De novo transcriptome assembly using Trinity

Robert Bukowski, Qi Sun Bioinformatics Facility Institute of Biotechnology

Slides: <u>http://cbsu.tc.cornell.edu/lab/doc/Trinity_workshop_Part1.pdf</u>

Exercise instructions: <u>http://cbsu.tc.cornell.edu/lab/doc/Trinity_exercise1.pdf</u>

Strategies for transcriptome assembly from RNA-Seq data



Trinity: state of the art de novo RNA-Seq assembly and analysis package





Variety of post-assembly tools

Assembly QC* Abundance estimation* Differential Expression analysis* Protein coding regions identification Functional annotation* Full-length transcript analysis*

(*next session)

(today's topic)

Developed at Broad Institute and Hebrew University of Jerusalem

N. G Grabherr et. al., Nature Biotechnology 29, 644–652 (2011) doi:10.1038/nbt.1883 B. J. Haas et. al., Nature Protocols 8, 1494–1512 (2013) doi:10.1038/nprot.2013.084

As other short read assemblers, Trinity works in K-mer space (K=25)

What are K-mers?



What an assembler is trying to do



- Short or long K-mers?
- No positional info on repeats (less important for transcriptomes)

Three stages of Trinity

0. Jellyfish

• Extracts and counts K-mers (K=25) from reads

1. Inchworm:

 Assembles initial contigs by "greedily" extending sequences with most abundant K-mers

2. Chrysalis:

 Clusters overlapping Inchworm contigs, builds de Bruijn graphs for each cluster, partitions reads between clusters

3. Butterfly:

 resolves alternatively spliced and paralogous transcripts independently for each cluster (in parallel)



From: N. G Grabherr et. al., Nature Biotechnology 29, 644–652 (2011) doi:10.1038/nbt.1883

Trinity programs

Trinity proper

- **Trinity** (perl script to glue it all together)
- Inchworm
- Chrysalis
- **Butterfly** (Java code needs Java 1.7)
- various **utility and analysis scripts** (in perl)

Bundled third-party software

- **Trimmomatic**: clean up reads by trimming and removing adapter remnants (Bolger, A. M., Lohse, M., & Usadel, B)
- Jellyfish: k-mer counting software
- Fastool: fasta and fastq format reading and conversion (Francesco Strozzi)
- ParaFly: parallel driver (Broad Institute)
- **Siclust**: a utility that performs single-linkage clustering with the option of applying a Jaccard similarity coefficient to break weakly bound clusters into distinct clusters (Brian Haas)
- **Collectl** : system performance monitoring (Peter Seger)
- Post-assembly analysis helper scripts (in perl)

External software Trinity depends on (needs to be in the search PATH):

- samtools
- Bowtie2
- **RSEM, eXpress**: alignment-based abundance estimation (Bo Li and Colin Dewey)
- kallisto, salmon: alignment-free abundance estimation
- Transcoder: identify candidate coding regions in within transcripts (Brian Haas Broad, Alexie Papanicolaou CSIRO)

Notation convention used on the following slides

Trinity commands will be abbreviated using the variable **TRINITY_DIR** to denote the location of the Trinity package

TRINITY_DIR=/programs/trinityrnaseq-Trinity-v2.4.0

(this is the latest version installed on BioHPC Lab)

Thus, a command

```
$TRINITY_DIR/Trinity [options]
```

really means

/programs/trinityrnaseq-Trinity-v2.4.0/Trinity [options]

The same notation is used in auxiliary shell scripts containing Trinity commands, used in the Exercise

Pipeline complicated, but easy to run

(in principle, just a single command)



Input: RNA-Seq reads

FASTA format: 2 lines for each read (">name", sequence)

>61DFRAAXX100204:1:100:10494:3070 ACTGCATCCTGGAAAGAATCAATGGTGGCCGGAAAGTGTTTTTCAAATACAAGAGTGACAATGTGCCCTGTTGTTT

FASTQ format: 4 lines per read ("@name", sequence, "+", quality string)

ASCII code of a letter in quality string - 33 equals phred quality score of the corresponding base.

For example, "B" stands for: 66 - 33 = 33, i.e., probability of the base (here: G) being miscalled is $10^{-3.3}$.

