De novo RNA-Seq Assembly and Transcriptome Studies using Trinity

with Applications towards Non-model Organism Studies

Brian Haas

Broad Institute

Workshop on Genomics, Cesky Krumlov, Jan 2016

RNA-Seq Empowers Transcriptome Studies



Generating RNA-Seq: How to Choose?

Many different instruments hit the scene in the last decade



Slide courtesy of Joshua Levin, Broad Institute.

RNA-Seq: *How to Choose?*



Slide courtesy of Joshua Levin, Broad Institute.

Generating RNA-Seq: How to Choose?

Popular choices for RNA-Seq today







Generating RNA-Seq: How to Choose?



RNA-Seq: How do we make cDNA?

Prime with Random Hexamers (R6)



Slide courtesy of Joshua Levin, Broad Institute.

Overview of RNA-Seq



From: http://www2.fml.tuebingen.mpg.de/raetsch/members/research/transcriptomics.html

Common Data Formats for RNA-Seq

FASTA format:

>61DFRAAXX100204:1:100:10494:3070/1 AAACAACAGGGCACATTGTCACTCTTGTATTTGAAAAACACTTTCCGGCCAT

FASTQ format:

Quality values

```
AsciiEncodedQual(x) = -10 * log10(Pwrong(x)) + 33
```

AsciiEncodedQual ('C') = 64

So, $Pwrong('C') = 10^{(64-33/(-10))} = 10^{-3.4} = 0.0004$

Paired-end Sequences



@61DFRAAXX100204:1:100:10494:3070/2
CTCAAATGGTTAATTCTCAGGCTGCAAATATTCGTTCAGGATGGAAGAACA
+

Overview of RNA-Seq



From: http://www2.fml.tuebingen.mpg.de/raetsch/members/research/transcriptomics.html

RNA-Seq reads



Advancing RNA-Seq analysis

Brian J Haas & Michael C Zody

Nature Biotech, 2010

New methods for analyzing RNA-Seq data enable de novo reconstruction of the transcriptome.













RNA-Seq reads

Assemble transcripts de novo End-to-end **Transcriptome**-based **RNA-Seq Analysis** Software Package Trinity NATURE PROTOCOLS | PROTOCOL *De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis Align transcripts to genome Brian J Haas, Alexie Papanicolaou, Moran Yassour, Manfred Grabherr, Philip D Blood, Joshua Bowden, Matthew Brian Couger, David Eccles, Bo Li, Matthias Lieber, Matthew D MacManes, Michael Ott, Joshua Orvis, Nathalie Pochet, Francesco Strozzi, Nathan Weeks, Rick Westerman, Thomas William, Colin N Dewey, Robert Henschel, Richard D LeDuc, Nir Friedman & Aviv Regev Affiliations | Contributions | Corresponding authors Nature Protocols 8, 1494–1512 (2013) | doi:10.1038/nprot.2013.084 Published online 11 July 2013

The General Approach to *De novo* RNA-Seq Assembly Using De Bruijn Graphs

Sequence Assembly via De Bruijn Graphs

a Generate all substrings of length k from the reads

ACAGC TCCTG GTCTC	AGCGC CTCTT GGTCG]
CACAG TTCCT GGTCT	CAGCG CCTCT TGGTC	
CCACA CTTCC TGGTC TGTTG	TCAGC TCCTC TTGGT	
CCCAC GCTTC CTGGT TTGTT	CTCAG TTCCT GTTGG	k more (k=E)
GCCCA CGCTT GCTGG CTTGT	CCTCA CTTCC TGTTG	- k-mers (k=5)
CGCCC GCGCT TGCTG TCTTG	CCCTC GCTTC TTGTT CGTAG	
CCGCC AGCGC CTGCT CTCTT	GCCCT CGCTT CTTGT TCGTA	
ACCGC CAGCG CCTGC TCTCT	CGCCC GCGCT TCTTG GTCGT	
ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG	CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG	- Reads

b Generate the De Bruijn graph



From Martin & Wang, Nat. Rev. Genet. 2011

b Generate the De Bruijn graph





From Martin & Wang, Nat. Rev. Genet. 2011

Contrasting Genome and Transcriptome Assembly

Genome Assembly

- Uniform coverage
- Single contig per locus
- Double-stranded

Transcriptome Assembly

- Exponentially distributed coverage levels
- Multiple contigs per locus (alt splicing)
- Strand-specific



Trinity Aggregates Isolated Transcript Graphs

Genome Assembly

Single Massive Graph



Entire chromosomes represented.

