍 Galaxy		
← → C # A https://m	ain.g2.bx.psu.edu	☆ ヽ
- Galaxy	Analyze Data Workflow Shared Data - Visualization - Admin Help - User -	Using 383.5 Gb
Tools	Tophat for Illumina (version 1.5.0)	History Ø
search tools Get Data Send Data ENCODE Tools Lift-Over Text Manipulation Convert Formats FASTA manipulation FILTER and Sort Join, Subtract and Group Extract Features Fetch Sequences Fetch Alignments Get Genomic Scores Operate on Genomic Intervals Statistics Graph/Display Data Regional Variation Multiple regression Multivariate Analysis Evolution Motif Tools Multiple Alignments Metagenomic analyses Human Genome Variation Ges Oc and manipulation NGS: OC and manipulation NGS: SAM Tools NGS: Indel Analysis	RNA-Seq FASTQ file: LisRR030882_lprain/fastq @ Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33 Will you select a reference genome from your history or use a built-in index?: Use a built-index Built-ins were indexed using default options Select a reference genome: Human (Homo sapines): hg19 Full @ If your genome of interest is not listed, contact the Calaxy team Is this library mate-paired?: Paired-end @ RAA-Seq FASTQ file: @ ERR030882_lprain/fastq @ Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33 Man Inner Distance between Mate Pairs: 10 TopHat settings to use: Commonity used @ Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33 Man Inner Distance between Mate Pairs: 10 TopHat settings to use: Commonity used @ Execute Tophat Overview TopHat settings to use: Commonly used settings. If you want full control use Full parameter list TopHat set splice junction mapper for RNA-Seq reads. It aligns RNA-Seq reads to mammalian-sized genomes using the ultra high-throughput short read aligner Bowtie, and then analyzes the mapping results to identify splice junctions between exons. Please cite: Trappell, C., Pachter, L. and Salzberg, S.L. TopHat: discovering splice junctions with RNA-Seq. Bioinformatics 25, 1105-1111 (2009). Know what you are doing Image: Pair is no such thing up to as an automated gearshift in splice junction identification. It is all like stick-shift driving in San Francisco. In	transcript expression 13: Cufflinks on data 8: ● Ø % gene expression 12: BodyMap-Brain 75bp ● Ø % SE mapped reads 11: Tophat for Illumina ● Ø % on data 4 and data 3: splice junctions 10: Tophat for Illumina ● Ø % on data 4 and data 3: deletions 9: Tophat for Illumina ● Ø % data 4 and data 3: deletions 9: Tophat for Illumina on ● Ø % data 4 and data 3: deletions 9: Tophat for Illumina on ● Ø % data 4 and data 3: insertions 8: BodyMap-Brain 50bp ● Ø % PE mapped reads 3.9 Gb format: bam, database: hg19 Info: TopHat v1.4.0 tophat -p 8 -r 110 -a 8 -m 0 -i 20 -1 500000 -g 40 -G // galaxy/main_pool/pool2/files/003 /634/dataset_3634785.dat Ilbrary-type fr-unstrandedmax-insertion-length 3coverage-searchmin-coverage-intron 20max-c Image: Image
NGS: Peak Calling	Input formats	

Galaxy

- tool integration framework
- heavy focus on usability
- sharing, publication framework

Genome Browser

- physical depiction of data
- visually identify correlations
- find interesting regions, features

Trackster

Trackster

Web-based visualization of your NGS data

- requires only a Web browser
- dynamic and customizable
- supports BAM, BED, GFF/GTF, WIG, VCF