Base qualities are ignored by the assembler.

Input: RNA-Seq reads

Paired-end case: we have two "parallel" files – one for "left" and another for "right" end of the fragment:

First sequence in "left" file

First sequence in "right" file

Note1 : Read naming convention: /1 and /2 in front of the same name core indicate read pairing

Note2 : Cassava 1.8 names will be handled properly (no re-formatting needed), e.g.,

@HWI-ST896:156:D0JFYACXX:5:1101:1652:2132 1:N:0:GATCAG

will be treated as

```
HWI-ST896:156:D0JFYACXX:5:1101:1652:2132/1
```

Trinity command line

Suppose we have paired-end reads from one individual and 4 conditions in FASTQ files

Sp_ds.left.fq.gz, Sp_hs.left.fq.gz, Sp_log.left.fq.gz, Sp_plat.left.fq.gz, Sp_ds.right.fq.gz, Sp_hs.right.fq.gz, Sp_log.right.fq.gz, Sp_plat.right.fq.gz

(In this example, the FASTQ files have been compressed with gzip; uncompressed files can also be used.)

A Trinity command may be long and tedious to type. It is convenient to create (using a text editor, like **nano** or **vi**) a **bash script**, like this (note: "\" characters break one long line into shorter pieces):

```
#!/bin/bash
$TRINITY_DIR/Trinity --seqType fq \
--left Sp_ds.left.fq.gz,Sp_hs.left.fq.gz,Sp_log.left.fq.gz,Sp_plat.left.fq.gz \
--right Sp_ds.right.fq.gz,Sp_hs.right.fq.gz,Sp_log.right.fq.gz,Sp_plat.right.fq.gz \
--SS_lib_type RF \
--max_memory 2G \
--CPU 2 \
--output /workdir/bukowski/trinity_out
```

Save the script (e.g., as my_trinity_script.sh) and run it, redirecting any screen output to a file on disk:

nohup ./my trinity script.sh >& my trinity script.log &

Basic Trinity command dissected...

```
#!/bin/bash
$TRINITY_DIR/Trinity --seqType fq \
--left Sp_ds.left.fq.gz,Sp_hs.left.fq.gz,Sp_log.left.fq.gz,Sp_plat.left.fq.gz \
--right Sp_ds.right.fq.gz,Sp_hs.right.fq.gz,Sp_log.right.fq.gz,Sp_plat.right.fq.gz \
--SS_lib_type RF \
--max_memory 2G \
--CPU 2 \
--output /workdir/bukowski/trinity_out
```

NOTES:

- Use all reads from an individual (all conditions) to capture most genes
- Read files may be gzipped (as in this example) or not (then they should not have the ".gz" ending)
- Paired-end reads specified with --left and --right. If only single-end, use --single instead.
- 2G is the maximum memory to be used at any stage which allows memory limitation (jellyfish, sorting, etc.)
- At most 2 CPU cores will be used in any stage
- Final output and intermediate files will be written in directory /workdir/bukowski/my_trinity_out
- --SS_lib_type RF: The PE fragments are <u>strand-specific</u>, with left end on the Reverse strand and the right end on Forward strand of the sequenced mRNA template
 - For non-strand specific reads, just skip the option --SS_lib_type

Strand-specific RNA-Seq



Strand specificity slightly increases computation time (more K-mers), but is very helpful for de novo assembly

• Helps disambiguate between overlapping genes on opposite strands of DNA, sense and nonsense transcripts

Most RNA-Seq protocols now in use are strand specific

Basic Trinity command with single-end reads



Basic Trinity command with mixture of PE and SE reads

```
#!/bin/bash
$TRINITY_DIR/Trinity --seqType fq \
--left Sp_ds.left.fq.gz,Sp_ds.leftUnpaired.fq.gz \
--right Sp_ds.right.fq.gz,Sp_ds.rightUnpaired.fq.gz \
--SS_lib_type RF \
--max_memory 2G \
--CPU 2 \
--output /workdir/bukowski/my_trinity_out
```

Final output: assembled transcriptome

FASTA file called **Trinity.fasta** (located in the run output directory)

FASTA headers contain gene/isoform ID for each transcript, as determined by Trinity

