Progress and Challenges in Developing a Web-based Platform for Computational Biomedical Research

Jeremy Goecks Depts. of Biology and Math & Computer Science Emory University











Dave Clements



Dannon Baker



Jeremy Goecks



Dan Blankenberg



Jennifer Jackson



Nate Coraor



Greg von Kuster



Kanwei Li



James Taylor



Kelly Vincent



Anton Nekrutenko

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

Overview

Genomics

Galaxy

- + accessible, reproducible, and transparent science
- on the cloud
- visual analytics

Reflections on Galaxy



http://en.wikipedia.org/wiki/Central dogma of molecular biology



http://en.wikipedia.org/wiki/Central dogma of molecular biology

Central Dogma of Molecular Biology : Eukaryotic Model



http://en.wikipedia.org/wiki/Central dogma of molecular biology



Goals of Genomics

Identify and annotate all functional genomic elements

• genes, promoters, enhancers, silencers, epigenetic modifications

Understand genomic regulation

interactions, networks, feedback systems

Apply knowledge of genome to address biomedical challenges

- personalized medicine
- + aging
- environment interactions
- pathogen analysis
- + ..

Trends in Genomics (1)



Trends in Genomics (2)





Will Computers Crash Genomics?

New technologies are making sequencing DNA easier and cheaper than ever, but the ability to analyze and store all that data is lagging

you-go service, accessible from one's own desktop, that provides rented time on a large cluster of machines that work together in parallel as fast as, or faster than, a single powerful computer. "Surviving the data deluge means computing in parallel," says Michael

"Will Computers Crash Genomics?", Pennisi, E., Science, Feb 11, 2011

Challenges in Genomics

Generating data is easy

- high-throughput sequencing (HTS) technologies improving rapidly
- datasets are hundreds of MBs to GBs

Analyzing data is THE bottleneck

- computation is essential due to dataset size
- \$1,000 genome, \$1,000,000 interpretation?

http://www.bio-itworld.com/2010/10/01/interpretation.html

Using Computation in Science?

Scientists often not trained in computation

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

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

Overview

Genomics

Galaxy

- accessible, reproducible, and transparent science
- on the cloud
- visual analytics

Reflections on Galaxy

Galaxy Project: Fundamental Questions

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

- how can those methods best be made accessible to scientists?
- how best to ensure that analyses are reproducible?
- how best to facilitate transparent communication and reuse of analyses?

Vision

Galaxy is an open, Web-based platform for accessible, reproducible, and transparent computational biomedical research

Connecting Users with Tools



Filter and Sort

- <u>Filter</u> data on any column using simple expressions
- <u>Sort</u> data in ascending or descending order
- <u>Select</u> lines that match an expression

GFF FILES

- Extract features from GFF file
- <u>Filter GFF file by attribute</u> using simple expressions
- Filter GFF file by feature count using simple expressions

Extract Features Fetch Sequences Fetch Alignments Get Genomic Scores Operate on Genomic Intervals Statistics Graph/Display Data Regional Variation Multiple regression Multivariate Analysis Evolution Metagenomic analyses EMBOSS

NGS TOOLBOX BETA

NGS: QC and manipulation NGS: Mapping NGS: SAM Tools NGS: Indel Analysis NGS: Peak Calling

RGENETICS

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

Workflows

ecting Users with Tools



Accessibility

Filter and Sort

- Filter data on any column using simple expressions
- Sort data in ascending or descending order
- Select lines that match an Operate on Genomic Intervals
 - Intersect the intervals of two queries
- E Subtract the intervals of two aueries Fi
 - SI Merge the overlapping intervals of a query Fi
 - u <u>Concatenate</u> two queries into one query
 - Base Coverage of all intervals
 - Coverage of a set of intervals on second set of intervals
 - Complement intervals of a query
 - Cluster the intervals of a query
 - Join the intervals of two gueries side-by-side
 - Get flanks returns flanking region/s for every gene
 - Fetch closest feature for every interval
 - Profile Annotations for a set of genomic intervals

ecting Users with Tools



Accessibility



Accessibility

Options v

002

O 0 22

· 0 ×

DO 22

C Q- Google

History

00

Genes

E18

Trimmed

Trimmed

2: E18 PE.2 Reads

1: E18 PE.1 Reads

Sample E18

than ref. base

10: Variants from sample

7: Map with Bowtie for

Illumina on data 6 and data 5

imported: SNP Pileup Analysis for

15: Variants from sample 🛛 👁 🖉 🕱

14: UCSC mm9 RefSeq Genes @ 0 12

13: Variants from sample @ 0 🕱 E18 where consensus base different

9: Generate pileup on data 8 @ 0 22

8: SAM-to-BAM on data 7 @ 0 2 %

6: E18 PE.2 Reads Groomed, @ 0 12

5: E18 PE.1 Reads Groomed, @ 0 22

4: E18 PE.2 Reads Groomed @ 0 12

3: E18 PE.1 Reads Groomed @ 0 12

E18, consensus different, in RefSeq

Filter and Sort

Accessibility Filter data on any column using simple expressions Filter pileup DOIS Sort data in ascending Select dataset: descending order 10: Variants from sample E18 + Select lines that match which contains: Operate on Genom e Pileup with six columns (simple) + Intersect the interview See "Types of pileup datasets" below for examples queries Do not consider read bases with quality lower than: Options -E Subtract the inte 20 queries orted: SNP Pileup Analysis for No variants with quality below this value will be reported Fi ple E18 si Do not report positions with coverage lower than: Merge the overla Variants from sample 🛛 👁 🖉 🕱 consensus different, in RefSeq of a query 3 Fi Pileup lines with coverage lower than this value will be skipped u UCSC mm9 RefSeg Genes 👁 🖉 💥 NGS: SAM Too Only report variants?: Variants from sample 🛛 👁 🖉 🕱 where consensus base different Filter SAM o Yes 🛟 ref. base values See "Examples 1 and 2" below for explanation Variants from sample 002 Convert coordinates to intervals?: Convert SAN No 🛟 enerate pileup on data 8 👁 🖉 🕱 SAM-to-BAN See "Output format" below for explanation AM-to-BAM on data 7 @ 0 🕱 format to BA Print total number of differences?: Map with Bowtie for · 0 × mina on data 6 and data 5 No 🛟 BAM-to-SAM 18 PE.2 Reads Groomed, @ 0 🕱 See "Example 3" below for explanation format to SA nmed Print quality and base string?: 18 PE.1 Reads Groomed, @ 0 🕱 Merge BAM Yes 🛟 med files togethe See "Example 4" below for explanation 18 PE.2 Reads Groomed @ 0 🕱 Generate pil 18 PE.1 Reads Groomed @ 0 12 Execute dataset 18 PE.2 Reads · 0 × Filter pileup on coverage and 1: E18 PE.1 Reads 00 U X SNPs aligner designed to be ultrafast and memory-efficient. It is developed by Ben apnell. Please cite: Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and nent of short DNA sequences to the human genome. Genome Biology 10:R25. Pileup-to-Interval condenses pileup format into ranges of bases

