

BINF6308: Lecture 8

Instructor - Andrej Savol (a.savol@northeastern.edu)

TA Section 12579 (in-person): Lindsey Alexanian (alexanian.l@northeastern.edu)

TA Section 12580 (online): Yicheng Zhang (zhang.yicheng@northeastern.edu)

Meeting times:

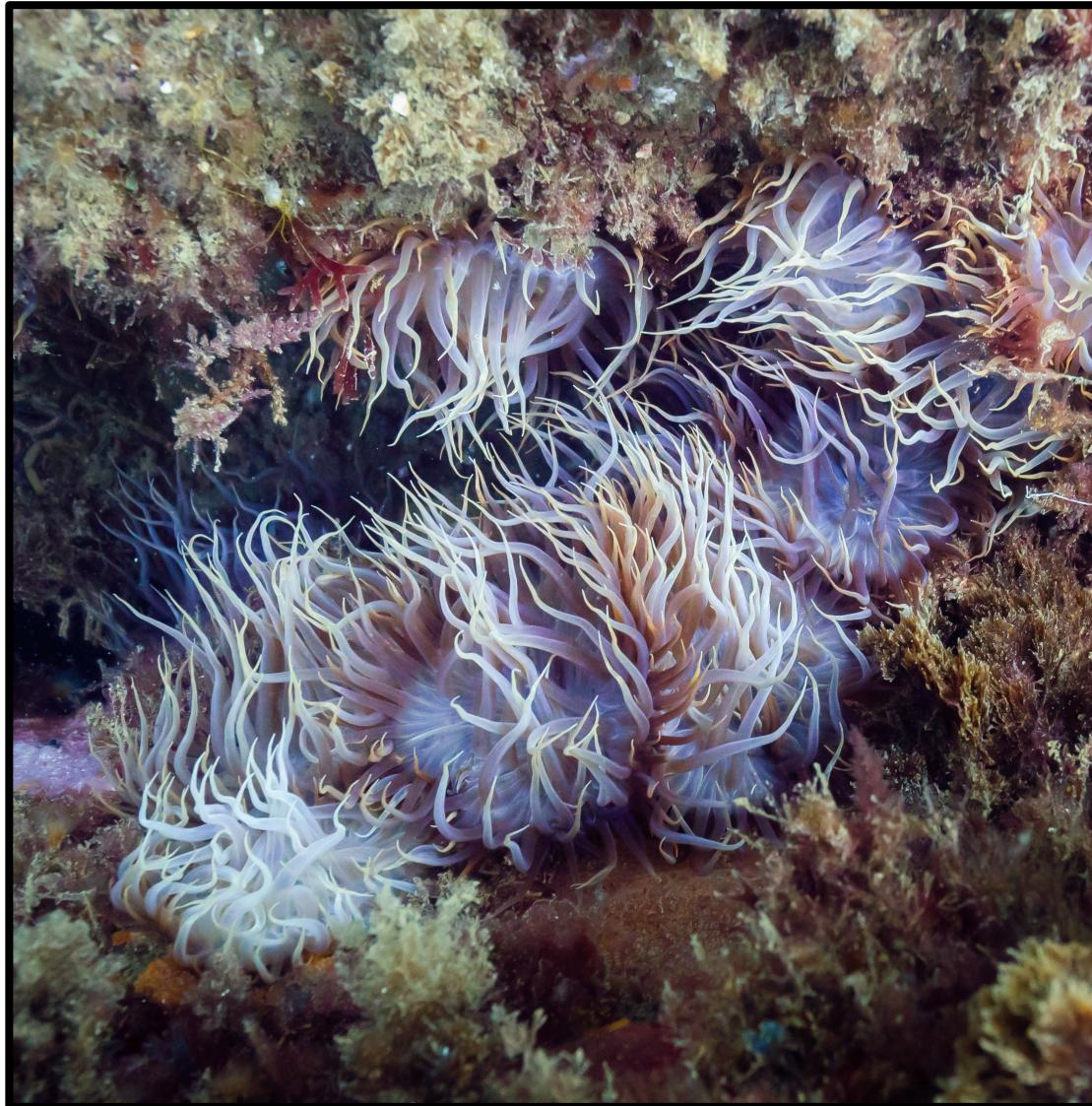
Section 1 (in-person): Monday 5:30-8:30pm ISEC 140