Trinity Transcriptome Assembly

Many Thousands of Small Graphs



Ideally, one graph per expressed gene.

Trinity – How it works:



Thousands of disjoint graphs



Decompose all reads into overlapping Kmers (25-mers)

Identify seed kmer as most abundant Kmer, ignoring low-complexity kmers.

Extend kmer at 3' end, guided by coverage.





GATTACA 9 T C



GATTACA 9 T C



GATTACA 9 T₀ C



GATTACA 9 C₄ C₄



 $\mathbf{GATTACA}_{9} \qquad \mathbf{C}_{4} \qquad \mathbf{A}_{1} \\ \mathbf{T}_{0} \\ \mathbf{C}_{4} \qquad \mathbf{C}_{4}$
















Report contig:AAGATTACAGA....

Remove assembled kmers from catalog, then repeat the entire process.



Expressed isoforms



















Chrysalis Re-groups Related Inchworm Contigs



Chrysalis uses (k-1) overlaps and read support to link related Inchworm contigs







Butterfly Example 1: Reconstruction of Alternatively Spliced Transcripts





Reconstruction of Alternatively Spliced Transcripts



Reconstructed Transcripts



Reconstruction of Alternatively Spliced Transcripts



Reconstructed Transcripts



Reconstruction of Alternatively Spliced Transcripts



Butterfly Example 2: Teasing Apart Transcripts of Paralogous Genes





Teasing Apart Transcripts of Paralogous Genes



Strand-specific RNA-Seq is Preferred

Computationally: fewer confounding graph structures in de novo assembly: ex. Forward != reverse complement

(GGAA != TTCC)

Biologically: separate sense vs. antisense transcription

NATURE METHODS | VOL.7 NO.9 | SEPTEMBER 2010 |



Comprehensive comparative analysis of strand-specific RNA sequencing methods

Joshua Z Levin^{1,6}, Moran Yassour^{1-3,6}, Xian Adiconis¹, Chad Nusbaum¹, Dawn Anne Thompson¹, Nir Friedman^{3,4}, Andreas Gnirke¹ & Aviv Regev^{1,2,5}

Strand-specific, massively parallel cDNA sequencing (RNA-seq) is a powerful tool for transcript discovery, genome annotation

Nevertheless, direct information on the originating strand can substantially enhance the value of an RNA-seq experiment. For

'dUTP second strand marking' identified as the leading protocol

computational pipeline to compare library quality metrics from any RNA-seq method. Using the well-annotated *Saccharomyces cerevisiae* transcriptome as a benchmark, we compared seven library-construction protocols, including both published and transcribed strand or other noncoung to tris, demarcate the exact boundaries of adjacent genes transcribed on opposite strands and resolve the correct expression levels of coding or noncoding overlapping transcripts. These tasks are particularly challenging in small microbial genomes, prokaryotic and eukaryotic, in which

dUTP 2nd Strand Method: Our Favorite



Modified from Parkhomchuk et al. (2009) Nucleic Acids Res. 37:e123

Slide courtesy of Joshua Levin, Broad Institute.

Overlapping UTRs from Opposite Strands



Schizosacharomyces pombe (fission yeast)



Antisense-dominated Transcription



Trinity output: A multi-fasta file

comp0 c0 seq1 len=5528 path=[1:0-3646 10775:3647-3775 3648:3776-5527]