Filter and Sort		Accessibility				
Filter data on any column using						
simple expressions	Filter pileup					
 Sort data in ascending descending order 	Select dataset:	History Options 👻				
descending order	10: Variants from sample E18					
Select lines that match Operate on Genom	which contains:	🕲 🖃 🛛 🖉 🖻				
 Intersect the intersect 	Pileup with six columns (simple)	Variant Analysis for Sample E18				
G queries	Do not consider read bases with quality lower than:	15: Intersect to get Variants @ 0 %				
E Subtract the inte	20	from sample E18, consensus different,				
<u>Fi</u> queries	No variants with quality below this value will be reported	in RefSeq Genes				
si • Merge the overla	Do not report positions with coverage lower than:					
<u>Fi</u> of a query	3 Dilaun lines with severage lower than this value will be skinned	14: UCSC mm9 RetSeq Genes @ 0 🕅				
NGS: SAM Too	Only report variants?	13: Filter to get Variants from \textcircled{O}				
Filter SAM o	Yes 🗘	sample E18 where consensus base				
I values	See "Examples 1 and 2" below for explanation	different than ref. base				
• <u>(</u> • <u>Convert SAN</u>	Convert coordinates to intervals?:	10: Filter pileup to get				
SAM-to-BAM	See "Output format" below for explanation	Variants from sample E18				
format to B/	Print total number of differences?:					
BAM-to-SAM		9: Generate pileup on data 8 👁 🖉 💥				
format to SA	Print quality and base string?	8: SAM-to-BAM on data 7 @ // \$2				
 J = <u>Merge BAM</u> files together 	Yes +	OF SAME OF BAM OF BUILD OF BAM				
nies togethe	See "Example 4" below for explanation	7: Map with Bowtie for 🛛 👁 🖉 💥				
 <u>Generate pil</u> dataset 	Execute	Illumina on data 6 and data 5				
Eiltor piloup	on coverage and	6: E18 PE-2 Reads Groomed. (D) 12				
i SNPs	aligner designed to be ultrafact and memory-afficient. It is develo	Trimmed				
F Pileup-to-Ini	terval condenses					
pileup forma	t into ranges of	5: E18 PE.1 Reads Groomed.				
bases	22	Irimmed				

ilter an	d Sort		This dataset is Show all I Save	large an	d only the firs	t megab	oyte is sh	nown below	<i>N</i> .		ŀ	Accessib	oility
Filter	data on any co		<u>show an</u> (<u>sare</u>										
simple	e expressions	chr10 chr10	6882036 6882037 14243075	A 14243076	A 107 G	0 G	60 96	32 0	.\$, 60	c 35 t			
Sort d	lata in ascendir	chr10 chr10	14243079 14465082	14243080 14465083	č	č	106 173	0 176	60 60	35 . 35 G	Gi	Optio	ns 👻
descer	nding order	chr10 chr10	14465083 14465084	14465084 14465085	Ğ T	К Т	144 117	144	60 60	35 . 38 .		Copilo	
Select	lines that mate	chr10 chr10	14465085 14465257	14465086 14465258	ē c	ē c	70 79	0	60 60	38 . 42 .			2
e Op	erate on Geno	chr10 chr10	14465258 14465263	14465259 14465264	A A	A A	137 136	0	60 60	46 . 61 .		(
	Intersect the in	chr10 chr10	14465366 14465371	14465367 14465372	A G	A G	101 137	0 0	60 60	38 gt 50 .1	nalysis for Samp	ole E18	
G	aueries	chr10 chr10	14465410 14465447	14465411 14465448	G T	G T	184 186	0 0	60 60	69 .9 65 .9	\$		0
E.		chr10 chr10	14465456 14465465	14465457 14465466	G T	G T	193 177	0 0	60 60	70 . 63 .s	sect to get Varia	nts 👁	$\theta \propto$
	Subtract the int	chr10 chr10	14465485 14465569	14465486 14465570	C T	T T	129 219	129 0	60 60	34 t: 84 .	nple E18, conser	isus diff	erent,
<u>Fi</u> (queries	chr10 chr10	14465581 14465586	14465582 14465587	G C	e ç	$240 \\ 248$	0	60 60	84 ,9 82 .9	<u>1 Genes</u>		
si 🛛	Merge the over	chr10 chr10	14465621 14465658	14465622 14465659	c c	ç	134 134	0	60 60	49 . 49 ,			
FI C	of a query	chr10 chr10	14465660 14465691	14465661 14465692	T G	T G	153 128	0	60 60	55 42 .9	mm9 RefSeq G	enes 👁	0 🛛
ū,		chr10 chr10	14465778 14465791	14465779	C G	с G	89 104	0	60 60	34 ,9 33 ,9	\$		
	NGS: SAM To	chr10 chr10	14465881 17445088	174450882	Å	A	103	0	60 60	41 . 34 .	to get Variants	from @	0 🛚
	Filter SAM	chr10 chr10	17445271 17731269	17731270	Ť	Ť	113	0	60 60	42	18 where conse	nsus bas	se
- 6	values	chr10 chr10	19928468	19928469	ç	Ť	132	132	60	35 T	than ref. base		
	C	chr10	19928494	19928495	Ĉ	Ť	138	138	60	37 T	ŕ.		
	 Convert SA 	chr10 chr10	19928538	19928539	Ğ	Ĝ	144	0	60	52 ,9 40 G	pileup to get	۲	0 %
	SAM-to-BA	chr10 chr10	19928741	19928742	T	T	80 117	0	60 60	30 ,	from sample E1	8	
= (format to I	chr10 chr10	28750217 28750397	28750218 28750398	Č	T	138 154	138 211	60 60	37 Ť 64 C	т. *	-	
0		chr10 chr10	28750401 28750423	28750402 28750424	Å C	Å	128	0	60 60	47 ,9 35 Ť	ate nileun on dat	ta 8 on	0 52
- (BAIM-to-SA format to 4	chr10 chr10	28750438 28750446	28750439 28750447	Ă	Â G	95 165	0 165	60 60	36 .9 46 62	and pricup on da	<u></u>	~~~
	ionnat to .	chr10 chr10	28750487 28750512	28750488 28750513	A G	Ā G	80 220	0	60 60	31 , 72 ,	-RAM on data	7 1	$D \otimes 2$
- 1	 Merge BAN 	chr10 chr10	28750548 28750574	28750549 28750575	G T	C T	255 237	255 0	60 60	97 C	brin on uala	<u> </u>	0 00
9	files toget	chr10 chr10	28750577 28750578	28750578 28750579	T T	T T	234 242	0 0	60 60	82 ,9 76 ,9	ith Pourtie for	0	$D \sim$
- (Ceperate n	chr10 chr10	28750593 28750640	28750594 28750641	G T	G C	220 165	0 165	60 60	75 És 46 Ce	an data 6 and d	e e e	0 23
	dataset	chr10 chr10	28750746 28750766	28750747 28750767	G A	A G	202 205	202 205	60 60	58 AJ 59 G			
		chr10	28750769	28750770	<u>T</u>	C	175	175	60	49 c			0 00
-	 Filter pileu 	<u>p</u> on c	overage and							0: E18 F	rE.2 Reads Groom	ied, 👁	0 23
- 1	SNPs			alig	ner designed to be u	ultrafast and	d memory-e	efficient. It is d	evelo	Trimme	<u>a</u>		
- F	Pileup-to-	Interva	l condenses	apn	nt of short DNA seq	uences to th	aphen C, Po he human g	enome. Genom	ne Bic				0.44
9	pileup forn	nat int	o ranges of		_		_	_		<u>5: E18 F</u>	PE.1 Reads Groom	<u>ed,</u> @	$0 \otimes$
	bases									Trimme	d		