>TRINITY_DN0_c0_g2_i1 len=995 path=[1183:0-52 1184:53-994] [-1, 1183, 1184, -2] >TRINITY_DN30_c0_g1_i1 len=2677 path=[2675:0-1921 2676:1922-2676] [-1, 2675, 2676, -2] >TRINITY_DN30_c0_g2_i1 len=1940 path=[2673:0-1921 2674:1922-1939] [-1, 2673, 2674, -2] >TRINITY_DN2_c0_g1_i1 len=386 path=[402:0-254 403:255-385] [-1, 402, 403, -2] >TRINITY_DN2_c0_g2_i1 len=279 path=[404:0-254 405:255-278] [-1, 404, 405, -2] >TRINITY_DN3_c0_g1_i1 len=1286 path=[2096:0-1271 2097:1272-1285] [-1, 2096, 2097, -2] >TRINITY_DN3_c0_g1_i2 len=2102 path=[2094:0-1271 2095:1272-2101] [-1, 2094, 2095, -2]



is a **gene** identifier is an **isoform** identifier

The rest of the header shows nodes of de Buijn graph traversed by the transcript.

Advanced (but important) Trinity options to consider

General issues with de novo transcriptome assembly of RNA-Seq data

RNA-Seq reads often need pre-processing before assembly

- remove barcodes and Illumina adapters from reads
- clip read ends of low base quality
- remove contamination with other species

Very non-uniform coverage

- highly vs lowly expressed genes
- high-coverage in some regions may imply more sequencing errors complicating assembly down the road
 - some "normalization" of read set needed

Gene-dense genomes pose a challenge (assembly may produce chimeric transcripts)

- overlapping genes on opposite strands (strand-specific RNA-Seq protocols may help)
- some overlap of genes on the same strand (harder to handle)

Pre-assembly read clean-up option

Trimmomatic: A flexible read trimming tool for Illumina NGS data (Bolger et al., http://www.usadellab.org/cms/?page=trimmomatic)

```
java -jar /programs/Trimmomatic-0.36/trimmomatic-0.36.jar PE -threads 2 -phred33 \
left.fq.gz right.fq.gz \
left.fq.gz.P.qtrim left.fq.gz.U.qtrim \
right.fq.gz.P.qtrim right.fq.gz.U.qtrim \
ILLUMINACLIP:TruSeq3-PE.fa:2:30:10 SLIDINGWINDOW:4:5 LEADING:5 TRAILING:5 MINLEN:25
```

Filtering operations (in order specified) performed on each read:

- Remove Illumina adapters (those in file **TruSeq3-PE.fa**) using "palindrome" algorithm
- Clip read when average base quality over a 4bp sliding window drops below 5
- Clip leading and trailing bases if base quality below 5
- Skip read if shorter than 25bp

Pre-assembly read clean-up: Trimmomatic

Output files (gzipped if requested file name ends with ".gz"):

- left.fq.gz.P.qtrim, right.fq.gz.P.qtrim: "parallel" files of reads such that both ends of a fragment survived filtering
- left.fq.gz.U.qtrim, right.fq.gz.U.qtrim: "left" and "right" reads that survived filtering, although the mate did not

Trimmomatic is bundled with Trinity

To invoke from within Trinity, add option

```
--trimmomatic
```

Also, one can request specific **Trimmomatic** settings by adding something similar to

```
--quality_trimming_params "ILLUMINACLIP:TruSeq3-PE.fa:2:30:10 SLIDINGWINDOW:6:5
LEADING:20 TRAILING:10 MINLEN:30"
```

Trinity will use **all surviving reads**, treating the un-paired ones as single-end.

How to choose trimming parameters: Read quality assessment with fastqc



Run fastqc before running Trinity (part of the Exercise)

Dealing with other species contamination: pre-assembly

Get rid of the reads from unwanted species. First, <u>check the raw reads for contamination</u>, **BLAST**ing the reads against <u>**nt**</u> database</u>.

As always, work on /workdir: copy the nt database files to local scratch directory

```
mkdir -p /workdir/myID/blast_db
cd /workdir/myID
cp /shared_data/genome_db/BLAST_NCBI/nt.* ./blast_db
cp /shared_data/genome_db/BLAST_NCBI/taxdb* ./blast_db
```

Then use our custom script to run BLAST agains nt using random 500 of your reads as a query, and then analyze the result, e.g.

