



### Exercise 1

#### Run shell script run\_tuxedo.sh (to be explained)

- 1. Copy ch4\_demo\_dataset.tar.gz from ~/Data to your Desktop and unzip
  - \$cp ~/Data/ch4\_demo\_dataset.tar.gz ~/Desktop
  - \$cd ~/Desktop
  - \$tar -xvf ch4\_demo\_dataset.tar.gz

\$chmod 755 run\_tuxedo.sh \$cp scripts/shell/run\_tuxedo.sh ch4\_demo\_dataset/

4. Change name of SL2.40ch04.fasta to SL2.40ch04.fa 5. run shell script from ~/Desktop/ch4\_demo\_dataset/

\$cd ch4\_demo\_dataset/

\$./run\_tuxedo.sh •

2. Copy shell script run\_tuxedo.sh to ~/Desktop/ch4\_demo\_dataset/ and make it executable

\$mv ch4\_demo\_dataset/bwt2\_index/SL2.40ch04.fasta ch4\_demo\_dataset/bwt2\_index/SL2.40ch04.fa

This will run ~15 minutes. We'll discuss what it is doing while it runs.

### Transcriptome Assembly First steps

- \* FASTQC give some basic stats about raw reads, can run again after cleaning
  - You can see bias of the random hexamer priming in per base sequence content (next slide)
  - Also GC bias can be present, can effect DE results
- Duplicate removal not recommended for RNA-seq
- To clean or not to clean? Cleaning only critical for de novo assembly, unless adapter contamination
  - \* can use fastq-mcf, trimmomatic, prinseq for this, be sure to run both mate files together if you have paired end data!







### Transcriptome Assembly De novo vs Reference-guided

<u>Reference-guided</u> - map to previously-assembled genome or transcriptome

- Pros: computationally easier
- transcripts
- <u>De novo</u> cluster reads into transcripts
  - Pros: no prior assemblies/annotations necessary
  - cons: more difficult, more memory needed, LOTS of transcripts

\* cons: must have good reference assembly to a related species, may miss some

### Transcriptome Assembly Reference-guided vs De novo

#### **Reference-guided**



- How many reads for DE?
  - ◆ 30 x coverage of CDS (Introduction to RNA-Seq)
- ♦ What length? SE or PE?
  - short (50 bp), SE ok for differential expression, multi-mapping more problematic
  - Index of the second seco

- index)

- bowtie2-build for tophat2)

TopHat Aligns RNA-Seq reads to the genome using Bowtie Discovers splice sites

Most use Burrows Wheeler transform (transform and compress into

◆ Tool must be able to deal with splice junctions in eukaryotes

Does tool accept paired or single-end reads (fastq format)?

map to reference fasta genome file (needs to be indexed with





Tophat:

- multi-mapping reads default tries to find best, or else
- genome, map 25 bp fragments to genome

arbitrarily assigns a location (mapping quality) NH:i:2 + 3 steps: aligns to transcriptome (if provided), aligns to

#### Tophat output:

- + 2 bam files: mapped and unmapped reads (bam = compressed .sam file)
- + bed files: junctions, insertions, deletions
- log files

#### Sam file sample

SN:SL2.40ch00 LN:21805821 @SQ @SQ SN:SL2.40ch01 LN:90304244 @SQ SN:SL2.40ch02 LN:49918294 @SQ @SQ SN:SL2.40ch03 LN:64840714 SN:SL2.40ch04 LN:64064312 @SQ SN:SL2.40ch05 LN:65021438 @SQ SN:SL2.40ch06 LN:46041636 @SQ SN:SL2.40ch07 LN:65268621 @SQ SN:SL2.40ch08 LN:63032657 @SQ SN:SL2.40ch09 LN:67662091 SN:SL2.40ch10 LN:64834305 @SQ @SQ SN:SL2.40ch11 LN:53386025 @SQ SN:SL2.40ch12 LN:65486253 HWI-EAS339 0011:2:92:9381:20526#0 PG:Z:novoalign ZO:Z:+-AS:i:249 UQ:i:249 NM:i:7 MD:Z:7T1A1T5A31A17A1T7 PQ:i:459 HWI-EAS339\_0011:3:88:16221:9366#0 99 SL2.30ch00 90 5 10S76M = FF PG:Z:novoalign ZO:Z:+- AS:i:88 UQ:i:88 NM:i:0 MD:Z:76 PQ:i:88 SM:i:0 AM:i:0