A Tool in Galaxy



Defined via abstract interface:

- inputs & outputs
- parameters
- how to generate command line

As simple as possible but allows for rigorous reasoning

Reproducibility in Genomics

18 *Nat. Genetics* experiments in microarray gene expression

<50% of 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)

Reproducibility in Genomics

18 *Nat. Genetics* experiments in microarray gene expression

<50% of 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) 14 re-sequencing experiments in *Nat. Genetics, Nature, Science*

0% reproducible?

Problems

- missing primary data (50%)
- tools unavailable (50%)
- missing parameter setting, tool versions (100%)

"Devil in the details," Nature, vol. 470, 305-306 (2011).

Metadata = Reproducibility



Automatic Metadata

7: Map with Bowtie for 9,073,928 lines, format: sam, database: mm9 Run this job again Run this job again 10 minute aligned.					
1.QNAME	2.FLAG	3.1			
HWI-EAS269:3:1:1449:913	99	chi			
HWI-EAS269:3:1:1449:913	147	chi			
HWI-EAS269:3:1:709:832	99	chi			
HWI-EAS269:3:1:709:832	147	chi			
HWI-EAS269:3:1:1422:1087	99	chi			
HWI-EAS269:3:1:1422:1087	147	chi			
)∢	•			

Map with Bowtie for Illumina

Will you select a reference genome from your history or use a built-in index?: Use a built-in index Built-ins were indexed using default options Select a reference genome: mm9 if your genome of interest is not listed - contact Galaxy team Is this library mate-paired?: Paired-end P

Forward FASTQ file:

5: E18 PE.1 Reads Gr..ed, Trimmed 🛟 Must have Sanger-scaled quality values with ASCII offset 33

Reverse FASTQ file:

6: E18 PE.2 Reads Gr..ed, Trimmed 🛟

Must have Sanger-scaled quality values with ASCII offset 33

Maximum insert size for valid paired-end alignments (-X):

1000

The upstream/downstream mate orientation for valid paired-end alignment against the forward reference strand (--fr/--rf/--ff):

FR (for Illumina) 🛟

Bowtie settings to use:

Commonly used

For most mapping needs use Commonly used settings. If you want full control use Full parameter list

Suppress the header in the output SAM file:

\checkmark

Bowtie produces SAM with several lines of header information by default

Execute

28



User Metadata

History	otions 🔻
	0 🖻
Variant Analysis for Sample E18	· —
Tags: snp x pileup x bowtie x demo x sample:e18 x C)
Annotation / Notes: Perform a variant analysis with defa parameters to identify variants in sa E18 that lie in annotated genes.	ult ample

10: Variants from sample E18 26,742 regions, format: interval, database: mm9 Info: ↓ ②	 X						
Tags:							
pileup × sample:e18 ×							
snps × Z							
Annotation:							
Find variants with coverage >= 30 and quality score >= 20.							
display at UCSC <u>main</u> view in <u>GeneTrack</u> display at Ensembl <u>Current</u>							
1.Chrom 2.Start 3.End 456							
chr10 6882036 6882037 A A 107	2.1						
chr10 14243075 14243076 G G 96	1						
chr10 14243079 14243080 C C 106	5.1						
chr10 14465082 14465083 T K 173):						
chr10 14465083 14465084 G K 144	ł :						
chr10 14465084 14465085 T T 117	2						

Data Provenance

Datasets are immutable in Galaxy

- associations between Dataset objects and their usage: HistoryDatasetAssociation, LibraryDatasetAssociation
- metadata associated with associations, not datasets

Exporting from Galaxy

histories and workflows can be exported in JSON format

Metadata

- Datatype defined in Python code
- Job/tool metadata has official
- JSON format in database and when exported

Galaxy Workflows

\varTheta 🔿 🔿 🛛 Galaxy						
Image: Image	du/ d	Q Google				
Galaxy	Analyze Data Workflow Shared Data Visualization Help Use	r				
Tools Options 👻		Histor				
search tools	This dataset is large and only the first megabyte is shown below. <u>Show all Save</u>	Saved Histories				
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	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	SNP P Histories Shared with Me Current History 10: V Create New Samp 26,74 Clone Clon				
Statistics Graph/Display Data Regional Variation Multiple regression	chri0 14455881 14455882 G G 110 0 60 41 chri0 17445281 14455882 G G 110 0 60 41 chri0 17445271 17445272 A A 103 0 60 34 chri0 17445271 17445272 A A 55 0 60 34 chri0 17731290 T T 113 0 60 42 chri0 19920468 19920469 C T 132 132 60 35 chri0 19920468 19920469 C T 132 132 60 35 chri0 19920468 19920469 A A 119 0 60 44	chr10 14465082 14465083 T K 173 : chr10 14465083 14465084 G K 144 : chr10 14465084 14465085 T T 117 (
Evolution Metagenomic analyses EMBOSS	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	9: Generate pileup on data 8 8: SAM-to-BAM on data @ 0 %				
NGS TOOLBOX BETA	opriv 28759402 A A 128 0 60 47 ohr10 28759423 28759424 C T 113 113 60 35 ohr10 28759438 28759439 A A 95 0 60 36 ohr10 28759438 28759439 A A 95 0 69 36					
NGS: QC and manipulation NGS: Mapping NGS: SAM Tools NGS: Indel Analysis NGS: Peak Calling	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Off 7: Map with Bowtie for ● Ø \$\$ 15 Illumina on data 6 and data 5 25 9,073,928 lines, format: sam, database: mm9 16 Info: Sequence file aligned. 08 10				
SNP/WGA: Data; Filters SNP/WGA: OC: LD: Plots	chr10 28750769 28750770 T C 175 175 60 49 chr10 28750787 28750798 T T 255 0 60 90 chr10 28750797 28750798 C C 180 0 60 64 chr10 28750781 28750798 C C 180 0 60 64 chr10 28750813 28750814 C C 195 0 60 67 chr10 28750813 28750814 A A 152 0 60 67	CO: 1.QNAME 2.FLAG 3.1 .5 HWI-EA3269:3:1:1449:913 99 cha .6 HWI-EA3269:3:1:1449:913 147 cha				
SNP/WGA: Statistical Models	chr10 28750835 28750836 A A 139 0 60 52 chr10 28750850 28750851 G G 101 0 60 38 chr10 28750873 28750874 C C 83 0 60 32	4 HWI-EA3269:3:1:709:832 99 cha				