Platform for visual analysis

tools integrated with visualization

Can share & publish fully-functional visualizations



Image: Service Control Computer-Lanazonaus com/tracket/h/22.29191945-29192235 Image: Service Control Computer-Lanazonaus com/tracket/h/22.99191945-29192235 Image: Service Control Computer-Lanazonaus com/tracket/h/22.99191945-29192235 Image: Service Control Contro Control Control Control Control Control Control Control	🔲 🌑 🖨 🧰 Tabuar			
Calaxy Addyce Data Text Og199 (333.00) (333.00) (333.00) (333.00) (11) (11) (11) (11) (11)	← → C ᡤ ③ ec2-50-17-146-90.compute-1	.amazonaws.com/tracks#chr22:29191945-29192	325	公 人
Tex (hg 19) (hg 2, hg 3, hg 4, hg 2, hg 3, hg 4, hg 2, hg 4,	- Galaxy	Analyze Data Workflow Shared D	ata - Visualization - Help - User -	Using 29.4 Gb
Dy 182,000 Dy 182,000 <td>Test (hg19)</td> <td>chr22</td> <td>29,191,945 - 29,192,325 👂 🔎</td> <td>02/20</td>	Test (hg19)	chr22	29,191,945 - 29,192,325 👂 🔎	02/20
Trephet for Elumina on data 3: accepted, MB III III IIII IIIIIIIIIIIIIIIIIIIIII	29,192,000 IUCSC Main on Human: refGene (chr22:1-51304566)	29,192,100	29,192,200	29,192,300
No reads spanning small intron	IITophat for Illumina on data 1: accepted_hits 🔍 🖃 🌒 🗇			
No reads spanning small intron				
29.192.000 29.192.100 29.192.200 29.192.300		No reads	s spanning small	intron
	29,192,000	29,192,100	29.192.200	29,192,300 <

🖶 🗘 🌒 🚎 Calaxy		
🗲 🔿 C 🕂 🚨 https://main.g	2.bx.psu.edu	☆ ヽ
- Galaxy	Analyze Data Workflow Shared Data Visualization Admin Help User -	Using 383.5 Gb
Galaxy Tools Image: Construct of the search tools Search tools Search tools Get Data Send Data ENCODE Tools Lift-Over Text Manipulation Convert Formats FASTA manipulation Filter and Sort Join, Subtract and Group Extract Features Fetch Sequences Fetch Alignments Get Genomic Scores Operate on Genomic Intervals Statistics Graph/Display Data Regional Variation Multiple regression Multivariate Analysis Evolution Motif Tools Multiple Alignments	Analyze Data Workflow Shared Data - Visualization - Admin Help+ User+ Tophat for Illumina (version 1.5.0) RNA-Seq FASTQ file: I: ERR030882_1_brain.fastq : Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33 Will you select a reference genome from your history or use a built-in index?: Built-ins were indexed using default options Select a reference genome: Human (Homo sapiens): hg19 Full * If your genome of interest is not listed, contact the Galaxy team Is this library mate-paired?: Paired-end : RNA-Seq FASTQ file: : ERR030882_2_brain.fastq : Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33 Mucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33 Mucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33 Mucleotide space: Must have Sanger-scaled quality values with ASCII offset 33 Mean Inner Distance between Mate Pairs: 110 North mapping need use Commonly used settings. If you want full control use Full parameter list Full parameter list : TopHat settings to use Full parameter list : TopHat will treat the reads as strand specific. Every read a	Using 383.5 Gb History Image: state stat
Metagenomic analyses Human Genome Variation Genome Diversity EMBOSS NGS TOOLBOX BETA NGS: QC and manipulation NGS: Mapping NGS: SAM Tools NGS: Indel Analysis	Std. Dev for Distance between Mate Pairs: 20 The standard deviation for the distribution on inner distances between mate pairs. Anchor length (at least 3): 8 Report junctions spanned by reads with at least this many bases on each side of the junction. Maximum number of mismatches that can appear in the anchor region of spliced alignment: 0 The minimum later least b	display at UCSC main display at Ensembl <u>Current</u> display with IGV <u>web current local</u> display in IGB <u>Local Web</u> Binary ben alignments file 7: Tophat for Illumina on O X data 2, data 4, and data 1: splice junctions
K Peak Calling		* II. >

³⁸ experimentation



³⁹ experimentation

🗧 单 🚔 🧮 Galaxy 🛛 🗶		
← → C 🕂 🔇 ec2-50-17-146-90.comp	ute-1.amazonaws.com/tracks#chr22:29191945-29192325	☆ 、
Galaxy	Analyze Data Workflow Shared Data - Visualization - Help - User -	Using 29.4 Gb
Test (hg19)	chr22 29,191,945 - 29,192,325 👂 🔊	0 4 4 6 9
29,192,000 UCSC Main on Human: refGene (chr22:1-51304566)	29,192,100 29,192,200	29,192,300
IAdrenal 75bp SE mapped reads 🗹 🖃 🗟 👘 🕩 🛪		
29,192,000	29,192,100	29,192,300
	40 SUCCESS	