HWI-EAS339\_0011:2:3:16931:15264#0 163 SL2.30ch00 92 11 3S73M1D9M1H = 139 132 GF PG:Z:novoalign ZO:Z:+- AS:i:143 UQ:i:143 NM:i:3 MD:Z:19T11T41^T9 PQ:i:203 SM:i:0 AM:i:0

\* reference-guided assembly output

99 SL2.30ch00 78 20 7S77M2S = DDD?DDDDDDD6:32-18:6<BCC>DDDDD>>A?A=CACCAC=CCC6@>@;>C@C;>>6@?CCACDD=D=>7BC8B?BAABACCCA SM:i:150 AM:i:21 

#### Sam format column definitions

| Ind | ex | Field Name | Description       |
|-----|----|------------|-------------------|
| 1   |    | QNAME      | Query pair NAM    |
| 2   |    | FLAG       | Bitwise FLAG      |
| 3   |    | RNAME      | Reference sequer  |
| 4   |    | POS        | 1-based leftmost  |
| 5   |    | MAPQ       | MAPping Qualit    |
| 6   |    | CIGAR      | Extended CIGA     |
| 7   |    | MRNM       | Mate Reference :  |
| 8   |    | MPOS       | 1-based leftmost  |
| 9   |    | ISIZE      | Inferred Insert S |
| 10  |    | SEQ        | Query SEQuence    |
| 11  |    | QUAL       | Query QUALity     |
|     |    |            |                   |

 Table 1.1: Brief summary of the SAM format

AE if paired; or Query NAME if unpaired

nce NAME POSition of the clipped sequence ty R string sequence NaMe; "=" if the same as RNAME Mate POSition of the clipped sequence SIZE

e

#### Sam flag field definitions

Excellent explanation for bitwise flags found here: http://seqanswers.com/ forums/showthread.php? t=2301

Tool to translate meaning of bitwise flag: <u>http://</u> picard.sourceforge.net/ explain-flags.html

| Field    | Hex Code | Description      |
|----------|----------|------------------|
| $f_0$    | 0x0001   | the read is pai  |
|          |          | a pair           |
| $f_1$    | 0x0002   | the read is ma   |
|          |          | mally inferred   |
| $f_2$    | 0x0004   | the query sequ   |
| $f_3$    | 0x0008   | the mate is un   |
| $f_4$    | 0x0010   | strand of the o  |
| $f_5$    | 0x0020   | strand of the n  |
| $f_6$    | 0x0040   | the read is the  |
| $f_7$    | 0x0080   | the read is the  |
| $f_8$    | 0x0100   | the alignment    |
|          |          | multiple prima   |
| $f_9$    | 0x0200   | the read fails   |
| $f_{1}0$ | 0x0400   | the read is eit. |

ired in sequencing, no matter whether it is mapped in

- apped in a proper pair (depends on the protocol, nor-
- during alignment)
- uence itself is unmapped
- nmapped
- query (0 for forward; 1 for reverse strand)
- mate
- e first read in a pair
- e second read in a pair
- is not primary (a read having split hits may have ary alignment records)
- platform/vendor quality checks
- her a PCR duplicate or an optical duplicate

 Table 1.2: Brief summary of the SAM format

#### Cigar String

| Op | BAM | Description      |
|----|-----|------------------|
| М  | 0   | alignment mate   |
| I  | 1   | insertion to the |
| D  | 2   | deletion from t  |
| N  | 3   | skipped region   |
| S  | 4   | soft clipping (c |
| Н  | 5   | hard clipping (  |
| Ρ  | 6   | padding (silent  |
| =  | 7   | sequence match   |
| X  | 8   | sequence mism    |
|    |     |                  |

- ch (can be a sequence match or mismatch)
- e reference
- the reference
- from the reference
- clipped sequences present in SEQ)
- (clipped sequences NOT present in SEQ)
- deletion from padded reference)
- 1
- atch

### Return to Exercise 1

the SRA in .sra format and extracted using the SRA toolkit (NCBI).