AATTGAATCCCTTTTTGTATCGAAAAATTGAAAGTGAAAGACATATACAGATTGAATGCGGTGATGGAATATAAGAATTTGGAAGATTAGAAAATATACAAAATTGACGGAGCACACCTAGGTTCG TOCACTOCCATCATOTOGAGATACTACAGAGGACTATCCGTCCACAGGACGTAACTGAACCCGATTCCTCCTTTCTTGCAAAGTCTTGACTTGACTAGGATCTCAGTAGAAAAAGCAGCAGCATTCTTTTTTCAGTCT TTOTGAATCCCAGACAGTTACGATAAGAATGCAATGGTGTGCCCCGGAGCAGCCATGGGAAGACCAGTCCTCACCAAGTCATCTTTCACCTTACAGTTACCCTCAGGAATAAAGCGACGGGAACAAGAACAGA GTGAACAACATGAACACCCTGATGCAGCAGTCTTAAGTGTCAACAGGACACCAACATCAGGECCATTATAAAACATACCTTTTTCAACCTAAAAACCTAGGTTAAAACCCATTTAAACCCTAGGTTAAACCTAGGTTAAACCTAGGTTAAACCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTAAACCCTAGGTTGGTTGGTTAGGTTAGGTTAGGTTAGGTTAGGTTAGGTTAGGTTAGGTTAGGTTAGGTTAGGTTGGTTAGGTTAGGTTAGGTTAGGTTAGGTTGGTTAG TCACAGTAACTGGACACCCAAAGGACAGAAATAGTCTCAACGAAGAAGACGAGGACTACCAGGGCTGGGGTCTTCACATTGCCATCTGTAAGAGGTCCCCCTTTACATGTCCCGAAGAACACCTCT GCTTCTCCCATACATCAATGAGCACATGAACAGCGAGCAGCAGCAGCAGTAATAGTCTGAGAACTGCAACTCTGTCTTCAAACAACAAGAGCGCCCCAAACCCGTGCTGGTGCTGGTACCTTCAGCACACTCTTTGACCACATCCAG

>comp0_c0_seq2 len=5399 path=[1:0-3646 3648:3647-5398]

ARTTGRATCCCTTTTTGTATCGRARASCTGRARGCATATACAGATGGATGGATGGATGGGATGGAAATATAATGCARATTAGAAAATTATGAAAATTGATGGAGGACGACGACGACGCCCCGGGTGTGG ASTTATCTCAARATGTAAGAATTAGACATTGAAAATGCACATTAGAAAATCAGCAAGTAACAAGAAGTAAACAAGCACATGAACAACAACAACAACAAGACCAGGCGCCCCACATGCAAGAACAAGACA TTOTGAATCCEAGACAGTTACGATAAAGAATGCAATGGTGTGCTGCGGGCAGTGGGAAGACCAGTCCTCACCAGTCTTTCACCTTACAGTTACCAGTACAGGAATAAAGTGGCGGCGCGGGGAACAAGAACAGA GTARACCCRGRTGRGGGTCCTGCTGCTGCTGCTGTTATATACAATTGCTGTATATTTGATACCCCCRARAATTGATTCACGATCCATGCATGCATGCATGCATGCAGGAAGTTCCGGATTAGAACAATGCCAGC ASCOCTCCAGAATCATGTAATAAAGTTCAACCTCAGCCTCCACCATCTTCTCCCACCATCTTCTCCGCCAGGGGCAGAAACATGGTTTTGGAGAGCCTCCACCGGGCATATAGAT TAAATGGGCCGGAGGGGCGGTCGTTAGGGTCCTGCACATGGCCCGGGGTCGCCATGATGACAAGCGCAGAACCTCAGT

nature protocols

NATURE PROTOCOLS | PROTOCOL

De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis

Brian J Haas, Alexie Papanicolaou, Moran Yassour, Manfred Grabherr, Philip D Blood, Joshua Bowden, Matthew Brian Couger, David Eccles, Bo Li, Matthias Lieber, Matthew D MacManes, Michael Ott, Joshua Orvis, Nathalie Pochet, Francesco Strozzi, Nathan Weeks, Rick Westerman, Thomas William, Colin N Dewey, Robert Henschel, Richard D LeDuc, Nir Friedman & Aviv Regev

Affiliations | Contributions | Corresponding authors

Nature Protocols **8**, 1494–1512 (2013) | doi:10.1038/nprot.2013.084 Published online 11 July 2013



Evaluating the quality of your <u>transcriptome</u> assembly

- Read representation by assembly
- Full-length transcript reconstruction
- Contig N50 leveraging expression data (ExN50)
- Detonate



Evaluating the quality of your transcriptome assembly

Read representation by assembly

Align reads to the assembled transcripts using Bowtie. A typical 'good' assembly has ~80 % reads mapping to the assembly and ~80% are properly paired.