Section 2(online): Asynchronous

Lecture 8: Outline

- Transcriptome Assembly
 - Assembly v. alignment
 - Reference-guided v. *de novo*
 - Index-based alignment
 - Using Trinity (demo)
- Links in Unix
- Using IGV

Aiptasia



Diego Delso, CC BY-SA 4.0 <<https://creativecommons.org/licenses/by-sa/4.0>>, via Wikimedia Commons

Using Trinity for transcriptome assembly

Reconstruct
this

CTAGGCCCTCAATTTT
CTCTAGGCCCTCAATTTT
GGCTCTAGGCCCTCATTTTT
CTCGGCTCTAGCCCCTCATTTT
TATCTCGACTCTAGGCCCTCA
TATCTCGACTCTAGGCC
TCTATATCTCGGCTCTAGG
GGCGTCTATATCTCG
GGCGTCGATATCT
GGCGTCTATATCT

From these

→ **GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTT**

Courtesy of [Ben Langmead](#). Used with permission.

<http://www.langmead-lab.org/teaching-materials/>

Understanding base coverage

(I.e. average coverage: the average number of reads covering a position in the genome)

CTAGGCCCTCAATTTT	
CTCTAGGCCCTCAATTTT	
GGCTCTAGGCCCTCATTTTT	
CTCGGCTCTAGCCCCTCATTTT	
TATCTCGACTCTAGGCCCTCA	177 nucleotides
TATCTCGACTCTAGGCC	
TCTATATCTGGCTCTAGG	
GGCGTCTATATCTCG	
GGCGTCGATATCT	
GGCGTCTATATCT	
GGCGTCTATATCTGGCTCTAGGCCCTCATTTTT	35 nucleotides

What is the approximate coverage here?

Understanding base coverage

(I.e. average coverage: the average number of reads covering a position in the genome)

The diagram shows a vertical stack of DNA sequences, each consisting of a string of letters (A, T, C, G) representing the bases. A red bracket is drawn vertically across the stack, spanning from the top sequence down to the bottom sequence. The bottom sequence is highlighted in red. An arrow points upwards from the bottom sequence towards the text below.

CTAGGCCCTCAATTTT
CTCTAGGCCCTCAATTTT
GGCTCTAGGCCCTCATTTTT
CTCGGCTCTAGCCCTCATTTT
TATCTCGACTCTAGGCCCTCA
TATCTCGACTCTAGGCC
TCTATATCTCGGCTCTAGG
GGCGTCTATATCTCG
GGCGTCGATATCT
GGCGTCTATATCT
GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTT

Coverage at this position = 6

Courtesy of [Ben Langmead](#). Used with permission. <http://www.langmead-lab.org/teaching-materials/>

What is the approximate coverage here?

Here are standard alignment steps; which steps are different in Assembly?

Pre-alignment QC
("alignment-free" QC)

Base-quality score summaries (FastQC). Base distributions

Cheap/Fast

Pre-processing

Trimming. Fastq interleaving. Demultiplex ("demux")

Reasonable/Medium

Alignment

BWA, STAR, HISAT2, GSNAP, bowtie, GMAP

Expensive/Slow

Sorting

Sort bam file by position or read name for easy access.

Indexing

Generate accessory file ("bam index") for fast read access.

Post-alignment
QC

% aligned. Intronic/exonic content. Insert size.

Reference-guided assembly: why assemble a genome if a reference already exists?

Say two reads truly originate from overlapping stretches of the genome. Why might there be differences?