Waypoint

Found good parameter values via experimentation

 can map and assemble RNA-seq reads from a single dataset

But there are 16 tissues, each with 3 RNA-seq datasets



Galaxy Workflow System

Galaxy		a Workflow	Shared Data					Us	ing 383.5 (
ools	Options -	Workflow Canva	as BodyMap	Mapping and Asse	mbly				Options
filter								Workflow	v Paramet
Filter and Sort									
 <u>Filter</u> data on any cousing simple expression 	olumn ssions								
Filter on ambiguitie polymorphism data	is in sets								
GFF		Input Dataset	82						
Filter GFF data by a using simple expre-	ttribute ssions	output	5						
Filter GFF data by fe	eature		/	Tophat for Illumina					
expressions			L.	RNA-Seq FASTQ fil	e		Filter GFF	data by attri	bute 🕱
Filter GTF data by a	ttribute		A	Gene Model Annot	ations		Filter		
values list				insertions (bed)	00		out_file1		DO
Fetch Alignments			1	deletions (bed)	ad				
<u>Filter MAF</u> by specif	fied			junctions (bed)	B		Filme		
attributes				accepted hits (harr			Filter	×	
Operate on Genomic	Intervals			accepted_into (bail	m	\mathcal{P}	O Filter		-
 Intersect the interva datasets 	als of two				1		out_file1	0	9
Subtract the interva	ls of two								
datasets				Cufflin	ks		×		
<u>Cluster</u> the interval dataset	s of a	Input dataset	×	SAM or reads	BAM file	of aligned F	RNA-Seq		
Graph/Display Data		output	¢.	Referer	nce Annot	ation			
VCF to MAF Custom	1 Track			Global	model (fo	r use in Tra	ckster)	8	
for display at UCSC				genes_	expressio	n (tabular)	00		
Regional Variation				transcr	ipts_expre	ession (tabu	ular) 🗆 🖸 🖸	-	
 Filter nucleotides b quality scores 	ased on			assemt	oled_isofo	rms (gtf)			
Fetch Indels from 3	-way			total_m	nap_mass	(txt)			

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

Facilitate reuse and provide precise reproducibility of a complex analysis

Galaxy Workflows

			C	C.d. raodie
red Data	a Visualizati	ion Help	User	
e first m	iegabyte is sho	wn below.		Histor History Lists
07 00 07 00 07 00 07 00 07 00 07	60 96 106 173 144 117 70 79 136 101	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	0. 555555888444588	EI8 Current History Gov 10: V Create New Samp 26,74 Clone datab Share or Publish Info Extract Workflow
001001010000100	- 137 - 184 - 196 - 197 - 129 - 219 - 240 - 240 - 134 - 153 - 120	0 60 0 60 0 60 0 60 0 60 129 60 60 0 60 0 0 0	50950034448899558	s Logic Extract worknow b Ldisp Dataset Security Gene Show Deleted Datasets Curre Show Hidden Datasets Show structure chr10 Delete
0004404040400	07 104 110 103 55 132 139 138 138 134 144	0 60 0 60 0 60 0 60 135 60 132 60 138 60 0 60	534343543543582	chr10 14465082 14465083 T K 173 chr10 14465083 14465084 G K 144 chr10 14465084 14465085 T 7 117 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9
0F0FCAF	80 117 138 154 128	0 60 0 60 138 60 211 60 0 60	30 37 37 64 47 35	Main o 1 1
	95 165 80 255 237 234 242 242 242 242 242 242 242 242 242	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	184312533225433	7: Map with Bowtie for 0 8 Illumina on data 6 and data 5 9,073,928 lines, format: sam, database: mm9 Info: Sequence file aligned. Co Max
00000440	175 255 180 195 152 139 101	175 60 0 60 0 60 0 60 0 60 0 60 0 60	490 47 53 53 53 53 53 53 53 53 53 53 53 53 53	I. ORACE 2. FLAG B HWT_EAS269:3:1:1449:913 99 cl B HWT_EAS269:3:1:1449:913 147 cl B HWT_EAS269:3:1:709:832 99 cl
0	83	0 60	32	.8 Y HWI-EAS269:3:1:709:832 147 cl