Given read pair: –

→ ←

Possible mapping contexts in the Trinity assembly are reported:



Assembled transcript contig is only as good as its read support.

% samtools tview alignments.bam target.fasta

911	921	931	941	951	961	971	981	991	1001	1011	1021	1031	1041	1051	1061	1071	
GTAGGTT	TAATTTCATC	TTCTAATTTAGAA	TCTTGCCAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	CTTAGATGO	CCAAGTACATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAAC	TCCTCCATT	CAAGACTTAATTGACTCT	GT
GT GTT	TAATTTCATC	TTCTAATTTAGAA	TCTTGCCAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	AAC	ctgcttctgagatt	tctaagtac	cttagatgo	ccaagtacatt	actataatt	ggtgttatcg	ggtcttcc c	tcctccatt	caagacttaattgactct	gt
GT	ATTTCATC	TTCTAATTTAGAA	TCTTGCCAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAAC	tgcttctgagatt	tctaagtac	cttagatgo	ccaagtacatt	actataatt	ggtgttatcg	ggtcttcca	cctccatt	caagacttaattgactct	gt
GT	atttcato	ttctaatttagaa	tcttgccaa	tcaagccctc	tcgaagttggc	aatatctata	actcaac	GCTTCTGAGAT	TCTAAGTAC	CTTAGATGO	CCAAGTACATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAA	cctccatt	caagacttaattgactct	gt
GT	atttcato	ttctaatttagaa	tcttgccaa	tcaagccctc	tcgaagttggc	aatatctata	actcaac	GCTTCTGAGAT	TCTAAGTAC	CTTAGATGO	CCAAGTACATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAA	cctccatt	caagacttaattgactct	gt
GTAGGTT	TAAT	aa	tcttgccaa	tcaagccctc	tcgaagttggc	aatatctata	actcaacc	tctgcttctgagatt	tcta	CTTAGATGO	CCAAGTACATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAAC	TCCTCCATT	CAAGACTTAA ct	lgt
GTAGGTT	TAATTT		tcttgccaa	tcaagccctc	tcgaagttggc	aatatctata	actcaacc	tctgcttctgagatt	tctaag	CTTAGATGO	CCAAGTACATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAAC	TCCTCCATT	CAAGACTTAA	
GTAGGTT	TAATTTCATC	Π	cttgccaa	tcaagccctc	tcgaagttggc	aatatctata	actcaacc	tctgcttctgagatt	tctaagt	TTAGATGO	CCAAGTACATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAAC	TCCTCCATT	CAAGACTTAAT	
GTAGGTT	TAATTTCATC	TTC	TGCCAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	ATGO	CCAAGTACATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAAC	TCCTCCATT	CAAGACTTAATTGAC	
GTAGGTT	TAATTTCATC	TTCTAAT	TGCCAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	GC GC	CCAAGTACATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAAC	TCCTCCATT	CAAGACTTAATTGACTC	
gtaggtt	taatttcatc	ttctaatttag	TGCCAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC		CATT	ACTATAAT	GGTGTTATCG	GGTCTTCCAAC	TCCTCCATT	CAAGACTTAATTGACTCT	GT
GTAGGTT	TAATTTCATC	TTCTAATTTAG	GCCAA	TCAAGCCTTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	C	catt	actataatt	ggtgttatcg	ggtcttccaac	tcctccatt	caagacttaattgactct	igt
GTAGGTT	TAATTTCATC	TTCTAATTTAG	CAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	C			tgttatcg	ggtcttccaac	tcctccatt	caagacttaattgactct	igt
GTAGGTT	TAATTTCATC	TTCTAATTTAG	CAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	CTT			g	ggtcttccaac	tcctccatt	caagacttaattgactct	igt
GTAGGTT	TAATTTCATC	TTCTAATTTAG		gccctc	tcgaagttggc	aatatctata	actcaacc	tctgcttctgagatt	tctaagtac	cttagatgo	CC		(GGTCTTCCAAC	TCCTCCATT	CAAGACTTAATTGACTCT	GT
GTAGGTT	TAATTTCATC	TTCTAATTTAGAA	T	CCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	CTTAGATGO	CCA		9	ggtcttccaac	tcctccatt	caagacttaattgactct	gt