The diagram illustrates two DNA sequences. The top sequence is TATCTCGACTCTAGGCC, with vertical lines below each base pair indicating alignment. The bottom sequence is TCTATATCTCGGCTCTAGG. A red arrow points to the second 'A' in the bottom sequence, which does not align with the 'T' in the top sequence at that position, highlighting a sequencing error or a different base being read.

Reference-guided assembly: why assemble a genome if a reference already exists?

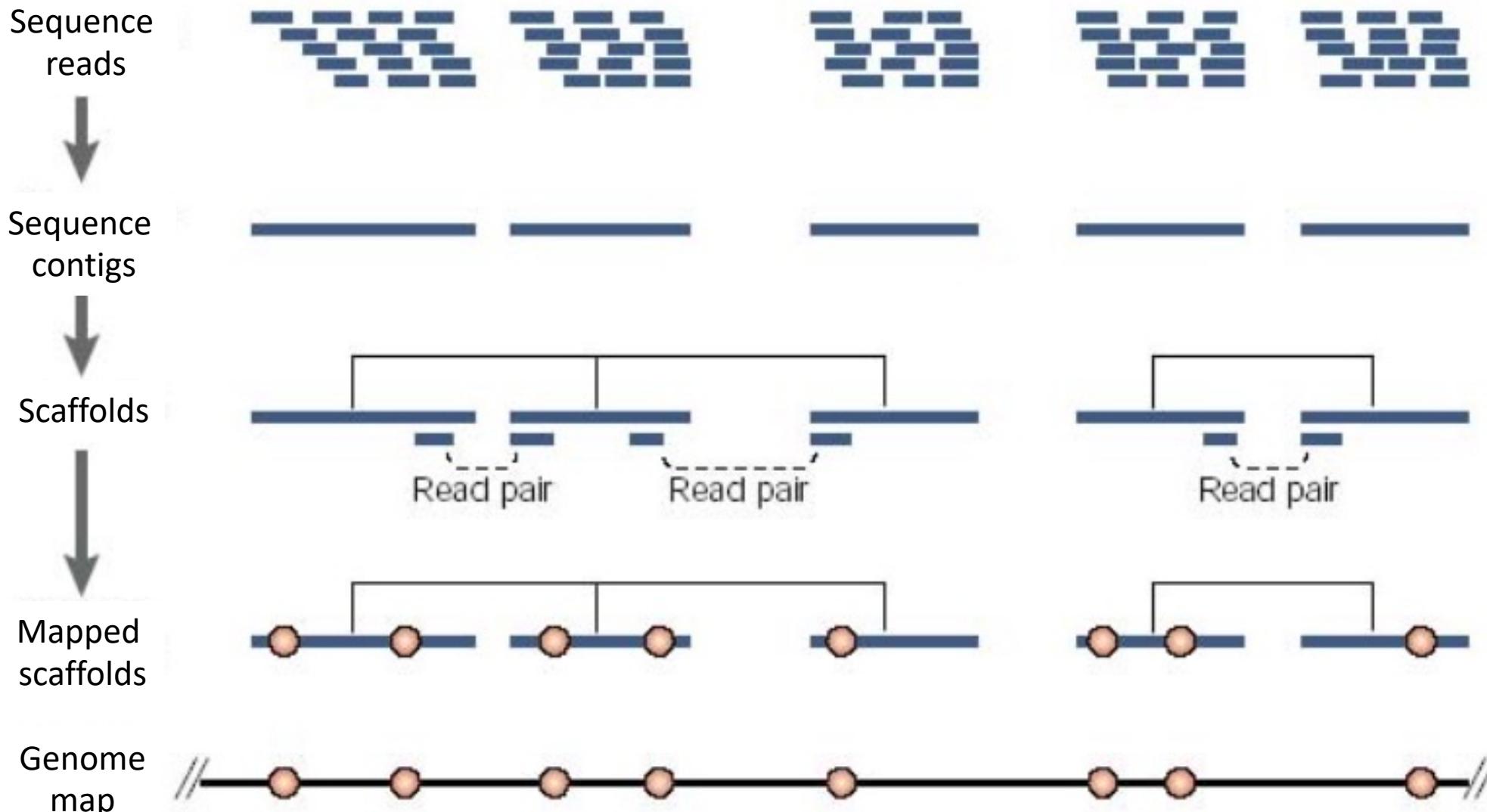
Say two reads truly originate from overlapping stretches of the genome. Why might there be differences?

