### Using RNA-seq for gene annotation, quantitation, and functional comparison in non-model organisms

Jeremy Goecks, Anton Nekrutenko, The Galaxy Team, and James Taylor

### The Question

What can I do with my Illumina RNA-seq reads from my non-model organism(s) using open source software?

### **The Question**

GBs of short, paired-end seq reads What can I do with my Illumina RNA-seq no gene annotation, perhaps no genome reads from my non-model organism(s) freely available to all, reproducible using open source software?



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

## **Current Projects**

Study of organism infection

- RNA-seq data from infected, uninfected individuals
- genome but no annotations

Study of dietary adaptations

- RNA-seq data from 4 related species; single lane of Illumina 76bp paired-end reads per species
- no genome or annotations

Study of parasite venoms

- RNA-seq data from 3 related parasite venoms, two lanes of Illumina 100bp paired-end reads per venom
- no genome or annotations



### **Desired Outputs**

For each organism

 assembled transcripts, quantitated and annotated

For pairs/groups of organisms

- \*log clusters
- differential expression amongst genes and/or \*logs

### What is Galaxy?

Web-based GUI for genomics that requires only a Web browser for everything: analysis, workflows, sharing, publication, and visualization

A public web service (http://usegalaxy.org) 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

# Galaxy in a Nutshell

### What you can do in Galaxy

- analysis interface, tools and datasources
- data libraries
- + workflows
- visualization
- + sharing
- + Pages

### Where you can use and build Galaxy

- public website
- local instance
- + on the cloud
- tool shed/contributing tools





### **Preprocessing Reads in Galaxy**



# Quality Statistics and Box NGS TOOLBOX BETA NGS: OC and manipulation ILLUMINA DATA FASTO Groomer convert

Graph/Display Data

and graph types

columns

Histogram of a numeric column

Plotting tool for multiple series

Quartiles

Scatterplot of two numeric

Boxplot of quality statistics

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



### FastQC



### FastQC



### FastQC



## **Quality Filtering**

### Filter FASTQ

### FASTQ File:

7: FASTQ Trimmer on data 2

Requires groomed data: if your data does not appear here try using the FASTQ groomer.

### Minimum Size:

Ľ	_	
L	n	
	v	

### Maximum Size:

0

A maximum size less than 1 indicates no limit.

### Minimum Quality:

0	$\sim$	
U	U.	
-	-	

### Maximum Quality:

0.0

0

A maximum quality less than 1 indicates no limit.

Maximum number of bases allowed outside of quality range:

This is paired end data:

Quality Filter on a Range of Bases

Add new Quality Filter on a Range of Bases

Execute

Quality Filter on a Range of Bases

Quality Filter on a Range of Bases 1

Define Base Offsets as:

### Absolute Values \$

Use Absolute for fixed length reads (Illumina, SOLiD) Use Percentage for variable length reads (Roche/454)

### Offset from 5' end:

0				

Values start at 0, increasing from the left

Offset from 3' end:

Values start at 0, increasing from the right

Aggregate read score for specified range:

min score 🛟

Keep read when aggregate score is:

>= ‡

Quality Score:

0.0

0

Remove Quality Filter on a Range of Bases 1

Add new Quality Filter on a Range of Bases

Execute

## **Overview: Single Organism, Genome but no Annotations**



# Genome but no Annotations (using Galaxy)



### NGS: RNA Analysis

RNA-SEQ

- <u>Tophat</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>Cuffdiff</u> find significant changes in transcript expression, splicing, and promoter use

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

### Tophat

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:

Human (Homo sapiens): hg18 Canonical

Is this library mate-paired?:

Single-end \$

### RNA-Seq FASTQ file:

1: imported: h1-hESC..ple Dataset 🛟

Must have Sanger-scaled quality values with ASCII offset 33

### TopHat settings to use:

Use Defaults You can use the default settings or set custom values for any of Tophat's parameters.

### Execute

### Cuffdiff

Transcripts: 29: Cuffcompare on da..transcripts

A transcript GTF file produced by cufflinks, cuffcompare, or other source.

Perform replicate analysis:

Perform cuffdiff with replicates in each group.