GTAGGTT	TAATTTCATC	TTCTAATTTAGAA	тст	ctc	tcgaagttggc	aatatctata	actcaacc	tctgcttctgagatt	tctaagtac	cttagatgo	ccaag		9	ggtcttccaac	tcctccatt	caagacttaattgactct	gt
GTAGGTT	TAATTTCATC	TTCTAATTTAGAA	тст	C	TCGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	CTTAGATGO	CCAAGTA			GTCTTCCAAC	TCCTCCATT	CAAGACTTAATTGACTCT	GT
GTAGGTT	TAATTTCATC	TTCTAATTTAGAA	тст		CGAAGTTGGC	AATATCTAT	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	CTTAGATGO	CCAAGTACA			gtcttccaac	tcctccatt	caagacttaattgactct	gt
GTAGGTT	TAATTTCATC	TTCTAATTTAGAA	тст		AAGTTGGC	AATATCTAT/	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	CTTAGATG	CCAAGTACATT			cttccaac	tcctccatt	caagacttaattgactct	gt
gtaggtt	taatttcatc	ttctaatttagaa	tcttgcc		C	AATATCTAT	ACTCAACO	TCTGCTTCTGAGAT	TCTAAGTAC	CTTAGATGO	CCAAGTACATT	ACTATAA		cttccaac	tcctccatt	caagacttaattgactct	gt
GTAGGTT	TAATTTCATC	TTCTAATTTAGAA	TCTTGCCA			CTAT/	ACTCAACC	TCTGCTTCTGAGAT	TCTAAGTAC	CTTAGATGO	CCAAGTACATT	ACTATAAT	GGTG	CTTCCAAC	TCCTCCATT	CAAGACTTAATTGACTC	GT
GTAGGTT	TAATTTCATC	TTCTAATTTAGAA	TCTTGCCAA					cttctgagatt	tctaagtac	cttagatgo	ccaagtacatt	actataatt	ggtgttatcg	ggtcttccaac	CTCCATT	CAAGACTTAATTGACTC	G
gtaggtt	taatttcatc	ttctaatttagaa	tcttgccaa	tcaagcc				cttctgagatt	tctaagtac	cttagatgo	ccaagtacatt	actataati	ggtgttatcg	ggtcttccaac	tccatt	caagacttaattgactct	gt
GTAGGTT	TAATTICATC	TICTAATITAGAA	TCTTGCCAA	TCAAGCC				cttctgagatt	tctaagtac	cttagatgo	ccaagtacatt	actataati	ggtgttatcg	ggtcttccaac	tccatt	caagacttaattgactct	gt
gtaggtt	taatttcato	ttctaatttagaa	tcttgccaa	tcaagccc				ttctgagatt	tctaagtac	cttagatgo	ccaagtacatt	actataatt	ggtgttatcg	ggtcttccaac	t tccatt	caagacttaattgactct	gt
GTAGGTT	TAATTICATC	TICTAATTTAGAA	TCTTGCCAA	TCAAGCCC				tgagati	сстаадтас	cttagatgo	ccaagtacatt	астатаат	ggtgttatcg	ggtcttccaac	tcc ccatt	caagacttaattgactct	gŢ
GTAGGTT	TAATTICATC	TICTAATTTAGAA	TCTTGCCAA	TCAAGCCCTC	T00110			tgagati	tctaagtcc	cttagatgo	ccaagtacatt	астатаат	ggtgttatcg	ggtcttccaac	tcct catt	caagacttaattgactct	gŢ
GTAGGTT	TAATTTCATC		TCTTGCCAA	TCAAGCCCTC	TCGAAG			tgagati	tctaagtac	cttagatgo	ccaagtacatt	actataati	ggtgttatcg	ggtcttccaac		caagacttaattgactct	91
GTAGGTT	ATTTCATC		TCTTGCCAA	TCAAGCCCTC	TCGAAG	ATATOTAT	ACTOAAC	gagati	tetaagtae	cttagatgo	ccaagtacatt	actataati	ggigilateg	ggtcttccaac	teete	AAGACTTAATTGACTC	G
	ATTCATC		TCTTGCCAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT		T ACAT	CTAAGTAC	CTTAGATO	ccaagtacatt	actataat	ggigilateg	ggtcttccaac		cttaattgactct	9
	TICATC	TICTAATTTAGAA	TCTTGCCAA	TCAAGCCCTC	TCGAAGTTGGC	AATATCTAT	ACTCAACC	AGAT	tetaagtac	CTTAGATG	CAAGTACATT	ACTATAAT	GGIGITATCG	GGICTICCAAC	teeteen	attgacter	gu
								yati	tetaagtac	cttagaty	coortocott	actataatt	gytyttatty	gglcllccaac			
								yac	tetaattae	cttagaty	coortocott	actataatt	gytyttatty	ggicilicidac			
								yacı		cttagaty	coortocott	actataatt	ggtgttattg	ggicilicidac	teeteeatt	C220	
									aaytat	.cctayaty	ladylatati	αιιαιααι	ggrgrrarcg	gylettecaac	teeteeatt	caagacttaattgactc	at
														TTCCAAC			21
														TCCAAC	TCCTCCATT	CAAGACTTAATTGACTC	Ğ
														Caac	tectecatt	caagacttaattgactc	at a
														caac	tectecatt	caagacttaattgactc	9
														aac	tcctccatt	caagacttaattgactc	ot e
														aac	tectecatt	caagacttaattgactc	9.
														000			a