/programs/fastq_species_detector/fastq_species_detector.sh reads.fq ./blast_db 1e-6

Example of what you may get:

Species	#Reads	%Reads	
Schizosaccharomyces pombe	478	97.154	contamination
Oryza sativa Japonica Group	12	2.439	
Schizosaccharomyces pombe 972h-	1	0.203	

Dealing with other species contamination: pre-assembly

If contamination detected, remove the unwanted reads by <u>aligning all to the transcriptome of the contaminant</u>, and taking only those reads that <u>do not align</u> to this transcriptome.

For example, a procedure based on **STAR aligner** works fine for this:

First, prepare and index the contaminant transcriptome

```
mkdir -p /workdir/myID/STARgenome
/programs/STAR/STAR --runMode genomeGenerate --runThreadN 2 \
--genomeDir STARgenome \
--genomeFastaFiles contam_genome.fa --sjdbGTFfile contam_annot.gff3 \
--sjdbGTFtagExonParentTranscript Parent --sjdbOverhang 49
```

The align the reads to it, saving the unmapped reads

```
/programs/STA/STAR --genomeDir STARgenome --runThreadN 2 \
--readFilesIn dirty.left.fq.gz dirty.right.fq.gz \
--readFilesCommand zcat --outFileNamePrefix oryza_ --outFilterMultimapNmax 1 \
--outReadsUnmapped Fastx --outSAMtype BAM SortedByCoordinate
```

Unmapped PE reads will be in files oryza_Unmapped.out.mate1 and oryza_Unmapped.out.mate2

USE THESE IN THE ASSMEBLY!

Dealing with other species contamination: post-assembly

- Assemble first, then compare contigs to existing databases
 - this is really a part of annotation procedure will be discussed in Part 2 of this workshop
- The only option if contaminant transcriptome not available

Dealing with very deep sequencing data

- Done to identify genes with low expression
- Several hundreds of millions of reads involved
- More sequencing errors possible with large depths, increasing the graph complexity

Suggested Trinity options/treatments:

• Use option

--min kmer cov 2

(singleton K-mers will not be included in initial Inchworm contigs)

- Perform K-mer based insilico read set normalization (now done by default)
 - May end up using just 20% of all reads reducing computational burden with no impact on assembly quality

K-mer based read normalization

Sample the reads (fragments) with probability

$$P = \min(1, \frac{T}{C})$$

Where *T* is the target K-mer coverage (k=25) and *C* is the median K-mer abundance along the read (or average over both fragment ends). Typically, T = 30 - 50.

(Also: filter out reads for which STDEV of K-mer coverage exceeds C)

Effect: poorly covered regions unchanged, but reads down-sampled in high-coverage regions.

Normalization is done by default with T=50. To change target coverage, add option

```
--normalize_max_read_cov 30
```

To skip read normalization step, add option **--no_normalize_reads**

To run read normalization separately, use the utility script:

\$TRINITY_DIR/util/insilico_read_normalization.pl

(run without arguments to see available options)

This normalization method has "mixed reviews"

Avoiding chimeric transcripts from gene-dense genomes

- Use strand-specific RNA-Seq protocol
- Try option -jaccard_clip (time-consuming; no effect in case of gene-sparse genomes, maybe just check with IGV after DE calculation for important genes)



Resource requirements

How many reads are needed?

How long will the assembly take?

How much RAM is needed?



а





b

mouse

S. Pombe

From: N. G Grabherr et. al., Nature Biotechnology 29, 644–652 (2011) doi:10.1038/nbt.1883

How long will it take?



From: B. J. Haas et. al., Nature Protocols 8, 1494–1512 (2013) doi:10.1038/nprot.2013.084

How much RAM and how many processors?