SAM or BAM file of aligned RNA-Seq reads: 11: Tophat on data 9:..cepted\_hits

SAM or BAM file of aligned RNA-Seq reads:

False Discovery Rate:

0.05 The allowed false discovery rate

Min Alignment Count:

The minimum number of alignments in a locus for needed to conduct significance testing or

### Perform quartile normalization:

No 
Removes top 25% of genes from FPKM denominator to improve accuracy of differential expre

Perform Bias Correction:

Yes 🛟

Bias detection and correction can significantly improve accuracy of transcript abundance est

Reference sequence data:

Set Parameters for Paired-end Reads? (not recommended):

No 🛟

Execute

### Cufflinks

SAM or BAM file of aligned RNA-Seq reads:

Max Intron Length: 300000

Min Isoform Fraction:

Pre MRNA Fraction:

Perform quartile normalization:

Removes top 25% of genes from FPKM denominator to improve accuracy of differential expression calls for low a

Use Reference Annotation:

Perform Bias Correction:

Yes 🗘

Bias detection and correction can significantly improve accuracy of transcript abundance estimates.

Reference sequence data:

Set Parameters for Paired-end Reads? (not recommended):



### Cuffcompare

GTF file produced by Cufflinks: 21: Cufflinks on data..transcripts

Additional GTF Input Files

Additional GTF Input Files 1

GTF file produced by Cufflinks: 18: Cufflinks on data..transcripts

(Remove Additional GTF Input Files 1)

\_\_\_\_\_

Add new Additional GTF Input Files

Use Reference Annotation:

Use Sequence Data:

Use sequence data for some optional classification

Choose the source for the reference list:

Execute

20





## Quantitation

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

# Quantitation and Differential Expression

- -> C fi 🕓 main.g2.bx.psu.edu

× +

Galaxy

- Galaxy

. . .