ccattcaagacttaattgactctg

ccattcaagacttaattgactctg

www.broadinstitute.org/igv/ C

☆ a



Overview

Can Examine Transcript Read Support Using IGV

🖲 🖸 🔵	Тгэг	Ka Ragions Tools CanomeSpace Hain	
Trinity.fasta	114		
L		254 bp	•
		p 100 bp 200 bp	-
GSNO_SRR1582647.bowtie.csoi am Coverage GSNO_SRR1582647.bowtie.csoi am	rt		
GSNO_SRR1582646.bowtie.csoi am Coverage GSNO_SRR1582646.bowtie.csoi am	rt		
GSNO_SRR1582648.bowtie.csoi am Coverage GSNO_SRR1582648.bowtie.csoi am	rt		
wt_SRR1582649.bowtie.csortec Coverage wt_SRR1582649.bowtie.csortec	d.		
wt_SRR1582650.bowtie.csortec Coverage wt_SRR1582650.bowtie.csortec	d.		
wt_SRR1582651.bowtie.csortec Coverage wt_SRR1582651.bowtie.csortec	d.		
Sequence → 3 tracks loaded TR	NITY	(_DN130_c0_g1_i]	

Can align Trinity transcripts to genome scaffolds to examine intron/exon structures

(Trinity transcripts aligned to the genome using GMAP)



Evaluating the quality of your transcriptome assembly

Full-length Transcript Detection via BLASTX



* Mouse transcriptome

Haas et al. Nat. Protoc. 2013

The Contig N50 statistic

"At least half of assembled bases are in contigs that are at least **N50** bases in length"

In genome assemblies – used often to judge 'which assembly is better'

In transcriptome assemblies – N50 is *not* very useful.

- Overzealous isoform annotation for long transcripts drives higher N50
- Very sensitive reconstruction for short lowly expressed transcripts drives lower N50

Often, most assembled transcripts are *very* lowly expressed (How many 'transcripts & genes' are there really?)



* Salamander transcriptome

Compute N50 Based on the Top-most Highly Expressed Transcripts (ExN50)

- Sort contigs by expression value, descendingly.
- Compute N50 given minimum % total expression data thresholds => ExN50



ExN50 Profiles for Different Trinity Assemblies Using Different Read Depths



Note shift in ExN50 profiles as you assemble more and more reads.

* Candida transcriptome

Detonate Software

"RSEM-EVAL [sic] uses a novel probabilistic model-based method to compute the joint probability of both an assembly and the RNA-Seq data as an evaluation score."



Li et al. Evaluation of de novo transcriptome assemblies from RNA-Seq data, Genome Biology 2014

Abundance Estimation (Aka. Computing Expression Values)

Calculating expression of genes and transcripts





Slide courtesy of Cole Trapnell
Calculating expression of genes and transcripts



Slide courtesy of Cole Trapnell

Normalized Expression Values

 Transcript-mapped read counts are normalized for both length of the transcript and total depth of sequencing.