TATCTCGACTCTAGGCC
||||||| |||||
TCTATATCTCGGCTCTAGG

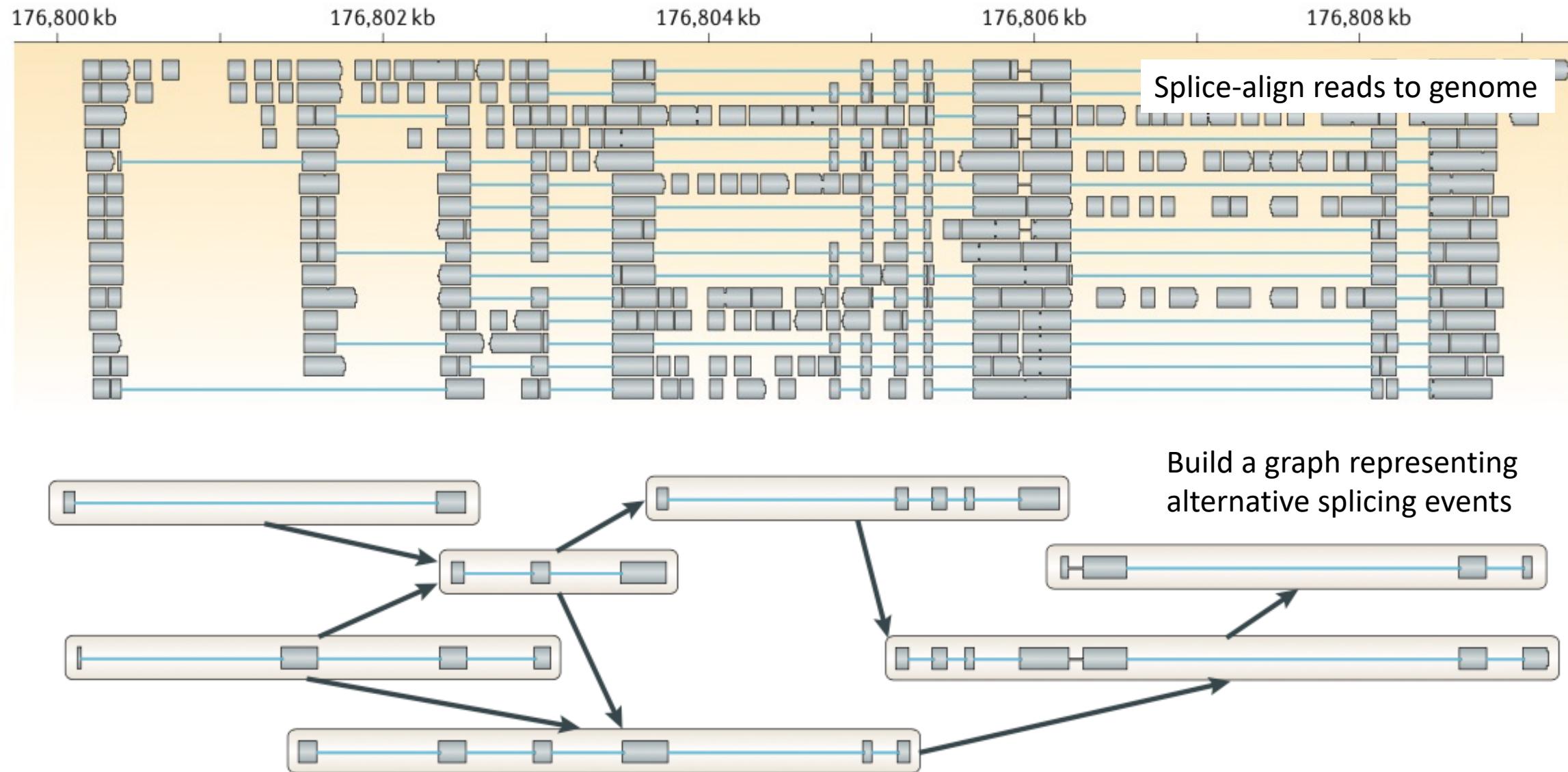
1. Sequencing error
2. Difference between inherited *copies* of a chromosome