Analyze Data Workflow Shared Data Visualization Admin Help Use

tracking_id	class_code nearest_ref_id g	e_id gene_short_name tss_id locu	us length coverage	status	q1_PPKM q1_con	f_10 q	1_conf_hi	q2_FPKM q2_conf_lo	q2_conf_hi
TCONS 00000001	XLOC 000001 -	<ul> <li>chr11:133735175-133</li> </ul>	735361 186 -	- OK	8142.9 0	20053.7 4	474.38 0	11703	
TCONS 00000002	o NM 001194 XLOC 0000	NM 001194 TSS1 chr	19:531860-541886	1112 -	OK nan	0 n	an nan	0 nan	
TCONS 00000003	0 NM 001194 XLOC 0000	NM 001194 7551 chr	19:531860-541886	1116 -	OK nan	0	an nan	0 nan	
TCONS 00000004	x NM 005035 XLOC 0000	NM 005035 7663 chr	19:572537-583491	1650 -	PAIL 0	0 0	nan	0 nan	
700NE 0000005	× NM 005035 ¥LOC 0000	NM 005035 7563 ohr	10.572537-593491	1654	PATE 0	0 0	0.00	0 000	
700%5 00000005	<ul> <li>NM 001040134</li> <li>XLOC 0000</li> </ul>	NN 001040134	10.601603_603150	475	08 240 54	0 0 5	80 307 672 002	0 1600 51	
TCONS_0000000	e NA_001040134 ALOC_0000	SA_ODIO40134 - Chr.	191001003-002130	973 -	08 240.34		69.391 612.991	0 1000.51	
TCONS_0000007	XLOC_000006 -	T554 Chr191/2/023-/31220	322 -	OK 1450.31	0 3472.1		0		
700NS_00000008	XLOC_000007 -	TSS5 Chr191/35956-740458	231 -	487.708	0 1181.4	10 0	U		
TCONS_0000009	XFOC_000008 -	TSS6 chr19:783179-785122	459 - 0	OK 221.808	0 543.43	30 0	0	2011 C	
TCONS_00000010	o NM_001928 XLOC_0000	NM_001928 TSS7 chr.	19:797406-812366	3305 -	PAIL 0	0 0	0	0 0	
TCONS_00000011	e NM_001972 XLOC_0000	NM_001972 - chr	19:797406-812366	171 =	PAIL 0	0 0	178.62	0 506.607	
TCONS_00000012	XLOC_000009 -	<ul> <li>chr19:797406-812366</li> </ul>	409 - 1	FAIL 0	0 0	39457.3 2	378.03 76536.5		
TCONS 00000013	o NM 001928 XLOC 0000	NM 001928 - chr	19:812949-813258	309 -	OK 1250.8	6 0 2	959.41 0	0 0	
TCONS 00000014	- XLOC 000010 -	TSS8 chr19:860913-867823	745 - 4	OK 171.777	0 429.18	6 87.8122 0	249.425		
TCONS 00000015	XLOC 000010 -	<ul> <li>chr19:860913-867823</li> </ul>	439 - 4	OK 504.984	0 1203.9	123.554 0	349.431		
TCONS 00000016	0 NM 138690 XLOC 0000	NM 138690 75510 chr	19:932490-971944	1047 -	OK nan	0 n	an nan	0 nan	
70018 00000017	0 NM 138690 YLOC 0000	NM 139690 75610 chr	19,932490-971944	1079	08 030	0	40 040	0 000	
70000 00000019	- NH 010113	NN 010112 20010 011	10.00/100 00/01/	1483	DATE O		un nun	C nan	
70085 00000018	0 NM_019112 AD0C_0000	NA_019112 TSB12 Chr.	191964330-994324	1483 -	PALL 0		04076 7	Tan and tanks of	
TCONS_0000019	6 NH_012292 XLOC_0000	NR_012292 T5513 Chr.	19:1026593-1039065	941 -	FAIL 0	0 0	20376.7	733.016 40019.5	
TCONS_00000020	o NM_012292 XLOC_0000	NM_012292 TSS13 chr.	19:1026593-1039065	2163 -	PAIL 0	0 0	nan	0 nan	
TCONS_00000021	o NM_012292 XLOC_0000	NM_012292 - chr	19:1026593-1039065	259 -	FAIL 0	0 0	8770.53	0 19819.6	
TCONS_00000022	e NM_012292 XLOC_0000	NM_012292 - chr:	19:1026593-1039065	280 -	PAIL 0	0 0	497.775	0 1261.4	
TCONS_0000023	x NM_002695 XLOC_0000	NM_002695 TSS17 chr:	19:1042075-1043800	686 =	OK 189.92	4 0 4	74.007 169.841	0 437.158	
TCONS 00000024	x NM 002695 XLOC 0000	NM 002695 75517 chr	19:1042075-1043800	691 -	OK 92.922	5 0 2	51.555 213.744	0 \$37.177	
TCONS 00000025	x NM 002695 XLOC 0000	NM 002695 TSS18 chr	19:1042075-1043800	579 -	OK 342.83	2 0 8	27.993 37.2074	0 105.934	
TCONS 00000027	XLOC 000017 -	TSS20 chr19:1050958-10540	64 620 -	OK 442.107	0 1049.4	2 136,17 0	338,512		
TCONS 00000028	- XLOC 000017 -	TSS21 chr19:1050958-10540	64 315 - 0	OK 679.659	0 1621.2	7 445,941 0	1111.04		
70015 0000030	0 NM 001039848 XLOC 0000	NM 001039848 75623 chr	19:1056073-1056456	297	08 334.04	4.0 8	40.68 264.54	0 706.481	