- Datasets :

  - breaker fruit (two files) - immature fruit (two files)

guided assembly

- Two RNA-seq datasets used in the tomato genome project were downloaded from
  - http://www.ncbi.nlm.nih.gov/sra
- They were already cleaned using fastq-mcf. All data is from S. pimpinellifolium.

In this exercise, we will map the reads to tomato chromosome 4 using reference-



### Exercise 1 (cont'd)

#### Run shell script run\_tuxedo.sh (to be explained)

- 1. Copy ch4\_demo\_dataset.tar.gz from ~/Data to your Desktop and unzip
  - \$cp ~/Data/ch4\_demo\_dataset.tar.gz ~/Desktop
  - \$cd ~/Desktop
  - \$tar -xvf ch4\_demo\_dataset.tar.gz

\$chmod 755 run\_tuxedo.sh \$cp scripts/shell/run\_tuxedo.sh ch4\_demo\_dataset/

4. Change name of SL2.40ch04.fasta to SL2.40ch04.fa 5. run shell script from ~/Desktop/ch4\_demo\_dataset/

\$cd ch4\_demo\_dataset/

\$./run\_tuxedo.sh •

2. Copy shell script run\_tuxedo.sh to ~/Desktop/ch4\_demo\_dataset/ and make it executable

\$mv ch4\_demo\_dataset/bwt2\_index/SL2.40ch04.fasta ch4\_demo\_dataset/bwt2\_index/SL2.40ch04.fa



It is a shell script that runs 6 linux commands (each can be typed in a terminal)

**STEP 1:** uses bowtie2-build to index the tomato reference file.  $\sim 2$ minutes

\$bowtie2-build bwt2\_index/SL2.40ch04.fa bwt2\_indexch04

**STEPS 2-5:** use tophat2 to map each fastq file in ch4\_demo\_dataset/ to the reference. ~3 minutes each step

\$tophat2 -o SRR404334\_tophat\_out/ --no-novel-juncs --no-coverage-search \ bwt2\_index/ SL2.40ch04 breaker/SRR404334/SRR404334\_ch4.fq

### Exercise 1 (cont'd)

#### What does run tuxedo.sh do?



| g | emacs@biodebian.redrover.cornell.edu 🗕 🗆 🗙   |
|---|--|
| F | ile Edit Options Buffers Tools Sh-Script Help  |
|   | 🕒 🖻 🗙 📥 🛣 🥱 💑 📭 📬 🔍 🚍 📧 💿  |
| ^ | ##Shell script for running tophat2 and cuffdiff<br>##Place in /home/bioinfo/Desktop/ch4_demo_dataset to use  |
|   | <pre>#Step 1: run index the reference fasta file for faster match searching<br/>bowtie2-build bwt2_index/SL2.40ch04.fa bwt2_index/SL2.40ch04</pre>   |
| Ш | <pre>#Step 2: run tophat2 on first breaker fastq file tophat2 -o breaker/SRR404334/SRR404334_ch4_thoutno-novel-juncsno-coverag fe-search bwt2_index/SL2.40ch04 breaker/SRR404334/SRR404334_ch4.fq</pre>  |
|   | <pre>#Step 3: run tophat2 on second breaker fastq file tophat2 -o breaker/SRR404336/SRR404336_ch4_thoutno-novel-juncsno-coverag fe-search bwt2_index/SL2.40ch04 breaker/SRR404336/SRR404336_ch4.fq</pre>   |
|   | <pre>#Step 4: run tophat2 on first immature fruit fastq file tophat2 -o immature_fruit/SRR404331/SRR404331_ch4_thoutno-novel-juncsno-₽ coverage-search bwt2_index/SL2.40ch04 immature_fruit/SRR404331/SRR404331_ch4.f₽ fo</pre>  |
|   | #Step 5: run tophat2 on first immature fruit fastq file<br>tophat2 -o immature_fruit/SRR404333/SRR404333_ch4_thoutno-novel-juncsno-₽<br>≤coverage-search bwt2_index/SL2.40ch04 immature_fruit/SRR404333/SRR404333_ch4.f<br>≤q  |
|   | <pre>#Step 6: run cufflinks to find putatively differential expressed genes<br/>cuffdiff -o cuffdiff_out -b bwt2_index/SL2.40ch04.fa -u annotation/ITAG2.3_gen<br/>e_models_ch4.gtf breaker/SRR404334/SRR404334_ch4_thout/accepted_hits.bam,brea<br/>ker/SRR404336/SRR404336_ch4_thout/accepted_hits.bam immature_fruit/SRR404331/S<br/>RR404331_ch4_thout/accepted_hits.bam,immature_fruit/SRR404333/SRR404333_ch4_th<br/>out/accepted_hits.bam</pre> |