Introduction Autor Marce Oldard Oldard Maddel Jd57 Mchr22 29189779-29196670 Image Oldard Marce Oldard <th></th> <th>Galaxy</th> <th>×</th> <th></th> <th></th> <th></th> <th></th> <th></th>		Galaxy	×					
Calaxy Analyza Data Wardfall Mark Mark Calax Wardfall Mark Calax War	e = 0	https://main.g2.bz	x.psu.edu/tracks/browser?id=01	ded04aabb43d57#chr22:2918977	9-29196670			☆ 4
BodyMay Transcript Assemblies for Multiple Transcript A	- Gal	axy		Analyze Data Workflow Shared	Data - Visualization - Admin	Help + User +		Using 383.5 Gb
Displace	BodyMap 1	Transcript Assemblies for	Multiple Tissues (hg19)	chr22	29,189,779 - 29,196,670	Q Q	_	0210
	29,19	90,000	29,191,000	29,192,000	29,193,000	29,194,000	29,195,000	29,196,000
	BodyMap As	sembled Transcripts						
	CLIFF CLIFF	M. 005088 36886 1 36886 2 MM 001274533			CUFF 36887.1 CUFF.36888.	1		
	Brain							
Protate 000000000000000000000000000000000000		CUFF.55307.2 5M_005080 5M_001073539 CUFF.55307.1			CUFF. 55388. 1			CVFF, 55389.1
	Prostate							
	9.1	NM_005080			CUFF.62444.1			corr.s
		CUFF.62443.1 MM_001079539						
	Ovary							
29.190.000 29.191.000 29.195.000 29.195.000 29.195.000 29.195.000	5443.1	NM_005030 NM_001070510						CUFF.SS383.
29,190,000 29,191,000 29,192,000 29,193,000 29,194,000 29,195,000 29,195,000 29,196,000								
	29.19	90.000	29,191,000	29,192,000	29,193,000	29,194,000	29,195,000	29.196.000

兽 单 🌒 🚍 Galaxy	×								
🗲 🔿 C 📅 🔒 https://main.g2	.bx.psu.edu								☆ २
- Galaxy		Analyze Data	orkflow Shared Da	ata - Visualization	- Admin Help - User -			Using	383.5 Gb
Tools ©	TCONS_00051704 TCONS_00051705 TCONS_00051706	j NM_00107953 j NM_00107953 = NM_00107953	 XLOC_027679 XLOC_027679 XLOC_027679 XLOC_027679 	XBP1 TSS36244 XBP1 TSS36245 XBP1 TSS36246	chr22:29190544-29196560 chr22:29190544-29196560 chr22:29190544-29196560	2051 - 1672 - 1794 -	1.54003 1 2.41066 1 1.2497 (History 1 60: Cuffmerge-Cuffdiff	• • 0 z 1
<u>Get Data</u> <u>Send Data</u> <u>ENCODE Tools</u> Lift-Over	TCONS_00051707	= NM_005080	XLOC_027679	XBP1 TSS36246	chr22:29190544-29196560	1820 -	14.4725	59: Cuffcompare-Cuffdiff	• <i>0</i> z
Text Manipulation Convert Formats FASTA manipulation								58: Cuffdiff on data 12, data 8, and others: transcri FPKM tracking	90 () X <u>PI</u>
Filter and Sort Join, Subtract and Group Extract Features								57: Cuffdiff on data 12, data 8, and others: transcri differential expression test	es () 22 Int ting
Fetch Sequences Fetch Alignments Get Genomic Scores								56: Cuffdiff on data 12, data 8, and others: gene FP tracking	®⊳0% <u>KM</u>
Operate on Genomic Intervals Statistics Graph/Display Data								55: Cuffdiff on data 12. data 8, and others: gene differential expression test	® Ø ¤ ting
Regional Variation Multiple regression Multivariate Analysis								54: Cuffdiff on data 12, data 8, and others: TSS gro FPKM tracking	90 / 33 11/25
Evolution Motif Tools Multiple Alignments								53: Cuffdiff on data 12, data 8, and others: TSS gro differential expression test	ab (/ 33 ups ting
Metagenomic analyses Human Genome Variation Genome Diversity								52: Cuffdiff on data 12, data 8, and others: CDS FPI tracking	®⊳0% <u>(M</u>
EMBOSS NGS TOOLBOX BETA NGS OC and manipulation								51: Cuffdiff on data 12, data 8, and others: CDS FPI differential expression test	®⊳∂% KM ting
NGS: SAM Tools NGS: Indel Analysis								50: Cuffdiff on data 12, data 8, and others: CDS overloading diffential expr testing	ession
A Calina								in .	