Memory and time needed for assembly strongly depend on data complexity (rather than amount of data)

<u>Rule of thumb:</u>	
Memory needed:	1GB of RAM per 1 million reads
Timing:	¹ ⁄ ₂ - 1 hours per 1 million reads (does not include trimming or normalization)

Other tips:

- Most (not all) parts of Trinity are parallelized makes sense to use most available CPUs (via - CPU option)
 - some programs may adjust it down (by default, Inchworm runs on at most 6 CPUs)
- Butterfly (last stage) benefits most from massive parallelization....
 -although running on too large a number of CPUs may lead to memory problems
 - controlled using -bflyCPU and -bflyHeapSpaceMax options
- Most runs should go through on BioHPC Lab <u>medium-memory machines (cbsumm*, 128 GB of memory ad 24 CPU cores</u>) with options

--CPU 20 --max_memory 100G

A couple of more "real" runs

```
Example 1:

Drosophila

Machine: cbsum1c1b016 (8 CPU, 16 GB RAM - small-memory)

Options: --CPU 8 (6 for Inchworm) --max_memory 12G --bflyCPU 6 --bflyHeapSpaceMax 2G

Example 2:

Drosophila

Machine: cbsumm02 (24 CPU, 128 GB RAM - medium-memory)

Options: --CPU 20 (6 for Inchworm) -max_memory 100G --bflyCPU 20 --bflyHeapSpaceMax 4G

Example 3:

Drosophila

Machine: cbsumm02 (24 CPU, 128 GB RAM - medium-memory)
```

Options: --CPU 20 (6 for Inchworm) -max_memory 100G --bflyCPU 20 --bflyHeapSpaceMax 4G -max_kmer_cov 2

	Initial	PE frags after		V	Vall-clock time	es [minutes]			
	PE Trags	normalztn	Normalization	Jellyfish	Inchworm	Chrysalis	Butterfly	Total (excluding normalization)	
Example 1	40M	8.6M	100	5	30	80	523	638	
Example 2	206M	19.6M	270	6	52	118	252	428	
Example 3	206M	Not norm	-	13	35	6,800	396	7,244	

NOTE: read normalization and trimming also take time!

Launching and monitoring a Trinity run

Launching a Trinity run

Assuming the input files and the script my_trinity_script.sh (previous slides) are all in /workdir/<yourID>, launch the programs using the commands

cd /workdir/<yourID>
nohup ./my_trinity_script.sh >& my_trinity_script.log &

NOTES:

- The program will be launched in the background (because of the "&" at the end)
- **nohup** ensures you can log out of the machine the program will still be running when you log back in
- screen output will be written to the file /workdir/<yourID>/my_trinityscript.log

Progress of the program can be checked in following ways:

- Monitor the screen output file my_trinityscript.log (will be located in /workdir/<yourID>)
- Monitor intermediate files written to the output directory /workdir/<yourID>/trinity_out
- Monitor running processes using top Linux utility

Look at the screen output file (in /workdir/<yourID>):

more my_trinityscript.log (will page through the file from the beginning)
tail -f my_trinityscript.log (will display new lines of the file as they appear)
nano my_trinityscript.log (will open the file, still incomplete, in text editor nano)

Monitor intermediate output files (in /workdir/<yourID>/trinity_out):