Galaxy Workflows

00	Tool	History items created	-
	Upload File	1: E18 PE.1 Reads	
Galaxy	This tool cannot be used in workflows	Treat as input dataset	
Tools Opt	·,		History Lists
search tools	Upload File	2: E18 PE.2 Reads	Saved Histories
Get Data	This tool cannot be used in workflows	Treat as input dataset	Histories Shared with Me
Send Data ENCODE Tools			Current History
Lift-Over	EASTO Croomer		Create New
Text Manipulation Convert Formats	FASTQ Groomer	3: E18 PE.1 Reads Groomed	Share or Publish
FASTA manipulation	Include "FASTQ Groomer" in workflow		Extract Workflow
Filter and Sort			Dataset Security
Extract Features	FASTQ Groomer		Show Deleted Datasets
Fetch Sequences	Include "EASTO Creemer" in workflow	4: E18 PE.2 Reads Groomed	Show Hidden Datasets
Fetch Alignments	Include FASTQ Groomer in worknow		Show structure
Operate on Genomic Intervals			Delete
Statistics	FASTQ Trimmer	5: E18 PE 1 Reads Groomed	14465082 14465083 T K 173 :
Graph/Display Data	Include "EASTO Trimmer" in workflow	Trimmed	14465083 14465084 G K 144 : 14465084 14465085 T T 117 •
Multiple regression) () () () ()
Multivariate Analysis			
Evolution Metagenemic analyses	FASTQ Trimmer	6: E18 PE.2 Reads Groomed,	erate pileup on 👁 🖉 🐹
EMBOSS	✓ Include "FASTQ Trimmer" in workflow	Trimmed	I-to-BAM on data 👁 🖉 🕱
NGS TOOLBOX BETA		-	
NGS: QC and manipulation	Man with Bowtie for Illumina		with Bowtie for
NGS: Mapping NGS: SAM Tools		7: Map with Bowtie for Illumina on	928 lines, format: sam,
NGS: Indel Analysis	✓ Include "Map with Bowtie for Illumina" in workflow	data 6 and data 5	se: mm9 equence file aligned.
NGS: Peak Calling	IN WORKHOW		Ø 🖻
RGENETICS			E 2.FLAG 3.1
SNP/WGA: Data; Filters	SAM-to-BAM	a call as Ball on data 7	\$269:3:1:1449:913 99 ch
SNP/WGA: QC; LD; Plots SNP/WGA: Statistical Models	✓ Include "SAM-to-BAM" in workflow	6: SAM-to-BAM on data /	5269:3:1:1449:913 147 chi . \$269:3:1:709:832 99 chi
Workflows			\$269:3:1:709:832 147 ch
TREATIONS			Carlotter and Carlotter
	Generate pileup	9 Cenerate nileun on data 8	
	☑ Include "Generate pileup" in workflow	or denerate preup on data o	

Galaxy Workflows



Transparency

Sharing, Collaborating, and Publishing with Galaxy

Sharing and Publishing

Sharing and Publishing History 'Variant Analysis for Sample E18'

Making History Accessible via Link and Publishing It

This history is currently restricted so that only you and the users listed below can access it. You can:

Make History Accessible via Link

Generates a web link that you can share with other people so that they can view and import the history.

Make History Accessible and Publish

Makes the history accessible via link (see above) and publishes the history to Galaxy's <u>Published Histories</u> section, where it is publicly listed and searchable.

Sharing History with Specific Users

You have not shared this history with any users.

Share with a user

Back to Histories List

Sharing and Publishing

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.g2.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 <u>Published Histories</u> 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
Galaxy | Published History | Variant Analysis for Sample E18 + Shttp://main.g2.bx.psu.edu/u/jgoecks/h/variant-analysis-for-sample-e18 C Q- Google 4 1 🚾 Galaxy Analyze Data Workflow Shared Data Visualization Help User About this History Published Histories | jgoecks | Variant Analysis for Sample E18 Import history Galaxy History ' Variant Analysis for Sample E18' Author Annotation: Perform a pileup analysis with default parameters to identify variants in sample E18. jgoecks Dataset Annotation **Related Histories** 1: E18 PE.1 Reads ۲ Forward reads from sample E18. All published histories Published histories by jgoecks 2: E18 PE.2 Reads Ð Reverse reads from sample E18. Rating 3: E18 PE.1 Reads Groomed ۲ Groom reads to convert quality scores from Solexa 1.0 to Solexa 1.3 Community ***** (1 rating, 4.0 average) Ð Groom reads to convert quality scores from Solexa 1.0 4: E18 PE.2 Reads Groomed Yours **** to Solexa 1.3 5: E18 PE.1 Reads Groomed, Trimmed ۲ Trim reads from 3' end to remove low-quality nts. Tags Community: 6: E18 PE.2 Reads Groomed, Trimmed ۲ Trim reads from 3' to remove low-quality nts. snp pileup bowtie demo sample 7: Map with Bowtie for Illumina on data 6 and data 5 ۲ Map paired-end reads with default parameters. Yours: 8: SAM-to-BAM on data 7 Ð Need to convert Bowtie SAM to BAM so that pileup snp x pileup x bowtie x analysis can be performed. demo 🗙 (sample:e18 🗙) 🆧 9: Generate pileup on data 8 ۲ Pileup analysis with default parameters 10: Filter pileup to get Variants from sample E18 ۲ Find variants with coverage >= 30. 13: Filter to get Variants from sample E18 where consensus base ۲ Filter pileup to find variants where the consensus base different than ref. base is different than the reference base. 14: UCSC mm9 RefSeq Genes ۲ UCSC mm9 RefSeq genes. 15: Intersect to get Variants from sample E18, consensus different, in ۹ Variants with consensus different that occur in RefSeq **RefSeq Genes** genes.