TCONE 0000031	× NH 001100122 XLOC 0000	NW 001100133 //0036 chr	10.1050670-1050005	124	08 1066 0	5 0 A	041 07 1606 05	0 3057 47	
10048 00000031	- NH 001100122 ALOC 0000	NH 001100122 10025 CHL	1911030079-1050005	493	08 2316 1	5 6 6	103 34 006 301	0 0100 55	
100NS_00000032	x NN_001100122 XLOC_0000.	NA_OUIIOOI22 TSS2/ Chr.	19:1064990-1065509		08 2213.3		103.74 070.201	0 2120.55	
TCONS_00000033	x NM_001100122 XLOC_0000.	NR_001100122 T5529 Chr.	19:106/336-10/4/23	1109 -	OK nan	o n	an nan	o nan	
TCONS_00000034	x NM_001100122 XLOC_0000	NM_001100122 TSB30 chr	19:1067336=1074723	765 -	OK 1813.5	20 4	158,93 3727.42	0 7926.98	
TCONS_0000035	XLOC_000024 -	TSS31 chr19:1074797-10829	36 916 - 0	OK nan	0 nan	nan 0	nen		
TCONS_00000036	1 NM_014963 XLOC_0000	NM_014963 TSS32 chr:	19:1074797-1082936	464 -	OK 863.28	70 2	049.43 1128.66	0 2673.69	
TCONS_00000037	i NM_014963 %LOC_0000	NM_014963 TSS34 chr:	19:1082990-1086597	286 -	FAIL 0	0 0	0	0 0	
TCONS_00000038	i NM_014963 XLOC_0000	NM_014963 75535 chr	19:1082990-1086597	1050 -	FAIL 0	0 0	0	0 0	
TCONS 00000039	i NM 014963 XLOC 0000	NM 014963 - chr:	19:1082990-1086597	948 -	FAIL 0	0 0	0	0 0	
TCONS 00000040	x NM 014963 XLOC 0000	NM 014963 TSS36 chr	19:1104001-1106745	894 -	OK nan	0 n	an nan	0 nan	
TCONS 00000041	x NM 014963 XLOC 0000	NM 014963 75537 chr	19:1104001-1106745	189 -	OK 4297.2	8 0 1	0297.2 3281.36	0 7773.26	
TCONS 00000042	O NR 023312 XLOC 0000	NR 023312 76639 chr	19-1206374-1220506	1137 -	08 040	0	an nan	0	
TCONS 00000043	0 NR 023312 XLOC 0000	NR 023312 75540 chr	19:1220567-1227635	- 022	PATL 0	0 0	1200.49	0 2848.08	
TCONS 0000044	0 NM 017914 TLOC 0000	NW 017914 75540 chr	19-1220567-1227635	768	PATT 0	0 0	89 8201	0 243 688	
20082 0000044	<ul> <li>NR 033310</li> <li>NR 033310</li> <li>NLOC 0000</li> </ul>	ND 033312	10.1000557 1007635	227	DATE O		4010 2	0 0530 07	
TCOAS_0000045	e ak_023312 ALOC_0000.	NR_023312 - Chr.	1911220367-1227635		FAIL 0		4010.2	9330.07	
TCONS_0000046	C NR_023312 XLOC_0000	NR_023312 - Chr.	1911220567-1227635	1108 -	FAIL 0	0 0	nan	0 nan	
TCONS_0000047	e NR_023312 XLOC_0000	NR_023312 - chr.	19:1220567-1227635	435 -	PAIL 0	0 0	289.05	0 715.339	
TCONS_0000048	e NM_017914 XLOC_0000	NM_017914 - chr	19:1227903-1228515	612 -	OK 203.13	3 0 4	83.744 413.155	0 1000.25	
TCONS_00000049	1 NM_001405 XLOC_0000	NM_001405 TSS42 chr:	19:1241801-1244814	628 -	OK 425.35	90 1	014.39 0	0 0	
TCONS_00000050	o NM_001405 XLOC_0000	NM_001405 TSS43 chr	19:1250355-1254081	453 -	OK 610.51	2 0 1	457.01 114.849	0 323.665	
TCONS 00000051	- XLOC 000029 -	<ul> <li>chr19:1250355-125400</li> </ul>	81 471 - 0	OK 564.127	0 1351.7	5 297.756 0	753.12		
TCONS 00000052	- XLOC 000030 -	TSS45 chr19:1254176-12572	85 1026 - 0	OK nan	0 nan	nan 0	nan		
TCONS 00000053	XLOC 000030 -	TS\$45 chr19:1254176-125721	85 1023 -	OK nan	0 nan	nan 0	nan		
TCONS 00000054	- XLOC 000031 -	TSS46 chr19:1269125-12748	07 578 -	PATE 0	0 0	0 0	0		
TCONS DODDODES	- XLOC 000031	T\$\$46 chr19:1269125-12740	07 517 -	PATE 0	0 0	0 0	0		
1000HE 00000055	×100 000031	20040 abold.1360135_13346	07 1667	DATE O	0 0	0 0	0		
20010 0000056		20040 che10.1260125-12740	07 1670	DATE O	0		0		
TCONS_00000057	XLOC_000031 -	T5549 Chr19:1269125-12748	0/ 10/9 - 1	FAIL 0	0 0	0 0	0		
TCONS_0000060	- XLOC_000032 -	T8852 chr19:1275597-12791	16 690 - 0	OK 2628.87	5873.2	6 4294.15 0	9493.66		
TCONS_0000061	XLOC_000033 -	TSS54 chr19:1356344-13607	18 477 - 1	PAIL 0	0 0	38.2618 0	110.037		
000000 00000000	110 110 11 FLOO 0000	NH 136313 36655 about	10.130/344 13/0310	1617	ILATE O			0	