Quantification

Read to transcript mapping difficult because read can map to multiple transcripts

FPKM = fragments per kilobase of exon per millions of reads

normalize by exon length and sample size

1 FPKM ~ 1 transcript per cell in mouse

Results

1000 CC	Adipose	Adrenal	Brain	Breast	Colon	Heart	Kidney	Liver	Lung
XBP1u	24.5 (CI: 22.0-27.0)	84.7 (77.6-91.7)	14.5 (11.9-1	7.0) 69.4 (63.0-7	5.8) 46.1 (41.0-51.	1) 39.9 (34.8-45.0)	121.1 (111.2-131.0)	287.1 (253.8-320.3)	89.2 (82.2-96.3)
XBP1s	11.7 (10.1-13.3)	15.5 (13.3-17.7)	1.2 (0.8-1.7)	12.7 (10.4-1	4.9) 23.9 (20.8-27.)	0) 14.1 (11.2-16.9)	10.3 (8.7-11.8)	25.3 (20.7-29.8)	40.9 (37.2-44.7)
Non-coding Isoform (uc003aec.2)		6.9 (6.0-7.9)	1.5 (1.0-2.1)						
			к	L	м	N	0	Р	Q
		Lymph I	Vode	Ovary	Prostate	Skeletal Muscle	Testes	Thyroid	White Blood Cells
		325.7 (30	10.4-350.9)	74.6 (66.9-82.4)	29.9 (27.4-32.4)	60.4 (55.0-65.8)	22.6 (20.5-24.8)	49.0 (44.9-53.1)	97.0 (88.9-105.1)
		58.5 (53.	3-63.7)	21.2 (18.0-24.4)	Not enough data	20.1 (17.6-22.5)	12.3 (10.8-13.8)	7.3 (5.9-8.8)	4.7 (3.4-6.0)

Can assemble XBP1 isoforms and quantify expression

• Quick check: *XBP1u* > *XBP1s* in all tissues

Pervasive transcription of both *XBP1u* and *XBP1s* in 15 (16?) tissues

very high in some tissues, lower in others

UPR is a continuous process rather than a discreet event

Questions

Can current RNA-seq tools be used to identify and quantify *XBP1*'s highly similar isoforms?

Can XBP1's translational pause point be identified?