Transparency

Transparency



Transparency

Galaxy Published Histories								
	//main.g2.bx.psu.edu/history/list_published				CQ* Google		2	
🔁 Galaxy	Analyze Data Workflow	Shared Data	Visualization	Help User				
Published Hist	ories							
Name	Annotation	Owner	Community Rating †	Commu Tags	unity Last Updated			
Galaxy vs MEGAN	Comparison of Galaxy vs. MEGAN pipeline.	aunl	****	metag megar galaxy	n Mar 19, 2010			
metagenomic analysis		aunl	****	netag galaxy	y Mar 19, 2010			
<u>5M 1186088</u>	Datasets correspond to our paper published in Science by Peleg et al. entitled : Altered histone acetylation is associated with age-dependent memory impairment. Experiment layout: This history contains 4 datasets in the form of BED files of uniquely mapped reads produced after chip-seq for histone modifications H4K12ac and H3K9ac in mouse hippocampus of 3 months (young) and 16 months (old) mice after fear conditioning. For detailed information please refer to supplementary materials and methods of the respective work by peleg et al.	fischerlab	****	*	Apr 19, 2010			
<u>Variant Analysis</u> for Sample E18	Perform a pileup analysis with default parameters to identify variants in sample E18.	jgoecks	****	snp bowtie demo sampl	pileup e 2 minutes ago le			
get longest exon		henri	****	chr22 longe: marc humai works	st) exon Sep 02, 2010 n thop			
FASTA to Tabular Test		n	****	*	Aug 26, 2010			
EKLF		yzc109	****	*	Aug 24, 2010		*	
Open "http://main.g2.bx.ps	u.edu/history/list_published?sort=rating&f-tags=AII" in a new tab						1	

Galaxy Pages

A web-based, interactive medium for presenting all aspects of an analysis: data, methods, and results



00	Galaxy Published Page Variant Analysis for sample E18							
+	In http://main.g2.bx.psu.edu/u/jgoecks/p/variant-analysis-for-sample-e18	oogle)					
🗧 Gala	Analyze Data Workflow Shared Data Visualization Help User							
Published Pa	ges jgoecks Variant Analysis for sample E18		About this Page					
Varian	t Analysis of Embryonic Mouse Brain Tissue		Author jgoecks					
Jeremy Goeci	ks, Anton Nekrutenko, James Taylor, and The Galaxy Team		Related Pages					
Results To demonstra analysis exper tissue from da	ite how Galaxy can support accessible, reproducible, and transparent NGS re-sequencing studies, we performed a simple variant riment. This experiment identifies variants from a set of 4,536,964 RNA-seq reads obtained from sequencing a sample of mm9 brain ay 18 of embryonic development.		All published pages Published pages by igoecks Rating					
The initial ana determined by 2796 occur in	alysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus baseas y the MAQ modeldiffers from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, n known RefSeq Genes. These potential variants are:		(0 ratings, 0.0 average) Yours					
•	Galaxy Dataset Intersect to get Variants from sample E18, consensus different, in RefSeq Genes Image: Consensus different that occur in RefSeq genes. Variants with consensus different that occur in RefSeq genes. Image: Consensus different that occur in RefSeq genes.		Tags Community: none					
Method			Yours:					
In the first ste were trimmed grooming and and was filter	ep of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After I trimming, the reads were mapped using the short-read mapper Bowtie [2]. A pileup analysis using SAMtools [3] was then performed ed to identify variants supported by 30+ reads. The complete analysis is contained in this history:	ľ	•					
•	Galaxy History Variant Analysis for Sample E18 Perform a pileup analysis with default parameters to identify variants in sample E18.							
Here is a worl	kflow for performing this analysis:							
•	Galaxy Workflow Variant identification within annotated genes from NGS PE Data Identify variants in annotated genes from NGS paired-end data.							
Referenc	es							
[1] Han, X. et al. Transcriptome of embryonic and neonatal mouse cortex by high-throughput RNA sequencing. Proceedings of the National Academy of Sciences 106, 12741–12746 (2009).								
[2] Langmead Genome Biol	, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. 10, R25 (2009).	Å						
121.11.11.44.44	The Connect Allowershillow from and California, Birleformatics 3P, 3030, 30300	121	1					

Transparency

\varTheta 🔿 🔿 Galaxy Published Page Variant Analysis for sample E18	
+ Chttp://main.g2.bx.psu.edu/u/jgoecks/p/variant-analysis-for-sample-e18	e
Galaxy Analyze Data Workflow Shared Data Visualization Help User	
Published Pages jgoecks Variant Analysis for sample E18 The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus baseas determined by the MAQ modeldiffers from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, 2796 occur in known RefSeq Genes. These potential variants are:	About this Page
<u>Galaxy Dataset Intersect to get Variants from sample E18, consensus different, in RefSeq Genes</u> Variants with consensus different that occur in RefSeq genes.	Igoecks Related Pages
Method In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper Bowtie [2]. A pileup analysis using SAMtools [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:	All published pages Published pages by igoecks Rating Community (0 ratings, 0.0 average)
Galaxy History Variant Analysis for Sample E18 Perform a pileup analysis with default parameters to identify variants in sample E18.	Yours ****
8: SAM-to-BAM on data 7 Image: SAM to BAM so that pileup analysis can be performed. 9: Generate pileup on data 8 Pileup analysis with default parameters 10: Filter pileup to get Variants from analysis from sample E18 Find variants with coverage >= 30. 13: Filter to get Variants from sample for the performed. Filter pileup to find variants where the consensus base is different than the reference base. 14: UCSC mm9 RefSeq Genes UCSC mm9 RefSeq genes. 15: Intersect to get Variants from sample E18, consensus different, in RefSeq genes. Variants with consensus different that occur in RefSeq genes.	Community: none Yours:
Here is a workflow for performing this analysis:	
References	1