# Quantitation and Differential Expression

tracking_1d	Class_	code nea:	rest_ref_	1d gene_1d	gene_sho	ort_name	p tss_10	1 locus	lengt	h covers	ge	status	Q1_PPKM	q1_conf	10	q1_conf.	hì	QZ_FPKM	q2_conf_lo	q2_cont_hi
TCONS_00000001	- 0.07	- XLO	C_000001		-	chr11:	13373517	5-133735	361	186	-	OK	8142.9	0	20053.7	4474.38	0	11703		
TCONS_00000002	0	NM_001194	XLOC	000003	NM_00119	94	<b>TSS1</b>	chr19	531860-	541886	1112	-	OK	nan	0	nan	nan	0	nan	
TCONS_00000003	0	NM 001194	XLOC.	000003	NM_00115	9-6	TSS1	chr19:	531860-	541886	1116	-	OK	nan	0	nan	nan	0	nan	
TCONS 00000004	x	NM 005035	XLOC	000004	NM 00503	35	TSS3	chr19:	572537-	583491	1650	-	PAIL	0	0	0	nan	0	nan	
TCONS_00000005	×	NM_005035	XLOC	000004	NM_00503	35	<b>TSS3</b>	chr19:	572537-	583491	1654	-	PAIL	0	0	0	nan	0	nan	
TCONS 00000006	•	NM 001040134	XLOC	000005	NM 00104	40134	-	chr19:	681683-	682158	475	-	OK	240.569	0	589.397	672.907	0	1600.51	
TCONS_00000007	-	- XLO	000006	-	TSS4	chr19:	27023-7	31220	322	-	OK	1450.31	0	3472.77	0	0	0			
TCONS_0000008	-	- XLOO	000007	-	T885	chr19:	135954-7	40458	231	-	OK	487.708	0	1181.41	0	0	0			
TCONS_0000009	-	- XLO	000008	-	TSS6	chr19:	183179-7	85122	459	-	OK	221.808	0	543.433	0	0	0			
TCONS 00000010	0	NM 001928	XLOC	000009	NM 00192	28	TSS7	chr19:	797406-	812366	3305	-	PAIL	0	0	0	0	0	0	
TCONS_00000011	e	NM_001972	XLOC	000009	NM_00197	72	-	chr19:	797406-	812366	171	-	PAIL	0	0	0	178.62	0	506.607	
TCONS 00000012	-	- XLO	000009	-	-	chr19:	797406-8	112366	409	-	FAIL	0	0	0	39457.3	2378.03	76536.5			
TCONS 00000013	0	NM_001928	XLOC	000009	NM_00192	28	-	chr19:	812949-	813258	309	-	OK	1250.86	0	2959.41	0	0	0	
200NE 0000014	62236	VI.O/	010000			abr 19.	60912-9	67973	745		OF	171 777	0	420 106	07 0177		240 475	16.20	52	

# Filtering

# Filter for differentially expressed elements

Filter combined transcripts for those that are differentially expressed

Filter
Filter:
130: Cuffcompare on datranscripts
Dataset missing? See TIP below.
With following condition:
c14=='yes'
Double equal signs, ==, must be used



### Overview: de novo



### **De Novo Assembly**

Assemble reads by looking for overlap amongst reads and contigs

k = starting overlap between reads

### Highly resource intensive

- lots of memory and multiple cores required
- difficult to do with desktop computer

Wrappers coming to Galaxy soon

# **De Novo RNA-seq Assemblers**

	Trans-ABySS	Trinity
basic idea	combine multiple DNA assemblies	single assembly oriented toward RNA
runtime (24 cores)	7-9 days	1-2 days
disk space required	huge	small
scaffolding via PE reads	yes	no
strand-specific	no	yes

Robertson et al. (2010). "De novo assembly and analysis of RNA-seq data." Nature Methods 7(11): 909-912. Grabherr et al. (2011). "Full-length transcriptome assembly from RNA-Seq data without a reference genome." Nat Biotechnol.