ls -al

	-					_		
1			[screen 2	2: bash] buk	cows	ki@	cbsum1	c1b016:/workdir/bukowski/trinity_out4
drwxr-x	4	bukowski	bukowski	4096	Mar	5	13:19	
drwxr-x	11	bukowski	bukowski	4096	Mar	5	13:14	
-rw-r	1	bukowski	bukowski	69337189	Mar	4	10:41	both.fa
-rw-r	1	bukowski	bukowski	7	Mar	4	10:41	both.fa.read count
drwxr-x	2	bukowski	bukowski	4096	Mar	4	10:43	chrysalis
-rw-r	1	bukowski	bukowski	1256285	Mar	4	10:41	inchworm.K25.L25.fa
-rw-r	1	bukowski	bukowski	0	Mar	4	10:41	inchworm.K25.L25.fa.finished
-rw-r	1	bukowski	bukowski	8	Mar	4	10:41	inchworm.kmer_count
-rw-r	1	bukowski	bukowski	0	Mar	4	10:41	jellyfish.1.finished
-rw-r	1	bukowski	bukowski	39526208	Mar	4	10:41	jellyfish.kmers.fa
-rw-r	1	bukowski	bukowski	17911	Mar	4	10:41	jellyfish.kmers.fa.histo
-rw-r	1	bukowski	bukowski	21553	Mar	4	10:43	partitioned_reads.files.list
-rw-r	1	bukowski	bukowski	0	Mar	4	10:43	partitioned_reads.files.list.ok
drwxr-x	3	bukowski	bukowski	4096	Mar	4	10:43	read_partitions
-rw-r	1	bukowski	bukowski	87926	Mar	4	10:43	recursive_trinity.cmds
-rw-r	1	bukowski	bukowski	87926	Mar	4	10:51	recursive_trinity.cmds.completed
-rw-r	1	bukowski	bukowski	0	Mar	4	10:43	recursive_trinity.cmds.ok
-rw-r	1	bukowski	bukowski	3564376	Mar	4	10:42	scaffolding_entries.sam
-rw-r	1	bukowski	bukowski	429405	Mar	4	10:41	<pre>tmp.iworm.fa.pid_15881.thread_0</pre>
-rw-r	1	bukowski	bukowski	418950	Mar	4	10:41	<pre>tmp.iworm.fa.pid_15881.thread_1</pre>
-rw-r	1	bukowski	bukowski	0	Mar	4	10:41	trimmomatic.ok
-rw-r	1	bukowski	bukowski	418919	Mar	4	10:51	Trinity.fasta
-rw-r	1	bukowski	bukowski	777	Mar	4	10:51	Trinity.timing
[bukowski@cbsum1c1b016			trinity_c	out4]\$				

- Files named after Trinity components
- Files will appear in succession (see modification times)
- Trinity.fasta final output file with transcripts; if present, program ended successfully
- ***.finished**, ***.ok** mark successful completion of a program stage

Trinity is restartable – just run <u>the</u> <u>same command</u>, possibly with some options updated.

Monitor processes using top -u bukowski (substitute your own userID)

K-mer counting stage (with Jellyfish)

Initial contiging stage (with Inchworm)

₽			[scre	een 0: b	oash] t	oukow	ski@cb	sum1c	1b016:/	workdir/	/bukowski/	/rnasec	₽			[scree	n 0: b	ash] b	ukow	ski@cbs	um1c1b	016:/work	dir/bukowsk	i/rnaseq_w
top -	14:19:36	up 37	/ da	ys, 22	2:23,	6 us	sers,	load	avera	re: 0.40), 0.13,	0.09	top -	14:19	:54 u	ip 37	day	3, 22	:23,	6 us	ers,	load a	verage: 0	.61, 0.19,	0.11
Tasks	Tasks: 225 total, 2 running, 223 sleeping, 0 stopped, 0 zombie													: 225	total	L, 1	2 ru	nning	, 223	slee	ping,	0 st	opped,	0 zombie	
Cpu (s)): 22.0%us	, 0.	3%s	у, О.	0%ni,	, 77.6	5%id,	0.1%w	7a, 0.	O%hi,	0.0%si,	0.0%	Cpu (s)): 25.	0%us,	0.1	7%sy	, o.	O%ni,	74.3	°%id,	0.1%wa	, 0.0%hi	, 0.0%si,	0.0%st
Mem:	16336144k	tota	ıl,	104907	/00 k 1	ised,	5845	444k f	ree,	234240	k buffer	3	Mem:	16336	144k	tota	1,	32720	68k u	used,	80640	76k fr	ee, 234	256k buffe	rs
Swap:	18579452k	tota	1,	1873	396 kr 1	ised,	18392	056k f	ree,	6719148	k cached	1	Swap:	18579	452 k	tota	1,	1873	96k u	ised,	183920	56k fr	ee, 6771	.548k cache	d
Π						-							Π							-			-		
PID	USER	PR	NI	VIRT	RES	SHR	S %CP	U %MEM	(T	ME+ CO	MMAND		PID	USER		PR 1	NI 1	/IRT	RES	SHR	S %CPU	%MEM	TIME+	COMMAND	
21742	bukowski	20	0	2537m	2.3g	1428	S 200	.0 14.	8 0:	21.23 j	ellyfish	1	21758	bukow	ski	20	0	203m	151m	1428	R 200.	0 1.0	0:24.1	9 inchworm	1
2959	bukowski	20	0	289m	13m	9660	s 0.	7 0.1	0:5	9.12 gr	nome-term	inal	2959	bukow	ski	20	0	289m	13m	9660	S 1.0	0.1	0:59.25	gnome-ter	minal
2462	bukowski	20	0	108m	44m	7884	s o.:	3 0.3	1:2	7.21 Xv	mc		2462	bukow	ski	20	0	108m	44m	7884	s 0.3	0.3	1:27.26	Xvnc	
			-	· · -							-		21747	bukow	ski	20	0 1	7208	1256	872	R 0.3	0.0	0:00.09	top	