00)	Galaxy	Dataset	Filter to get	Variants fr	rom san	nple E18 v	where consens	us base di	fferent than ref. base	•
	+ Shttp://r	main.g2.bx.psu.edu/u/jg	oecks/d/fdo	dbcfd1207f	bd2					C Q. Goog	le
- c	alavy		An:	lyze Data	Workflow	Shar	red Data	Visualization	Help	liser	
0	апалу			ayee bata		0.1161	co oata				-
Dataset	ataset Filter to get Variants from sample E18 where consensus base different than ref. base										About this Dataset
Galaxy	Galaxy Dataset ' Filter to get Variants from sample E18 where consensus base different than ref. base' 🛛 📟 😋 🚺										Author
Annotativ	on: Filter pileup t	to find variants where the		hase is diff	erent than t	he refer	ence base				igoecks
Annotatio	n. riter pileup i	to find variants where the	e consensos	Dase is unit	crem man	ine reren	ence base.				Jgoccus
	This dataset is la	arge and only the first m	egabyte is s	hown below	I Show all						
••											
chr10	14465082	14465083		×	173	176	60	35	00000.00	000000.00	
chr10	14465083	14465084	â	ř.	144	144	60	35	00000.00	T. TTTTT	
chr10	14465485	14465486	č	T	129	129	60	34	tSTTTTT	TTTTTTTTTTTTTTTTT	
chr10	19928287	19928288	G	A	135	135	60	36	AAAAAAA	алалалалалала	
chr10	19928468	19928469	с	T	132	132	60	35	T\$ttttt	TTttttttTTT	
chr10	19928494	19928495	с	T	138	138	60	37	TTTTtttt	ttttttTtttT	The second se
chr10	19928543	19928544	A	G	147	147	60	40	GgGgGgg	gGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	
chr10	28750217	28750218	c	T	138	138	60	37	TTttttt	TTTTTTTTTTTTTT	Tags
chr10	28750397	28750398	A	c	154	211	60	64	C\$.\$C\$C	scscscscsccccc	
chr10	28/50423	28750424	ç	ž	113	113	60	35	TSEETTT	CELETIATITTETT	Community: none
chr10	28750446	28750549	C.	č	255	255	60	90	cscscsc	soggooggooooogg	contracting. Home
chr10	28750640	28750641	T	č	165	165	60	46	Cececee	CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC	Yours:
chr10	28750746	28750747	ĉ	Ä	202	202	60	58	λλααααα	aaaaaAaaaaaaaaa	12
chr10	28750766	28750767	A	G	205	205	60	59	G\$g\$G\$g	SG\$ggggggggggGG	S.
chr10	28750769	28750770	т	C	175	175	60	49	ccccCccd	CCCCCcCccccccCc	*
chr10	28750924	28750925	с	T	182	217	60	64	TŞTŞttT	tttTtTtTtTtTtTTg	
chr10	28751092	28751093	a	A	147	0	60	123		,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	
chr10	28751096	28751097	a	A	212	0	60	119	g\$g\$,,,		
chr10	28751114	28751115	g	A	225	235	60	85	aaaa,aaa		
chr10	28918972	28918973		N N	114	114	60	75	6,6,6,6	a , ese ese a a a a a a a a a a a a a a a a	
chr10	28918975	28918976	a	Ä	177	0	60	63	.5.5.5		
chr10	28918995	28918996	č	м	154	154	60	48	a\$aaaaa	100000000000000000000000000000000000000	
chr10	33613489	33613490	G	A	82	114	60	30	.SASASAS	SAAAAAAAAAAAAAA	
chr10	36721501	36721502	G	K	129	129	60	43	TTTT		
chr10	36721507	36721508	с	Y	51	51	60	54	.\$.,,,,		
chr10	36721695	36721696	T	A	120	120	60	31	aşaşaşa	SaŞaŞaŞaŞaŞaŞaŞaŞ	
chr10	36805412	36805413	A	G	126	126	60	33	Geogeoge	SCECCECCECCECCEC	
chr10	36805605	36805606	A	G	120	120	60	31	<u>aaaaaaaa</u>	33333333333333333	
chr10	360531/6	36053366	2	I D	17	17	60	33			
chr10	36854675	36854676	7	K	99	00	60	48			
chr10	36854678	36854679	Â	ñ	94	94	60	48			
chr10	36855346	36855347	a	R	159	159	60	50		GGG	
chr10	36855350	36855351	a	R	156	156	60	56			
chr10	36855356	36855357	a	A	134	0	60	49		,	
chr10	36855366	36855367	g	G	52	0	60	40			
chr10	36855370	36855371	g	G	25	0	60	32	.\$.\$.\$.	\$.\$.\$	
chr10	36855665	36855666	A	м	11	11	60	39	c\$c\$c\$c	\$c\$,,,	1
chr10	39144591	39144592	G	C	126	126	60	33	ccccccc		1
chr10	42722938	42722939	С	Y	63	63	60	34	t\$t\$t\$t\$	St\$	1
										14 1	

Transparency

😝 🔿 🤭 Galaxy Publish	ed Page Variant Analysis for sample E18		
Http://main.g2.bx.psu.edu/u/jgoecks/p/variant-analysis-for-samelysis-	ple-e18	C Qr Google	
Galaxy Analyze Data Workfile	w Shared Data Visualization Help User		
Published Pages jgoecks Variant Analysis for sample E18 The initial analysis produced support for 27,742 possible variants. Of these possible determined by the MAQ modeldiffers from the reference base and (b) read covera 2796 occur in known RefSeq Genes. These potential variants are: • Galaxy Dataset Intersect to get Variants from sample E18, Variants with consensus different that occur	e variants, there are 5,625 where (a) the consensus b ge at the base is 30x or greater. Of these potential v consensus different, in RefSeq Genes ur in RefSeq genes.	About this Pag About this Pag Author Jgoecks	e Vice
Method In the first step of this analysis, the reads were groomed to convert their quality scores; se grooming and trimming, the reads were mapped using the short-read mapper Bow and was filtered to identify variants supported by 30+ reads. The complete analysis	res from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next [1] for this step's rationale and parameter choices. le [2]. A pileup analysis using SAMtools [3] was then is contained in this history: Impor	, the reads After performed thistory Yours	ages s by jgoecks rage)
Galaxy History Variant Analysis f Perform a pileup analysis with default parameters to	o <u>r Sample E18</u> identify variants in sample E18.	Tags	
8: SAM-to-BAM on data 7 Image: Convert Bow performed. 9: Generate pileup on data 8 Image: Convert Bow performed.	ie SAM to BAM so that pileup analysis can be lefault parameters	Community: no Yours:	ne
10: Filter pileup to get Variants from sample E18 Find variants with constraints with constraint	erage >= 30.		
E18 where consensus base different than ref. base	ariants where the consensus base is different than th	e	
14: UCSC mm9 RefSeq Genes UCSC mm9 RefSeq g 15: Intersect to get Variants from Image: Comparison of the second se	nes. us different that occur in RefSeq genes.		
Sample E18, consensus different, in Refseg Genes		÷	
Galaxy Workflow Variant identification within ann Identify variants in annotated genes from	otated genes from NGS PE Data NGS paired-end data.	00	
References		4	