# **De Novo RNA-seq Assemblers**

	Trans-ABySS	Trinity
basic idea	combine multiple DNA assemblies	single assembly oriented toward RNA
runtime (24 cores)	7-9 days	1-2 days
disk space required	huge	small
scaffolding via PE reads	yes	no
strand-specific	no	yes

Robertson et al. (2010). "De novo assembly and analysis of RNA-seq data." Nature Methods 7(11): 909-912. Grabherr et al. (2011). "Full-length transcriptome assembly from RNA-Seq data without a reference genome." Nat Biotechnol.

# **De Novo RNA-seq Assemblers**

	Trans-ABySS	Trinity
basic idea	combine multiple DNA assemblies	single assembly oriented toward RNA
runtime (24 cores)	7-9 days	1-2 days
disk space required	huge	small
scaffolding via PE reads	yes	no
strand-specific	no	yes

Robertson et al. (2010). "De novo assembly and analysis of RNA-seq data." Nature Methods 7(11): 909-912. Grabherr et al. (2011). "Full-length transcriptome assembly from RNA-Seq data without a reference genome." Nat Biotechnol.

>comp17428\_c0\_seq1\_FPKM\_all:10.443\_FPKM\_rel:10.443\_len:198\_path:[0] AGAAAACTTTTTTTTGGTAAAAACAAAACAATTTTCTATTTTTGCAGTAAAATTTACGAAT CATAATTTTGGGGGGAAAATTTTTGATATTAATGTAAACTTGTCATAAGAGGGGGAAAATTG TGATACATTTTCCGCCATCTTGTGGGACAAGCAGAAAGTATTATTTGACATTTCGTAAAT TTTACTGTGCGGCTTGTT

>comp17416\_c0\_seq1\_FPKM\_all:13.273\_FPKM\_rel:13.273\_len:642\_path:[0] TCCTGACAAACAACGCCCCCTGCTTCGTTTGGCGAACGTTTTGGTACAACTTTAACTCCA CAATCCGCAGATCTTGATAAACAAAATCCAAGAAAAATAGGATCAGCAATTTTAACAAGA ATATTGTCTACTGAAAATACATGTACACTATTTATTCCGCGCTTCTTCATATTATCCAAT ACGCCTTGGATTTTTAATGCACGATAAAGTCCTCCATTTCCATCTGGTGCTTTTGATATG TGACATTTTTCATCCAATATTATTTTTACCATCAAAAGTAAAACAGGGGTAGCATCTCTTGT TTAAATGCTTTTATATTTTCTTTCTTTAATCCAAAAGTAAAACAGGGGTAGCATCTCTTGT TTGAATTGTTCGTGAGTTGCTTCACTTGTCATAATATACCATGTAATATGACCATGTTTT TTGAATTTTTCTTCAGCAAGTTCTTGTAGTTTACGAATTCTTAAAGCTTGCAAACGAAAT AGAGTGCTGTGTGAAGGAAGTCCAACATCATACATTCCCTTGGGATAATTTACACCAAGT CGAGTACCTTGACCTCCAGCCAATAAAAGTATAGCGACTCGACCATCGGAAATTTCCTTT

### **Evaluating Assemblies**

### Expectations

80% reassembly = success

### Statistics

hgʻ

n	n:100	n:N50	min	median	mean	N50	max	sum
650976	64671	6276	100	198	618	1744	17063	4.00E+07

### Transcript length distribution

0	mean	stdev	0%	25%	50%	75%	100%
	1722	1352	30	832	1497	2208	81928

兽 🔮 🌒 🧮 Galaxy   Accessible Workflow 🛪	+		
← → C n Smain.g2.bx.psu.edu/u/	jeremy/w/genemrna-length-statistics		む く
T Galaxy	Analyze Data Workflow Shared Data Visualization Admin Help User		
Accessible Workflow   Gene/mRNA length statistics			
Galaxy Workflow ' Gene/mRNA length statistic	's'	Author	1
Step	Annotation	jeremy	10
Step 1: Input dataset	Requirements: (a) one record per coding element; (b) element names have form [parent_name]]delimiter][element_name]	Related Workflows	1
BED file of coding elements select at runtime		All published workflows Published workflows by jeremy	
Step 2: Convert	Separate parent name and element name. Typically choose '.' for gene	Rating	
Convert all	annotation files and '_' for mRNA annotation files.	Community (0 ratings, 0.0 average)	*****
in Query		Yours	$(a,a) \in (a,a)$
Output dataset 'output' from step 1		Tags	
Step 3: Compute	Compute length of each element.	Community: none	
Add expression c3-c2		Yours:	
as a new column to			
Round result?			
Step 4: Group	Create table of parent element lengths. For gene datasets, typically want to sum on column 8, for mRNA datasets, typically want to column 13		
Select data Output dataset 'out_file1' from step 3	sum on column a, for move datasets, typicary want to sum on column 23.		
Group by column 4 (value not yet validated)			
Ignore case while grouping? False			
Operations			
Operation 1			
Sum			
On column			
Round result to nearest integer? NO			
Step 5: Summary Statistics	Compute statistics across all parent elements,		
Summary statistics on Output dataset 'out_file1' from step 4			
Column or expression			