Chrysalis stage

Transcript stage (Butterfly)

₽			[scre	en 0: b	oash] k	oukow	ski@cb	sum1c	1b016:/wor	kdir/bukows	ki/rnaseq	₽		[scre	en 0: b	ash] t	oukows	ki@cb	sum1	c1b016:/v	workd	dir/bukowski	/rnaseq_w
top - Tasks: Cpu(s)	14:20:37 u 225 total : 27.7%us,	up 3 , 0	7 da 2 r .0%s	ys, 22 unning y, 0.	2:24, g, 223 .0%ni,	6 us 3 slee , 72.1	ers, ping, %id,	load 0 ៖ 0.0%	average: stopped, va, 0.0%h	1.21, 0.41 0 zombie i, 0.1%si	, 0.19 , 0.0%:	top - Tasks Cpu(s	14:22:13 226 tota): 33.7%us	up 37 1, 2 , 4.0	day 2 ri 0%sy	ys, 22 unning y, 0.	:25, , 224 O%ni,	6 us 4 slee , 62.0	ers, ping, %id,	10ad 0 0.3%	average stopped, wa, 0.0	e: 2. , (O%hi,	.03, 0.88, D zombie , 0.0%si,	0.37 0.0%st
Mem:	16336144k	tot	al,	85780)92 kr 1	used,	7758	052k f	free, 23	4332k buff	ers	Mem:	16336144k	tota	1,	84255	12 k ı	ised,	7910	632 k	free,	2346	600k buffe	rs
Swap:	18579452k	tot	al,	1873	396kr 1	used,	18392	056k f	free, 680	8816k cach	ed	Swap:	18579452k	tota	1,	1873	96 kr 1	ised,	18392	056k	free, '	70275	504k cache	d
PID	USER	PR	NI	VIRT	RES	SHR	S %CP	U %MEN	1 TIME+	COMMAND		PID	USER	PR 1	NI	VIRT	RES	SHR	S %CP	U %ME	M TI	ME+	COMMAND	
21791	bukowski	20	0	499m	413m	1568	R 199	.5 2	6 0:13.	36 GraphFr	omFasta	24924	bukowski	20	0	1469m	27m	9.9m	S 11.	30.	2 0:0	0.34	java	
2959	bukowski	20	0	289m	13m	9660	s 0.	7 0.1	0:59.4	8 gnome-te	rminal	24941	bukowski	20	0	5240m	27m	10m	S 6.	60.	2 0:0	0.20	java	
3097	root	20	0	446m	14m	1280	s 0.	3 0.1	39:41.9	2 glusterf	3	2399	root	20	0	380m	47m	1352	S 4.	ο ο.	3 6:3	7.28	glusterfs	
												24763	bukowski	20	0	131m	8220	2056	s 3.	7 0.	1 0:00	0.11	perl	
												24862	bukowski	20	0	131m	8180	2048	S 3.	ο ο.	1 0:0	0.09	perl	

In case Butterfly crashes.....

If Trinity completes, but the final file output **Trinity**. **fasta** is not produced, scan the log file (screen output) for Butterfly errors:

Butterfly fails with java Error: Cannot create GC thread. Out of system resources

<u>Reason</u>: each Butterfly process (java VM) tries to allocate certain amount of heap space (by default: 4GB). With a large **--CPU** setting, this may exhaust all memory on a machine.

<u>Remedy</u>: restart the job with the same command as before, but with reduced **--CPU** setting. <u>It will pick up from where</u> it crashed (and hopefully run to completion).

Also, Butterfly CPU and memory options may be set independently (from --CPU, --max_memory):

```
--bflyCPU 20 --bflyHeapSpaceMax 4G
```

Assessing quality of the assembly

This session:

Check basic contig statistics

Check read representation of the assembly.