Transparency





Galaxy Published Page Variant Analysis for sample E18									
I + Ontrol http://main.g2.bx.psu.edu/u/jgoecks/p/variant-analysis-for-sample-e18									
- Galaxy Analyze Data Workflow Shared Data Visualization Help User									
Published Pages jgoecks Variant Analysis for sample E18	About this Page								
Variant Analysis of Embryonic Mouse Brain Tissue	Author								
Jeremy Goecks, Anton Nekrutenko, James Taylor, and The Galaxy Team	innerks								
Results	Related Pages								
To demonstrate how Galaxy can support accessible, reproducible, and transparent NGS re-sequencing studies, we performed a simple variant analysis experiment. This experiment identifies variants from a set of 4,536,964 RNA-seq reads obtained from sequencing a sample of mm9 brain tissue from day 18 of embryonic development.	All published pages Published pages by jgoecks								
The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus baseas determined by the MAD modeldifference base and (b) read coverage at the base is 30x or greater. Of these potential variants	Rating								
2796 occur in known RefSeq Genes. These potential variants are:	Community (0 ratings, 0.0 average)								
Galaxy Dataset Intersect to get Variants from sample E18, consensus different, in RefSeq Genes Variants with consensus different that occur in RefSeq genes. □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □	Yours ****								
	Tags								
Method	Community: none								
In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper Bowtie [2]. A pileup analysis using SAMtools [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:	Yours:								
Galaxy History Variant Analysis for Sample E18 Perform a pileup analysis with default parameters to identify variants in sample E18.									
Here is a workflow for performing this analysis:									
Galaxy Workflow Variant identification within annotated genes from NGS PE Data Identify variants in annotated genes from NGS paired-end data.									
References									
[1] Han, X. et al. Transcriptome of embryonic and neonatal mouse cortex by high-throughput RNA sequencing. Proceedings of the National Academy of Sciences 106, 12741-12746 (2009).									
[2] Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biol 10, R25 (2009).									
[3] Li, H. et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics 25, 2078 –2079 (2009).									
Open "http://main.g2.bx.psu.edu/workflow/imp?id=58d16d45527990b7" in a new tab	· · · · · · · · · · · · · · · · · · ·								



Creating a Page



47



Creating a Page

								Galaxy											
< <u>►</u> +	🛃 http://main.g2	.bx.psu.e	du/page	e/edit_co	ontent?id=d2	523e005e	lec427	7				Reader	C Q+	Google				_	2
- <mark>_</mark> Gala>	ĸy				Analyze Dat	a Wor	kflow	Shared Data	i Visi	alization	Help	User							
Page Editor T	Title : Variant Anal	ysis for s	ample E	18													Save	Clos	e
B $I \times^2 \times_2$	j≣ i≣ 48 48	₽ ¢ ¶	8 🤱 🛙	3	Paragraph t	ype 🔻 🔳	insert Li	ink to Galaxy C	bject ⊽	Embed (ialaxy Ob	ect 🔻							
To demonst identifies va The initial a from the rel	To demonstrate how Galaxy can support accessible, reproducible, and transparent NGS re-sequencing studies, we performed a simple variant analysis experiment. This experiment identifies variants from a set of 4,536,964 RNA-seq reads obtained from sequencing a sample of mm9 brain tissue from day 18 of embryonic development.											(
			Er	nbedded	i Galaxy Dat	aset 'Vari	iants fr	rom sample E	18, con:	ensus dif	ferent, in	RefSe	q Genes'						
				[Do not	t edit this blo	ck; Galax;	y will fi	ill it in with th	e annota	ted datase	t when it	is disp	layed.]						
Method	I																		
In the first s exclude bas Bowtie [2].	step of this analys se pairs with low o A pileup analysis	sis, the re quality so using SAI	ads wer ores; se Mtools [e groom e [1] for 3] was ti	ed to convert this step's ra hen performe	their qua tionale ar d and was	ality sco nd para s filtere	ores from Sole ameter choices ed to identify v	xa 1.0 ti . After g ariants :	o Solexa 1 rooming a supported	.3/Fastqs ind trimm by 30+ r	inger. ing, th ads. T	Next, the re e reads were he complete	ads were mapped analysis	trimmed f using the is contain	rom 36bp short-rea ed in this	to 27bp ad mapper history:	to	
					Embedde	d Galaxy	History	y 'Variant Pile	up Ana	lysis for S	ample E	8'							
				[Do no	t edit this blo	ck; Galax	y will fi	ill it in with th	e annota	ted histor	y when it	is disp	layed.]						
Here is a w	orkflow for perfor	rming this	s analysi	is:															
			E	Embedde	ed Galaxy We	rkflow 'S	iNP ide	ntification w	thin an	notated ge	enes fron	NGS	PE Data'						
				(Do not	edit this bloc	k; Galaxy	will fil	ll it in with the	annotat	ed workflo	w when i	t is dis	played.]						
Referen	ices																		
[1] Han, X. (2009).	et al. Transcriptor	me of em	bryonic	and neo	natal mouse	cortex by	high-t	hroughput RN	A seque	ncing. Pro	ceedings	of the	National Aca	demy of :	Sciences 1	06, 1274	1-12746		
[2] Langmei	ad, B., Trapnell, C	., Pop, M	. & Salzi	berg, S.L	. Ultrafast an	d memory	y-efficio	ent alignment	of short	DNA sequ	uences to	the hu	man genom	e. Genom	e Biol 10,	R25 (200	9).		Ų
[3] Li, H. et	al. The Sequence	Alignme	nt/Map	format a	nd SAMtools.	Bioinform	matics 2	25, 2078 -202	9 (2009).									Ŧ

Al Opportunity: "Now What?"

What can I do with this dataset? When should I use this tool? What should be the next step in my analysis? What are the "best practice" workflows for my analysis?

"Now What?" Factors

Past work

- what have I already done? (personal history)
- what have other people already done? (community history)

Approach

- exploration vs. focused analysis
- Google approach vs. scientist approach?