### Sequence Annotation de novo

- 1. If necessary, extract genomic DNA
  - using GFF/features!

### 2. Annotation via HMMScan

- PFAM database
- http://hmmer.janelia.org/





				<u>A</u>
€⇒Cπ <u>©</u>	nmmer.janella.org/results/score/11529264-B180-11E0-95F6-9859998A/913/ptam=1		_	12
	Fram Domains			
	Show hit details			
	_ Distribution of Significant Hits @			
	Bacteria Eukaryota Archaea Viruses Unclassified Sequences Other Sequences     E	irst « Previous Page )	1 of 31 No	xt » Last »
Query Matches (30)	26)			Customize
Target	Description	Species	E-value	Show Alignment
A2RQD6_HUMAN	Bcr-abl1 e6a2 chimeric protein (Fragment) (gene: BCR-ABL1)	Homo sapiens	9.0e-107	show
ABL_FSVHY@	Tyrosine-protein kinase transforming protein Abl (gene: ABL)	Feline sarcoma virus (strain Hardy- Zuckerman 2)	2.2e-106	show
A2RQD7_HUMAN@	Bcr-abl1 e19a2 chimeric protein (Fragment) (gene: BCR-ABL1 e19a2)	Homo sapiens#	2.5e-106	show
A2AV22_MOUSE	V-abl Abelson murine leukemia oncogene 1 (gene: Abl1)	Mus musculus 🗗	2.6e-106	show
Q3SYK5_MOUSE	Abl1 protein (gene: Abl1)	Mus musculus 🗗	2.6e-106	show
D3ZGM3_RAT	Uncharacterized protein (gene: Abl1)	Rattus norvegicus 🛱	2.7e-106	show
Q2PYT4_RATE	ABL1 (Fragment) (gene: Abl1)	Rattus norvegicus @	2.7e-106	show
ABL1_HUMANE	Tyrosine-protein kinase ABL1 (gene: ABL1)	Homo sapienst?	2.7e-106	show
D2H177_AILMELS	Putative uncharacterized protein (Fragment) (gene: PANDA_003253)	melanoleucat	2.86-106	snow
OS9EK4 HUMANI	V-abl Abelson murine leukemia viral oncorene homolog 1 isoform b variant (Fragment)	Homo sapiens	2.8e-106	show
Q4SJH9_TETNG	Chromosome 4 SCAF14575, whole genome shotgun sequence. (Fragment) (gene: GSTENG00017201001)	Tetraodon nigroviridis 🗗	3.7e-106	show
A9UF02_HUMAN	BCR/ABL fusion protein isoform X9 (gene: BCR/ABL fusion)	Homo sapienst?	4.5e-106	show
A9UF07_HUMANE	BCR/ABL fusion protein isoform Y5 (gene: BCR/ABL fusion)	Homo sapiens	5.0e-106	show
Q2PYT3_RATE?	ABL1 (Fragment) (gene: Abl1)	Rattus norvegicus 🗗	1.7e-105	show
BOUXN6_DANRE	Novel protein similar to vertebrate Abelson murine leukemia viral oncogene homolog 2 (Arg, Abelson-related gene) (ABL2) (gene: abl2)	Danio rerio 🗗	1.9e-94	show
BOUXN7_DANRE	Novel protein similar to vertebrate Abelson murine leukemia viral oncogene homolog 2 (Arg, Abelson-related gene) (ABL2) (gene: abl2)	Danio rerio 🗐	5.3e-94	show
Q4RUE3_TETNG	Chromosome 1 SCAF14995, whole genome shotgun sequence. (Fragment) (gene: GSTENG00028837001)	Tetraodon nigroviridis 🗗	5.4e-94	show
Q6P282_XENLA®	MGC69056 protein (gene: abl2)	Xenopus laevis@	6.4e-94	show
D6RBS4_HUMAN	Uncharacterized protein (gene: ABL2)	Homo sapiens	1.0e-93	show
D6RIE2_HUMANE	Uncharacterized protein (gene: ABL2)	Homo sapiens	1.1e-93	show
B5MEB6_HUMAN	Uncharacterized protein (gene: ABL2)	Homo sapienst	1.3e-93	show
ABL2 HUMANIC	ADIZ ISOTOTTI IBLUIS	Homo sapiens®	2.08-93	show
BTUEFS HUMANIA	Abi2 isoform 18C/TS	Homo sapiensia	3.10-02	show
DTUETS HUMANIS	Auto London Alegore	Home capience	3.10-93	show