Ribosome Profiling

Sequence RNA undergoing translation + ~35bp

In this work, data is colorspace



Ingolia et al., Cell 2011



	🖨 🌒 🌒 📃 Galaxy	* 101									
	← → C 🔒 https://n	main.g2.bx.psu.edu/tr	acks/browser?id=83	86f9e6817	ef93#X8P1u:0-9	17					☆ 4
Name: XP2 Translational Parce Point (premy shp1-0) Day: Day: <thday:< th=""> Day: Day: <t< th=""><th>- Galaxy</th><th></th><th></th><th>Analyze Da</th><th></th><th>Shared Data - Vi</th><th>sualization - Adr</th><th></th><th></th><th></th><th>Using 383.5 Gb</th></t<></thday:<>	- Galaxy			Analyze Da		Shared Data - Vi	sualization - Adr				Using 383.5 Gb
100 100 <td>Human XBP1 Translational</td> <td>Pause Point (jeremy:xl</td> <td>bp1-t)</td> <td></td> <td>XBP1u</td> <td></td> <td>0 - 917</td> <td>e e</td> <td></td> <td></td> <td>02/20</td>	Human XBP1 Translational	Pause Point (jeremy:xl	bp1-t)		XBP1u		0 - 917	e e			02/20
	0 100	0	200	300		400	500	600	200	800	900
	EXBP1 Transcriptome Annotation										
	start_coden)		eson II		111001111000		exon 26bp_intro	• 1000000		leu246 proline trp256]	
	IER Fostprints mapped reads, see	ed Jen=28, mismatches=2 P	PRE & x							stop_coson	
				-			-				
0 100											
0 100 200 100											
Inter Foregrers 5' Read Coverage 0 30 1											
100 1											
100 200 1											
Burk Foregreis 37 1											
33 0 1 </td <td>ER Postprints 5' Read Coverage</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td>000</td>	ER Postprints 5' Read Coverage										000
0 1	33									10	
UCytosol Footgarints mapped reads, seed_lem=2 Image: International control of the set of the se	0		1.1			1.	1		1. 1	l. i	
	Cytosol Footprints mapped reads	s, seed_len=28, mismatches	1-2								
9 <u>1109 1209 1209 1400 1509 1600 1700 1800</u>											
0 180 200 100 1400 1500 15	14										
0 <u>100</u> 200 <u>300</u> 409 <u>500</u> 600 <u>700</u> 800 <u>70</u>											
9 <u> 100 200 300 409 509 600 700 800 </u> 600 											
0 100 200 200 400 500 600 1700 600 600											
0. <u>1200 200 400 500 600 700 800</u>											
0 100 200 300 400 500 600 700 800 600											
0 100 200 300 400 500 600 700 800 600											
0 100 200 300 1400 500 600 700 800											
0 100 200 300 400 500 600 700 800											
0 100 200 300 400 500 600 700 800 K											
0 1200 1200 1400 1509 1600 1700 1820 K											
0 100 200 300 1400 500 600 700 800											
0 100 200 300 400 500 600 700 800 C											
	0 100	0	200	300		400	500	1600	700	800	(

Everything can be Shared

兽 🧶 💭 📆 Galaxy	×	an a				and the second sec
🔶 🤿 😋 🔒 https://main	.g2.bx.psu.edu/root#	e				☆ *
- Galaxy Ana	lyze Data Workflow	Shared Data -	Visualization +	Admin Help	+ User +	Using 383.5 Gb
Share or Publish Hist	tory 'BodyMap	Brain'			History	0
Making History Accessible This history is currently accessib	via Link and Publis	shing It			De BodyMap Brain	2 🖻 7.4 Gb
https://main.g2.bx.psu.edu/	u/jeremy/h/bodymap-	brain			XBP1 Isoform Trackin	ш Ф (/ Ж 9
You can:	~				59: Cuffcompare-Cuff XBP1 Isoform Tracking	' <u>diff</u> ● 0 %
Disable Access to History Link Disables history's link so that it is	s not accessible.				58: Cuffdiff on data 1 data 8, and others: tra FPKM tracking	2,
Publishes the history to Galaxy's	<u>Published Histories</u> sec	tion, where it is p	ublicly listed and se	archable.	57: Cuffdiff on data 1. data 8, and others: tra differential expression	2, @ () & nscript n testing
Sharing History with Specie You have not shared this history	fic Users with any users.				56: Cuffdiff on data 1 data 8, and others: ge tracking	2, @ 0 % ne FPKM
Back to Histories List					55: Cuffdiff on data 1 data 8, and others: ge differential expression	2. @ () 💥 ne n testing
>					Ш	>