Reported as: Number of RNA-Seq Fragments
 Per Kilobase of transcript
 per total Million fragments mapped
 FPKM

RPKM (reads per kb per M) used with Single-end RNA-Seq reads FPKM used with Paired-end RNA-Seq reads.

Transcripts per Million (TPM)

$$TPM_{i} = \frac{FPKM_{i}}{\sum_{j} FPKM} *1e6$$

Preferred metric for measuring expression

- Better reflects transcript concentration in the sample.
- Nicely sums to 1 million

Linear relationship between TPM and FPKM values.

TPM

Both are valid metrics, but best to be consistent.

FPKM

Multiply-mapped Reads Confound Abundance Estimation



Blue = multiply-mapped reads Red, Yellow = uniquely-mapped reads

Multiply-mapped Reads Confound Abundance Estimation



Isoform **B**

Blue = multiply-mapped reads Red, Yellow = uniquely-mapped reads

New fast alignment-free methods now available! eg. Kallisto

Use Expectation Maximization (EM) to find the most likely assignment of reads to transcripts.

Performed by:

- Cufflinks and Cuffdiff (Tuxedo)
- RSEM
- eXpress

Comparing RNA-Seq Samples

Some Cross-sample Normalization May Be Required

Why cross-sample normalization is important





Figure 1 Normalization is required for RNA-seq data. Data from [6] comparing log ratios of **(a)** technical replicates and **(b)** liver versus kidney expression levels, after adjusting for the total number of reads in each sample. The green line shows the smoothed distribution of log-fold-changes of the housekeeping genes. **(c)** An M versus A plot comparing liver and kidney shows a clear offset from zero. Green points indicate 545 housekeeping genes, while the green line signifies the median log-ratio of the housekeeping genes. The red line shows the estimated TMM normalization factor. The smear of orange points highlights the genes that were observed in only one of the liver or kidney largely attributable for the overall bias in log-fold-changes.

Robinson and Oshlack, Genome Biology, 2010

Normalization methods for Illumina high-throughput RNA sequencing data analysis.



From "A comprehensive evaluation of normalization methods for Illumina high throughput RNA sequencing data analysis" Brief Bioinform. 2013 Nov;14(6):671-83 <u>http://www.ncbi.nlm.nih.gov/pubmed/22988256</u>

Differential Expression Analysis Using RNA-Seq

Diff. Expression Analysis Involves

- Counting reads
- Statistical significance testing

	Sample_A	Sample_B	Fold_Change	Significant?
Gene A	1	2	2-fold	Νο
Gene B	100	200	2-fold	Yes

Observed RNA-Seq Counts Result from Random Sampling of the Population of Reads

Technical variation in RNA-Seq counts per feature is well modeled by the Poisson distribution



See: http://en.wikipedia.org/wiki/Poisson_distribution

Example: One gene*not* differentially expressed

SampleA(gene) = SampleB(gene) = 4 reads



Beware of concluding fold change from small numbers of counts

Poisson distributions for counts based on **2-fold** expression differences



From: http://gkno2.tumblr.com/post/24629975632/thinking-about-rna-seq-experimental-design-for

More Counts = More Statistical Power

Example: 5000 total reads per sample. Observed 2-fold differences in read counts.

	SampleA	Sample B	Fisher's Exact Test (P-value)
geneA	1	2	1.00
geneB	10	20	0.098
geneC	100	200	< 0.001

Tools for DE analysis with RNA-Seq





edgeR	ROTS
ShrinkSeq	TSPM
DESeq	DESeq2
baySeq	EBSeq
Vsf	NBPSec
Limma/Voom	SAMsec
mmdiff	NoiSeq
cuffdiff	

(italicized not in R/Bioconductor but stand-alone)

See: http://www.biomedcentral.com/1471-2105/14/91

Visualization of DE results and Expression Profiling

Plotting Pairwise Differential Expression Data



Significantly differently expressed transcripts have FDR <= 0.001 (shown in red)

Comparing Multiple Samples



Heatmaps provide an effective tool for navigating differential expression across multiple samples.

Clustering can be performed across both axes: -cluster transcripts with similar expression patters.

-cluster samples according to similar expression values among transcripts.

Examining Patterns of Expression Across Samples

Can extract clusters of transcripts and examine them separately.