Compute the ExN50 profile and E90N50 value

Other methods (some will be discussed next week):

Examine the **representation of full-length reconstructed protein-coding genes**, by searching the assembled transcripts against a database of known protein sequences.

Use **BUSCO** (Benchmarking Universal Single-Copy Orthologs) to explore completeness according to conserved ortholog content (<u>http://busco.ezlab.org/</u>)

Compute **DETONATE scores** (**DE** novo **T**ranscript**O**me r**N**a-seq **A**ssembly with or without the **T**ruth **E**valuation). DETONATE provides a rigorous computational assessment of the quality of a transcriptome assembly (<u>http://deweylab.biostat.wisc.edu/detonate/</u>)

Assessing quality of the assembly: basic contig statistics

More about it next week. Today we introduce a few simple methods, included in Trinity package

\$TRINITY_DIR/util/TrinityStats.pl Trinity.fasta >& stats.log

Total trinity 'genes': 1388798 Total trinity transcripts: 1554055 Percent GC: 44.52

Contig N10: 5264 Contig N20: 3136 Contig N30: 1803 Contig N40: 989 Contig N50: 606

Example output

The script also provides stats computed using only the longest contig from each gene

Nx: x% of all bases are in contigs of length at least Nx

For transcriptome assembly, large N50 not necessarily the best!

Median contig length: 288 Average contig: 511.61 Total assembled bases: 795066996

Assessing quality of the assembly: Read content

Typical Trinity transcriptome assembly will have the **vast majority of all reads mapping back to the assembly**, and **~70-80% of the mapped fragments found mapped as proper pairs**. Here is how to check this:

Map RNA-Seq reads back to the assembly (using the **bowtie2** aligner)

```
# First, build the bowtie2 index of the assembly
bowtie2-build Trinity.fasta Trinity.fasta
# Align reads to the assembly
bowtie2 --local --no-unal -x Trinity.fasta -q -1 left_reads.fq -2 right_reads.fq \
    samtools view -Sb -@ 2 | samtools sort -@ 2 -m 1G -n -o bowtie2.nameSorted.bam
```

This will produce the alignment BAM file **bowtie2_out.nameSorted.bam**.

Now produce alignment statistics:

\$TRINITY_DIR/util/SAM_nameSorted_to_uniq_count_stats.pl bowtie2_out.nameSorted.bam

Assessing quality of the assembly: Read content

Example output

```
Stats for aligned rna-seq fragments (note, not counting those frags where
neither left/right read aligned)
328489 aligned fragments; of these:
  328489 were paired; of these:
    3848 aligned concordantly 0 times
    324641 aligned concordantly exactly 1 time
    0 aligned concordantly >1 times
    3848 pairs aligned concordantly 0 times; of these:
    1555 aligned as improper pairs
    2293 pairs had only one fragment end align to one or more contigs; of
these:
       1023 fragments had only the left /1 read aligned; of these:
            1023 left reads mapped uniquely
            0 left reads mapped >1 times
       1270 fragments had only the right /2 read aligned; of these:
            1270 right reads mapped uniquely
            0 right reads mapped >1 times
Overall, 98.83% of aligned fragments aligned as proper pairs
```

Assessing quality of the assembly: ExN50 profile



Curves correspond to different sequencing depths

Each point shows **N50** computed from top most highly expressed transcripts that represent **x%** of the total normalized expression data

Short, incomplete contigs representing lowexpressed transcripts make the curves drop on the rhs

The peak indicates which transcripts are reasonably complete (typically about 1% of all Trinity transcripts)

The peak shifts to the right with increasing sequencing depth (more low-expressed transcripts are assembled)

Assessing quality of the assembly: ExN50 profile

Computation of ExN50 profile requires expression quantification – the subject of next week's session

- Map reads to transcriptome assembly
- Count reads mapping to transcripts (one read can map to more than one transcript!)
- Evaluate expression measures
- Produce ExN50 profile

Auxiliary script provided as a part of the Exercise

For more information visit Trinity site:

https://github.com/trinityrnaseq/trinityrnaseq/wiki

For exercise-related questions contact

bukowski@cornell.edu