👙 🗢 🍧 🧮 Galaxy Accessible History 🛛 🛪		
C A https://main.g2.bx.psu.edu/u/jeremy/h/bodymap-brain		4 4
Galaxy Analyze Data	Workflow Shared Data - Visualization - Admin Help - User -	Using 402.0 Gb
Accessible History BodyMap Brain	Import history	About this History
Galaxy History ' BodyMap Brain' Dataset <u>1: ERR030882 1 brain.fastq</u> <u>2: ERR030882 2 brain.fastq</u>	Annotation Illumina Hi-Seq, forward reads, 50bp. Illumina Hi-Seq, reverse reads, 50bp.	Author jeremy Related Histories All published histories Published histories by jeremy Rating
3: ERR030890 brain.fastq Image: Constraint of the second seco	Illumina Hi-Seq, single-end reads, 75 bp.	Community (0 ratings, 0.0 average) Yours Tags Community: none Yours:
8: BodyMap-Brain 50bp PE mapped reads 3.9 Gb format: bam, database: hg19 Info: TopHat v1.4.0 tophat -p 8 -r 110 -a 8 -m 0 -i 20 -i 500000 -g 40 -C /galaxy/main_pool/pool2/files/003/634/dataset_3634785.datlibrary-type fr- unstrandedmax-insertion-length 3max-deletion-length 3coverage-search min coverage-intron 20max-c ↓ ↓ display at LCSC main display at Ensembl <u>Current</u> display in IG8 Local Web Binary ban alignments file	Using RefSeq gene model.	
9: Tophat for Illumina on data 4 and data 3: insertions		
11: Tophat for Illumina on data 4 and data 3: splice junctions		
12: BodyMap-Brain 75bp SE mapped reads 2.3 Gb	Using RefSeq gene model.	

Exact Reproduction Possible

Galaxy Accessible Hit	itory ×							
C A https://main	.g2.bx.psu.edu/u/jeremy/h/bo	dymap-brain						(1) (1) (1) (1) (1) (1) (1) (1) (1) (1)
- Galaxy	Ar	alyze Data Workflow	Shared Data -	Visualization +	Admin	Help + User +		Using 402.0 Gb
Accessible History BodyMap Brain	6						Import history About this History	γ
Galaxy History ' BodyMap Bra	兽 🌑 🌢 / 🚍 https://main.g2.bx.psu.ed	×						
Dataset	C A https://main.g2	bx.psu.edu/datasets/12d116	/fbedcbSe9/show_pa	irams?username=j	eremy&slug=b	odymap-brain&use_p	anels=True	2
1: ERR030882 1 brain.fastq	Name: Created:	BodyMap-Bra Feb 13, 2012	ain 50bp PE mapped rea	ds				
2: ERR030882_2_brain.fastq	Filesize: Dbkev:	3.9 Gb ha19						
3: ERR030890 brain.fastg	Format: Tool Version:	bam						k #
4: UCSC hg19 refseg genes	Tool Standard Output: Tool Standard Error:	stdout stderr						**
5: Tophat for Illumina on data 2, c	Full Path:	/galaxy/mair	n_pool/pool1/files/003	/775/dataset_37750	065.dat		Value	
6: Tophat for Illumina on data 2, c	RNA-Seq FASTQ file						1: ERR030882_1_brain.fastq	
7: Tophat for Illumina on data 2, c	Select a reference genome Conditional (singlePair						/galaxy/data/hg19/bowtie_index/hg19 1	
8: BodyMap-Brain 50bp PE mappe	RNA-Seq FASTO P Mean Inner Cance between Mate Pairs						2: ERR030882_2_brain.fastq 110	
format: bam, database: hg19	Conditional (prarams) ary Type Std. Dev for Distance between Mate Pairs						I FR Unstranded 20	
tophat -p 8 -r 110 -a 8 -r -i 20	Anchor length (at least 3) Maximum number of mismatches that ca	n appear in the anchor region o	f spliced alignment				8	
unstrandedme insertion-length	The minimum intron length The maximum intron length						20 500000	
	Conditional (indel_search) Max insertion length.						3	
display at UCSC main display at Ensembl <u>Current</u>	Max deletion length. Maximum number of alignments to be al	owed	16 - 162				3 40	
display with IGV web current local display in IGB Local Web	Maximum intron length that may be foun Maximum intron length that may be foun Number of mismatches allowed in the ini	d during split-segment (default) d during split-segment (default tial read manning) search				20 500000	
Binary ban alignments file	Number of mismatches allowed in each s Minimum length of read segments	egment alignment for reads ma	pped independently				2 25	
9: Tophat for Illumina on data 4 a	Conditional (own_junctions) Conditional (gene_model_ann)						0 1 4: USEC hallo rafing constr	
10: Tophat for Illumina on data 4	Conditional (raw_juncs) Only look for supplied junctions						0 No	
11: Tophat for Illumina on data 4	Conditional (closure_search) Conditional (coverage_search)						1 0	4
12: BodyMap-Brain 75bp SE mapp	Minimum intron length that may be foun ed reads	d during coverage search		Using RefSeq ge	ne model.		20	5