### Exercise 1 (cont'd)

**STEP 6:** runs cufflinks, but we will discuss this part later

1. Convert accepted\_hits.bam to accepted\_hits.sam \$samtools view -ho accepted\_hits.sam accepted\_hits.bam 2. How many reads are in the tophat2 output file? \$grep -v "^@" accepted\_hits.sam |wc 2. How many chromosomes are in the reference file? \$grep "^@SQ" accepted\_hits.sam |wc 3. How many reads map to each chromosome?

## Exercise 1 (cont'd)

\$grep -v "^@" accepted hits.sam |cut -f3 |sort |uniq -c

#### De novo Assembly Must clean (quality trim) reads well before de novo assembly!

- Error correction
- In the second second
- Many "isoforms"!!! :(
- Also Velvet + Oases



### De novo Assembly

Highly accessed **Open Access** Method **Corset: enabling differential gene expression analysis for** *de* novo assembled transcriptomes

Nadia M Davidson<sup>1</sup> and Alicia Oshlack<sup>12\*</sup>

" a method that hierarchically clusters contigs using shared reads and testing"

### expression, then summarizes read counts to clusters, ready for statistical

#### Tablet (<u>http://bioinf.scri.ac.uk/tablet/</u>)

#### **Indexing files for Tablet:**

1. reference:

samtools faidx SL2.40ch04.fa

2. bam:

samtools index

#### Load into Tablet:

- 1. reference .fa file
- 2. mapped reads .bam file
- 3. annotations .gff3 file

### 2. Assembly Viewing Output



### 2. Assembly Viewing Output

#### • Samtools tview

| 0 | • |   | C | 2 | ( | 0 | 5 |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   | 7 | Ē6     | er.    | n | ni | n | a |
|---|---|---|---|---|---|---|---|---|--------|---|--------|---|---|---|---|---|---|---|--------|--------|--------|---|---|---|---|---|---|---|--------|--------|---|----|---|---|
| c | A | A | G | т | С | т | С | т | 1<br>T | Z | 1<br>T | G | A | A | т | т | A | A | 1<br>C | 3<br>C | 1<br>C | A | G | т | С | A | G | A | 1<br>C | 4<br>A | 1 | A  | A | A |
|   | ÷ | ÷ | ÷ | ÷ | ÷ | ÷ |   | ÷ | ÷      | ÷ | ÷      | ÷ | ÷ | ÷ | ÷ | ÷ | ÷ | ÷ | ÷      | ÷      | ÷      | ÷ | ÷ | ÷ | î | ÷ | ÷ | ÷ | ÷      | ÷      | ÷ |    | • |   |
| ٠ | ٠ | ٠ |   |   |   |   |   | ٠ | ٠      | ٠ | ٠      | ٠ | ٠ | ٠ | ٠ | ٠ | ٠ | • | •      | •      | ٠      | • | ٠ | ٠ | ٠ | ٠ | ٠ | ٠ | ٠      | •      | • | •  | • | • |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   | -  | - |   |
|   |   | ٠ | ٠ | ٠ |   |   |   |   |        |   |        |   | ٠ | ٠ | ٠ | ٠ |   |   |        |        | А      |   |   | ٠ |   | ٠ | ٠ | ٠ | ٠      |        |   | -  | • |   |
| ٠ | • | ٠ | ٠ | • | • |   |   |   |        |   |        |   | • | ٠ | ٠ | • | • |   |        |        | ٠      |   |   | ٠ | ٠ | • | ٠ | ٠ | •      |        |   | •  | • | • |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   | -  | - |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   | -  | - |   |
|   | ٠ | ٠ |   | • |   |   |   |   |        |   |        |   |   | ٠ | ٠ |   |   |   |        |        |        |   |   | ٠ |   | ٠ | ٠ |   |        |        |   |    | - |   |
|   |   |   |   |   |   |   |   | G |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   | 2 | 2 | 1 | 1 | 2 | 2 | 1 | 2 | 2      | 2 | 2      | 2 | 2 | 2 | 1 |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   | • | • | ٠ | ٠ | • |   |   | • | •      | • | •      |   | • | • | ٠ | ٠ | • |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   |   |   |   |   |   |   |   |        |   |        |   |   |   |   |   |   |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
| ٠ | ٠ | ٠ | ٠ | ٠ | ٠ | • |   |   | ٠      |   | ٠      | ٠ | ٠ | ٠ | ٠ | ٠ | ٠ |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
| ٠ | ٠ | • | ٠ | ٠ | ٠ | • | ٠ | ٠ | ٠      | ٠ | ٠      | ٠ | ٠ | ٠ | ٠ | ٠ | ٠ | • | •      | ٠      |        |   |   |   |   |   |   |   |        |        |   |    |   |   |
|   |   | ٠ | ٠ | ٠ | ٠ |   |   |   |        |   |        |   |   | ٠ | ٠ | ٠ | ٠ |   |        |        |        |   |   |   |   |   |   |   |        |        |   |    |   |   |



### 2. Assembly Viewing Output JBrowse



#### also IGV

### 3. Analysis How good is the assembly?

- \*
- •
- Length of contigs compare to known long transcripts? •
- How many reads map (reference-guided)? \*
- check no DNA contamination •
- coverage of features, biases (introns, 3', ribosomal) \*
- sequencing saturation (new gene discovery) \*

Check for contaminants by using BLAST to search appropriate databases (SeqClean)

How many contigs are there vs how many genes expected - total sequence length? CEGMA

Method

Highly accessed Open Access

**Evaluation of** *de novo* transcriptome assemblies from RNA-Seq data

Bo  $Li^{1^+}$ , Nathanael Fillmore<sup>2+</sup>, Yongsheng Bai<sup>3</sup>, Mike Collins<sup>4</sup>, James A Thomson<sup>456</sup>, Ron Stewart<sup>4</sup> and Colin N Dewey<sup>27\*</sup>







RPKM/FPKM (reads/fragments per 1k of exon per 1M mapped reads.

- •cufflinks
- •RSEM

•Only used as an expression metric within or between samples, not for de testing. Counts - number of reads overlapping a feature (use htseq-count or featurecount)

- •DESeq
- •edgeR

•These don't take read mapping uncertainty into account Other:

•cuffdiff - attempts to correct multimap reads (not default), computes separate variance model for isoforms, may produce high number of false positives (Rapaport et al., 2013), no support for factorial designs

•EBSeq - empirical Bayesian methods, de of isoforms



### Normalization???????

Total Count, Upper Quartile, Median (Med), DESeq (geometric), Trimmed Mean of M values implemented in the edgeR, Quantile, and the Reads Per Kilobase per Million mapped reads (RPKM)!!!!!