Overview

Genomics

Galaxy

- accessible, reproducible, and transparent science
- on the cloud
- visual analytics

Reflections on Galaxy

Three Ways to Use Galaxy

1. Download and Run Locally

2. Public Website (<u>http://usegalaxy.org</u>)

3. Run on the Cloud

1. Download and Run Locally

No configuration needed but everything can be configured

Prominent local installations at:

- + Cold Spring Harbor Lab
- Dept. of Energy's Joint Genome Institute
- Harvard School of Public Health
- U of Texas System
- U of Oklahoma
- Netherlands Bioinformatics Center
- Oxford College

Requires computing resources, technical expertise, and maintenance

2. 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

3. Galaxy on the Cloud

For extended or particular resource needs

- customization necessary
- oscillating data volume

Limited informatics expertise or infrastructure

Data production and (no?) sharing

The big picture



Galaxy CloudMan

Complete solution for instantiating, running and scaling cloud resources with automatically configured Galaxy

Tools and reference datasets exceed Galaxy Main

No computational expertise needed

no configuration needed but completely configurable

Reproducibility ensured

- Automated configuration for machine image, tools, and data
- Dynamic but persistent storage

Start an Instance



Configure Your Instance



$\leftarrow \rightarrow G$

© ec2-50-16-1-149.compute-1.amazonaws.com/cloud

🚾 Galaxy Cloudman

Info: report bugs | wiki | screencast

☆ 🌂

Galaxy Cloudman Console

Welcome to Galaxy Cloudman. This application will allow you to manage this cloud instance and the services provided within. If this is your first time running this cluster, you will need to select an initial data volume size. Once the data store is configured, default services will start and you will be able to add and remove additional services as well as 'worker' nodes on which jobs are run.

	Terminate cluster		Add nodes 🔻	Remove nodes	Acces	s Galaxy
Sta	itus					
с	luster name:	Heteroplas	my study 🛃			
D	isk status:	50M / 100	0G (1%) 🕵			Autoscaling is off.
W	/orker status:	Idle: 0 Av	ailable: 0 Requested: 0			Turn <u>on</u> ?
s	ervice status:	Application	s 😑 Data 👄			
E	xternal Logs:	Galaxy Log	I			
	luster status	log				Û



The importance of sharing

Share entire Galaxy CloudMan cluster instances

Fully automated solution

Publish an analysis

In progress or otherwise

Use CloudMan as PaaS

Deploy your own tool and make it available

Snapshot your instance

- + Data
- Configuration

Scaling the infrastructure with the computation



Exercising elasticity



Al Opportunity: Smart Resource Usage

How long will this tool run?

How much will it cost?

How much computing power do I need for this analysis?



http://en.wikipedia.org/wiki/Project_triangle

User Support vs. Automated

Users like (or require) control

 cognitive models of dynamic computing: parallelization and autoscaling?

Automation requires tool profiling

- local vs. global, parallelization
- linear vs. non-linear (e.g. graphs)

Overview

Genomics

Galaxy

- accessible, reproducible, and transparent science
- on the cloud
- visual analytics

Reflections on Galaxy

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



• • •	Galaxy Published Visualization	GCC2011-1: Viewing and Navigating	C Goodle
- Galaxy	Analyze Data Workflow Share	d Data Visualization Help User	
Published Visualizations jeremy GCC2011-1: Viewing 630,000 [640,000 UCSC Main on Human: knownGene (chr19) ♥ UCSC Main on Human: all_est (chr19) ♥ UCSC Main on Human: phyloP46wayPrimates (chr19) ♥ 1 -1 h1-hESC Tophat Mapped Reads ♥	and chr19		670,000 680,000 Auto (Squish) Dense Histogram Auto (Squish)
630,000 640,000	650,000	660,000	670,000 680,000
Display a menu			1



HTS Analysis Challenges

Complex tools, parameter dependent

Analysis and visualization not integrated

Want to be able to experiment
Dynamic Filtering



Integrating Tools and Visualization

- Galaxy	Analyze Data Work	kflow Shared Data	Visualization A	Admin Help	User	
GCC3: Running Tools (hg19)	chr19 1,530,000	• 1	1,523,098 - 1,545,23	32 👂 🗩	1	1,540,000
III UCSC Main on Human: knownGene 🔻					L	
221tj.2	······		······		······	
h1-hESC Tophat mapped reads 👻						
Max Intron Length 150000 Min Isoform Fraction 0.5 Pre MRNA Fraction 0.05 Perform quartile normalization No Run on complete dataset Run on visible region	50000, 0.5, 0.05, No] v					
**************************************	······•	CUFF.139.	1 >>>>>>>> C	CUFF.140.1 >> CUFF.141.1 CUFF.141.1	••••••••••••••••••••••••••••••••••••••	

III b1 bECC accombled transgripts _ reg	ion-falli parameters-f150000 0 5 0 05 Nol		
Cufflinks	ion=[aii], parameters=[150000, 0.5, 0.05, No]	*	
Max Intron Length	150000		
Min Isoform Fraction	0.05		
Pre MRNA Fraction	0.05		
Perform quartile normalization			
EF 138 1		·····	
>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>		CUFF.139.1 >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>	CUFF.140.1 25 CUFF.141.1 2555 CUFF.142.1 555
Cufflinks - region=[chr19:15230	98-1545232], parameters=[150000, 0.05, 0.05,	, No] 🗢	
CUFF.3.1			
		CUFF.5.1 >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>	CUFF.7.1 >>>>
		••••••	CUFF.8.1 >>>



Experimentation in Trackster

Tools integrated in visualization environment

Dynamically filter:

- visually identify features that match ranges
- on whole dataset

Run tools:

- quickly on visible region
- on whole dataset

Al Opportunity: Smart Visual Analytics

What datasets should I include in my visualization?

What are the interesting areas of my visualization?

What tools should I use in my visualization?

Combining Models

User modeling

interests + past activities + community activities

Data modeling

relationships amongst datasets

"Visual acuity" modeling: what can and should be seen in a visualization

Overview

Genomics

Galaxy

- accessible, reproducible, and transparent science
- on the cloud
- visual analytics

Reflections on Galaxy

Accessibility

Accessibility drives usage

- transparency and esp. reproducibility are secondary
- few incentives for reproducible research (for now)

A win-win for everyone

- tool developers: free GUI, more exposure
- users: easy to run tools, consistent interface
- everyone: amplification b/c tools can be chained to create complex analyses

Approach

Computer science research is driven by scientific needs

- we do biology research ourselves with Galaxy
- publish largely in biology journals

Don't need do everything, but what is done is done well

software engineering, community support are priorities

Galaxy used in and cited by papers in *Science, Nature, Genome Research, Genome Biology,* and more

Challenges Going Forward

Promoting community involvement

- tools, assays, analyses growing too fast for us alone
- facilitate community contributions and usage of contributions

Scaling to many, many Galaxies

- moving objects between Galaxies while ensuring reproducibility
- enabling users to find useful "stuff"

Novel application areas

genomics ideal application area -- what next?

Thanks! Questions, discussion?

https://bitbucket.org/galaxy/galaxy-central/wiki/Citations

Public server: http://usegalaxy.org Download and run: http://getgalaxy.org On the cloud: http://usegalaxy.org/cloud

jeremy.goecks@emory.edu