Galaxy Published Histories									
► + http:	//main.g2.bx.psu.edu/history/list	_published					C Q+ Go	ogle	
🚾 Galaxy		Analyze Data	Workflow	Shared Data	Visualization	Help	User		
Published Hist	ories								
search	Advanced Search								
Name	Annotation			Owner	Commun Rating †	<u>ty</u>	Community Tags	Last Updated	
<u>Galaxy vs MEGAN</u>	Comparison of Galaxy vs. MEGA	N pipeline.		aunl	***	k sk	metagenomics megan galaxy	Mar 19, 2010	
metagenomic analysis				aunl	***	**	(metagenomics) (galaxy)	Mar 19, 2010	
<u>SM_1186088</u>	Datasets correspond to our pap Peleg et al. entitled : Altered his associated with age-dependent Experiment layout: This history form of BED files of uniquely mi- chip-seq for histone modification mouse hippocampus of 3 month (old) mice after fear conditionin please refer to supplementary mi- respective work by peleg et al.	er published in Sc tone acetylation is memory impairme contains 4 datase upped reads produ- ins H4K12ac and is (young) and 16 g. For detailed infi- naterials and meth	ience by sent. ts in the iced after H3K9ac in months ormation ods of the	fischerlab	***	k ik		Apr 19, 2010	
Variant Analysis for Sample E18	Perform a pileup analysis with o variants in sample E18.	lefault parameters	to identify	jgoecks	***	h ik	snp pileup bowtie demo sample	2 minutes ago	
get longest exon				henri	***	kik	chr22 longest marc exon human workshop	Sep 02, 2010	
FASTA to Tabular Test				u	shakaka	kik –		Aug 26, 2010	
EKLF				yzc109	***	inir -		Aug 24, 2010	

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 + annotations
- anyone can import your work and immediately reproduce or build on it

Three Ways to Use Galaxy

1. Public Website (http://usegalaxy.org)

2. Download and run locally

3. Run on the cloud

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

Public Website, 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

Galaxy on the Cloud

For extended or particular resource needs

- customization necessary
- oscillating data volume

For when informatics expertise or infrastructure is limited









Enis Afgan



Dave Clements



Dannon Baker



Jeremy Goecks



Dan Blankenberg



Jennifer Jackson



Anton Nekrutenko



Nate Coraor



Greg von Kuster



James Taylor

Supported by the **NHGRI** (HG005542, HG004909, HG005133), **NSF** (DBI-0850103), Penn State University, Emory University, and the Pennsylvania Department of Public Health

Thanks! Questions?

http://galaxyproject.org

💳 Galaxy			
]	Data intensive bi	ology for everyon	e.
<u>Galaxy</u> is a research. W perform, re	n open, web-based platforn /hether on the <u>free public s</u> produce, and share comple	n for data intensive biomed <u>erver</u> or <u>your own instance</u> te analyses.	ical , you can
Use Galaxy	Get Galaxy	Learn Galaxy	Get Involved
	\checkmark	Advanced fastQ manipulation Provide Advanced fastQ	
Use the free public server	Install locally or in the cloud	Screencasts, Galaxy 101,	Mailing lists, Tool Shed, wiki

Galaxy publications: http://galaxyproject.org/wiki/Citing Galaxy is hiring! http://galaxyproject.org/wiki/GalaxyisHiring

jeremy.goecks@emory.edu