#### A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis

Marie-Agnès Dillies<sup>\*</sup>, Andrea Rau<sup>\*</sup>, Julie Aubert<sup>\*</sup>, Christelle Hennequet-Antier<sup>\*</sup>, Marine Jeanmougin<sup>\*</sup>, Nicolas Servant<sup>\*</sup>, Céline Keime<sup>\*</sup>, Guillemette Marot, David Castel, Iordi Estelle, Gregory Guernec, Bernd Jagla, Luc Jouneau, Denis Laloë, Caroline Le Gall, Brigitte Schaëffer, Stéphane Le Crom<sup>\*</sup>, Mickaël Guedj<sup>\*</sup> and Florence Jaffrézic<sup>\*</sup> on behalf of Fhe French StatOmique Consortium

"Based on three real mRNA and one miRNA-seq datasets, we confirm previous observations that RPKM and TC, both of which are still widely in use [40, 41], are ineffective and should be definitively abandoned in the context of differential analysis. "

"The other normalization methods (UQ, Med, DESeq and TMM) perform similarly on the varied datasets considered here, both in terms of the qualitative characteristics of the normalized data and the results of DE analyses."



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

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

Bowtie Extremely fast, general purpose short read aligner

TopHat Aligns RNA-Seq reads to the genome using Bowtie Discovers splice sites

Cufflinks package

Cufflinks Assembles transcripts

Cuffcompare Compares transcript assemblies to annotation

Cuffmerge Merges two or more transcript assemblies

Cuffdiff

Finds differentially expressed genes and transcripts Detects differential splicing and promoter use

> CummeRbund Plots abundance and differential expression results from Cuffdiff

- Cufflinks can also be used for:
- strand-specific RNA-seq
- In a novel transcript discovery in annotated genomes
- ✦ identification of novel splice variants
- detecting transcripts in genomes without annotation

#### For protocols see:

Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks (Trapnell et al., 2012).

#### <u>Challenges</u>

- paralogs/reads that map to more than one place
- Cufflinks divides these evenly across the matches or tries to assign them a location based on expression of surrounding reads.
- Counts-based methods do nothing (although counting software might). Filter on mapping quality?

### Exercise 2

| ₫ | emacs@biodebian.redrover.cornell.edu 🗕 🗖 🗙   |
|---|--|
| F | ile Edit Options Buffers Tools Sh-Script Help  |
|   |  |
| ^ | ##Shell script for running tophat2 and cuffdiff<br>##Place in /home/bioinfo/Desktop/ch4_demo_dataset to use  |
|   | <pre>#Step 1: run index the reference fasta file for faster match searching<br/>bowtie2-build bwt2_index/SL2.40ch04.fa bwt2_index/SL2.40ch04</pre>   |
| Ш | #Step 2: run tophat2 on first breaker fastq file<br>tophat2 -o breaker/SRR404334/SRR404334_ch4_thoutno-novel-juncsno-coverag<br>⊈e-search bwt2_index/SL2.40ch04 breaker/SRR404334/SRR404334_ch4.fq   |
|   | <pre>#Step 3: run tophat2 on second breaker fastq file tophat2 -o breaker/SRR404336/SRR404336_ch4_thoutno-novel-juncsno-coverag fe-search bwt2_index/SL2.40ch04 breaker/SRR404336/SRR404336_ch4.fq</pre>   |
|   | <pre>#Step 4: run tophat2 on first immature fruit fastq file tophat2 -o immature_fruit/SRR404331/SRR404331_ch4_thoutno-novel-juncsno-@ Gcoverage-search bwt2_index/SL2.40ch04 immature_fruit/SRR404331/SRR404331_ch4.f@ Gq</pre>   |
|   | <pre>#Step 5: run tophat2 on first immature fruit fastq file tophat2 -o immature_fruit/SRR404333/SRR404333_ch4_thoutno-novel-juncsno-? coverage-search bwt2_index/SL2.40ch04 immature_fruit/SRR404333/SRR404333_ch4.f? q</pre>   |
|   | <pre>#Step 6: run cufflinks to find putatively differential expressed genes<br/>cuffdiff -o cuffdiff_out -b bwt2_index/SL2.40ch04.fa -u annotation/ITAG2.3_gen<br/>e_models_ch4.gtf breaker/SRR404334/SRR404334_ch4_thout/accepted_hits.bam,brea<br/>ker/SRR404336/SRR404336_ch4_thout/accepted_hits.bam immature_fruit/SRR404331/SP<br/>RR404331_ch4_thout/accepted_hits.bam,immature_fruit/SRR404333/SRR404333_ch4_thP<br/>out/accepted_hits.bam</pre> |