# \*log Clustering de novo

(Quasi-)Reciprocal BLA(s)T

- find best matches between two sets of transcripts
- repeat for more species

LastZ

- designed for long reads/contigs
- high-quality mapping amongst \*log clusters

### Quantitation de novo

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

### RSEM de novo

### Expectation-maximization model

### Uses TPM, not FPKM

 when mean transcript length is 1k, 1 TPM=1 FPKM

### Be consistent

Li, B., V. Ruotti, et al. (2010). "RNA-Seq gene expression estimation with read mapping uncertainty." Bioinformatics 26(4): 493-500.

k71 <b>:</b> 94571	34285.52	0.00129541301452086	k71 <b>:</b> 94571
k70 <b>:</b> 110756u	34714.33	0.00147532810537204	k70 <b>:</b> 110756u
k72 <b>:</b> 79128u	35163.91	0.00139980345826473	k72 <b>:</b> 79128u
k71 <b>:</b> 95729u	36039.45	0.00164266197886051	k71 <b>:</b> 95729u
k73 <b>:</b> 62819u	36888.41	0.0018459751145641	k73 <b>:</b> 62819u
k73 <b>:</b> 62364u	37274.56	0.00180992945323558	k73 <b>:</b> 62364u
k73 <b>:</b> 63082u	37919.99	0.00136337550194011	k73 <b>:</b> 63082u
k72:78714	38233.64	0.00237886076778222	k72 <b>:</b> 78714
k73 <b>:</b> 62636	38362.56	0.000375271843972168	k73 <b>:</b> 62636
k71 <b>:</b> 95841	39041.42	0.00189572242891473	k71 <b>:</b> 95841
k73 <b>:</b> 63028u	39759.97	0.0015827631212604	k73 <b>:</b> 63028u
k72:78108	40155.68	0.00184053788631075	k72 <b>:</b> 78108
k71:93994	40967.55	0.00177806349073446	k71 <b>:</b> 93994
k74 <b>:</b> 38786u	44067.87	0.00220525090976878	k74 <b>:</b> 38786u
k72:78714	38233.64	0.00237886076778222	k72 <b>:</b> 78714
k73 <b>:</b> 62636	38362.56	0.000375271843972168	k73 <b>:</b> 62636
k71 <b>:</b> 95841	39041.42	0.00189572242891473	k71 <b>:</b> 95841
k73 <b>:</b> 63028u	39759.97	0.0015827631212604	k73 <b>:</b> 63028u
k72:78108	40155.68	0.00184053788631075	k72:78108
k71:93994	40967.55	0.00177806349073446	k71:93994

### Differential Expression de novo

Simple: combine results from \*log clustering and quantitation

## Summary

For each organism

assembled transcripts, quantitated and annotated

For pairs/groups of organisms

- \*log clusters
- differential expression amongst genes and/or \*logs







Enis Afgan



**Dave Clements** 

Kanwei Li



Dannon Baker



Jeremy Goecks

**James Taylor** 





Jennifer Jackson



**Kelly Vincent** 



Nate Coraor



Greg von Kuster



Anton Nekrutenko

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

## **Thanks! Questions?**

Slides: <u>http://dl.dropbox.com/u/4745528/goecks\_rnaseq\_workshop\_2011.pdf</u>

Galaxy

- public server: http://usegalaxy.org
- download: http://getgalaxy.org

jeremy.goecks@emory.edu