Functional Annotation of Transcripts



RNA-Seq
Trinity
Transcripts/Proteins
Functional Data
Discovery

Automated Higher Order Biological Analysis

http://trinotate.sf.net

Trinotate Web for Interactive Analysis

Blast Hits, Pfam Domains, etc. **TrinotateWeb Entry Point** Trinotate Web for Annotation and Expression Analysis Pote per pege: [Starting plan w.PREn [] of 55 Center expression values: __dverage __median __more Server Transcript Annotations (Gene: comp3142_c0, Transcript: comp3142_c0_seq2) Stats land, Insura/diffEran (PO.001, C3 -matrix Rell, State sharters, fixed, P., 2 Various summary stats go here Got 8694 genes and 9299 transcripts Searcl hs_rep1 Text search of transcript annotation: log_rep1 Pfam for m.2492 ds_rep1 PD1398.16 1AB1/Mov34/MPN/PAD-1 ubioutin on Still needed: search based on specific attribute: pfam, go, kegg, etc BLAST for m 2491 pression Comparisons (Volcano and MA pl plat_rep1 AMSH3 ARATHIPerID:40.321E:2e-49 Re Sp_ds Sp_hs Sp_log hsive Spinla GF30(STALP_HUMAN(PerID: 33.41)E:2e-48 log vs. Sp_plat SR558ISTALP_PONABIPerID:33.411E:5e-48 Re Multi-sample Comparisons (Expression Profiling) Go to the interactive heatmap for all DE transcripts. Analyses of clusters of expression profiles edgeR_trans/diffExpr.P0.001_C2_matrix_R_all_RData_dusters_fixed_P_20 with 55 clusters 0 7 24 caraciti42_cd_seail Very Early Release and nig villeart • Ma Mgniffcart • No • Tea Just Scratching the Surface pectate sequence Heatmaps Volcano Plots MA-Plots Individual Transcript

Clustered Expression Profiles

Transcript/Protein Annotation Report

Expression Profiles

Transcript and **Protein Sequences**

Deciphering the Cell Circuitry of Limb Regeneration Via Single Cell Transcriptome Studies





Axolotl (Ambystoma mexicanum) Transcriptomics

Axolotl "water monster", aka Mexican salamander or Mexican walking fish.

- Model for vertebrate studies of tissue regeneration
- Short generation time
- Can fully regenerate a severed limb in just weeks.
- Genome estimated at ~30 Gb (not yet sequenced)



Key morphological steps during limb regeneration







Jessica Whited, Mark Mannucci, Ari Haberberg

1. Building a reference Axolotl transcriptome



1.3 billion of 100 bp paired-end Illumina reads





limb tissues and select other tissues with biological replicates

Framework for De novo Transcriptome Assembly and Analysis





Axolotl Transcriptome De novo Assembly Statistics And Quality Assessment

In silico Normalization



Counts of Transcripts				
Trinity contigs (transcripts)	1,701,035			
Trinity components (genes)	1,327,843			

Min. length 200 bases



ExN50 looks good!

Percent of Non-normalized Fragments Mapping as Properly Paired to Transcriptome



Biological Replicates Cluster According to Sample



Pearson Correlation Matrix for Tissue Replicates

2. Identification of Tissue-enriched Expression



EdgeR, min 4-fold change, FDR <= 1e-3

Identification of Tissue-enriched Gene Expression



Most Highly Expressed Blastema-enriched Genes



Functional Characterization of Blastema-enriched KAZD1



Viral-based Delivered Over-expression of KAZD1 Leads to Regeneration Defects



Summary of Key Points

- RNA-Seq is a versatile method for transcriptome analysis enabling quantification and novel transcript discovery.
- Expression quantification is based on sampling and counting reads derived from transcripts
- Fold changes based on few read counts lack statistical significance.
- Trinity assembly and supported downstream computational analysis tools facilitate transcriptome studies.
Acknowledgements



Aviv Regev Brian Haas Timothy Tickle Asma Bankapur



Jill Mesirov James Robinson



BRIGHAM AND WOMEN'S HOSPITAL

Nathalie Pochet



Salamander limb regeneration

Jessica Whited Tia DiTommaso Tae Lee Anna Guzikowski



Thomas Doak Carrie Ganote Robert Henschel Ben Fulton

Trinotate & TrinoateWeb

Brian Couger Leonardo Gonzalez