#### What does run\_tuxedo.sh do? (cont'd)

**STEP 6:** Use cuffdiff to detect differentially expressed genes in immature and breaker fruit using known chromosome 4 tomato gene models

\$cuffdiff -o cuffdiff\_out -b bwt2\_index/SL2.40ch04.fa -u annotation/ ITAG2.3\_gene\_models\_ch4.gtf breaker/SRR404334/SRR404334\_ch4\_thout/ accepted\_hits.bam,breaker/SRR404336/SRR404336\_ch4\_thout/accepted\_hits.bam immature\_fruit/SRR404331/SRR404331\_ch4\_thout/ accepted\_hits.bam,immature\_fruit/ SRR404333/SRR404333\_ch4\_thout/accepted\_hits.bam

#### ~5 minutes



#### Cufflinks Output:

| - | rw-rr         | 1 | bioinfo | bioinfo | 12   | Apr | 16  |
|---|---------------|---|---------|---------|------|-----|-----|
| - | rw-rr         | 1 | bioinfo | bioinfo | 115  | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 124  | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 91   | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 115  | Apr | 16  |
|   | rw-rr         | 1 | bioinfo | bioinfo | 343K | Apr | 16  |
|   | rw-rr         | 1 | bioinfo | bioinfo | 183K | Apr | 16  |
|   | rw-rr         | 1 | bioinfo | bioinfo | 333K | Apr | 16  |
|   | rw-rr         | 1 | bioinfo | bioinfo | 564K | Apr | 16  |
|   | rw-rr         | 1 | bioinfo | bioinfo | 348K | Apr | 16  |
|   | rw-rr         | 1 | bioinfo | bioinfo | 189K | Apr | 16  |
|   | rw-rr         | 1 | bioinfo | bioinfo | 346K | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 586K | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 115  | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 466  | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 451  | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 115  | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 124  | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 12   | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 91   | Apr | 16  |
| - | rw-rr         | 1 | bioinfo | bioinfo | 115  | Apr | 16  |
|   | · · · · · · · |   | 1.1     |         |      | -   | ı . |

#### extract genes that have a significant value for differential expression

\$awk -F "\t" 'BEGIN {OFS = "\t"} \$14 = "yes" {print \$0}' gene\_exp.diff > significant\_genes.txt

We will analyze these results in R tomorrow!

11:17 cds.count tracking ll:17 cds.diff 11:17 cds\_exp.diff 11:17 cds.fpkm tracking 11:17 cds.read\_group\_tracking 11:17 gene exp.diff 11:17 genes.count\_tracking 11:17 genes.fpkm\_tracking 11:17 genes.read\_group\_tracking 11:17 isoform exp.diff 11:17 isoforms.count\_tracking 11:17 isoforms.fpkm\_tracking 11:17 isoforms.read\_group\_tracking 11:17 promoters.diff 11:17 read\_groups.info ll:17 run.info 11:17 splicing.diff 11:17 tss\_group\_exp.diff 11:17 tss\_groups.count\_tracking 11:17 tss\_groups.fpkm\_tracking 11:17 tss\_groups.read\_group\_tracking

### Other useful tools

- ✦ Bedtools used to compare features.
- snpeff gives predicted effect of SNP (synonymous, nonsynonymous, etc).
- ✦ sra toolkit for converting NGS data files from the SRA.
- ✦ Picard duplicate detection, etc.
- Prinseq duplicate removal
- Dindel indel calling from NGS mapping data.
- ✦ seqanswers