E.g. humans are diploid; we have two copies of each chromosome, one from mother, one from father. The copies can differ:

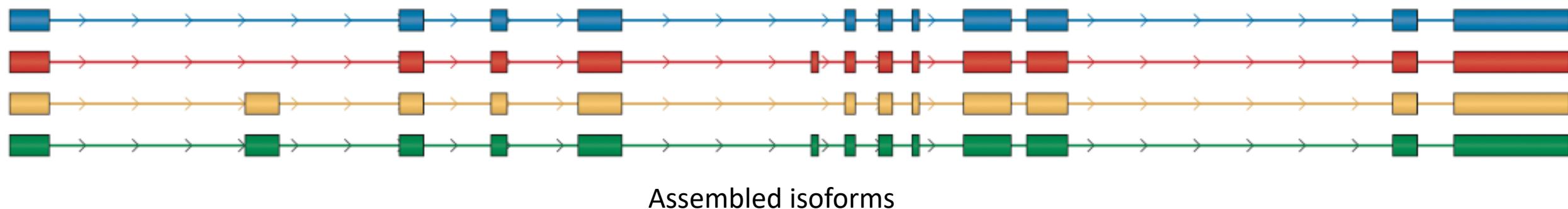
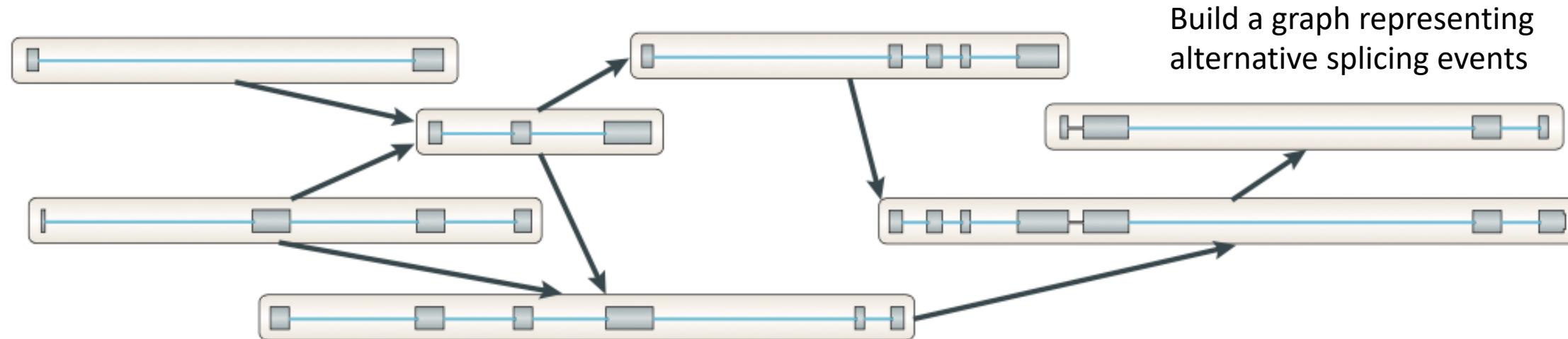
de novo whole genome assembly



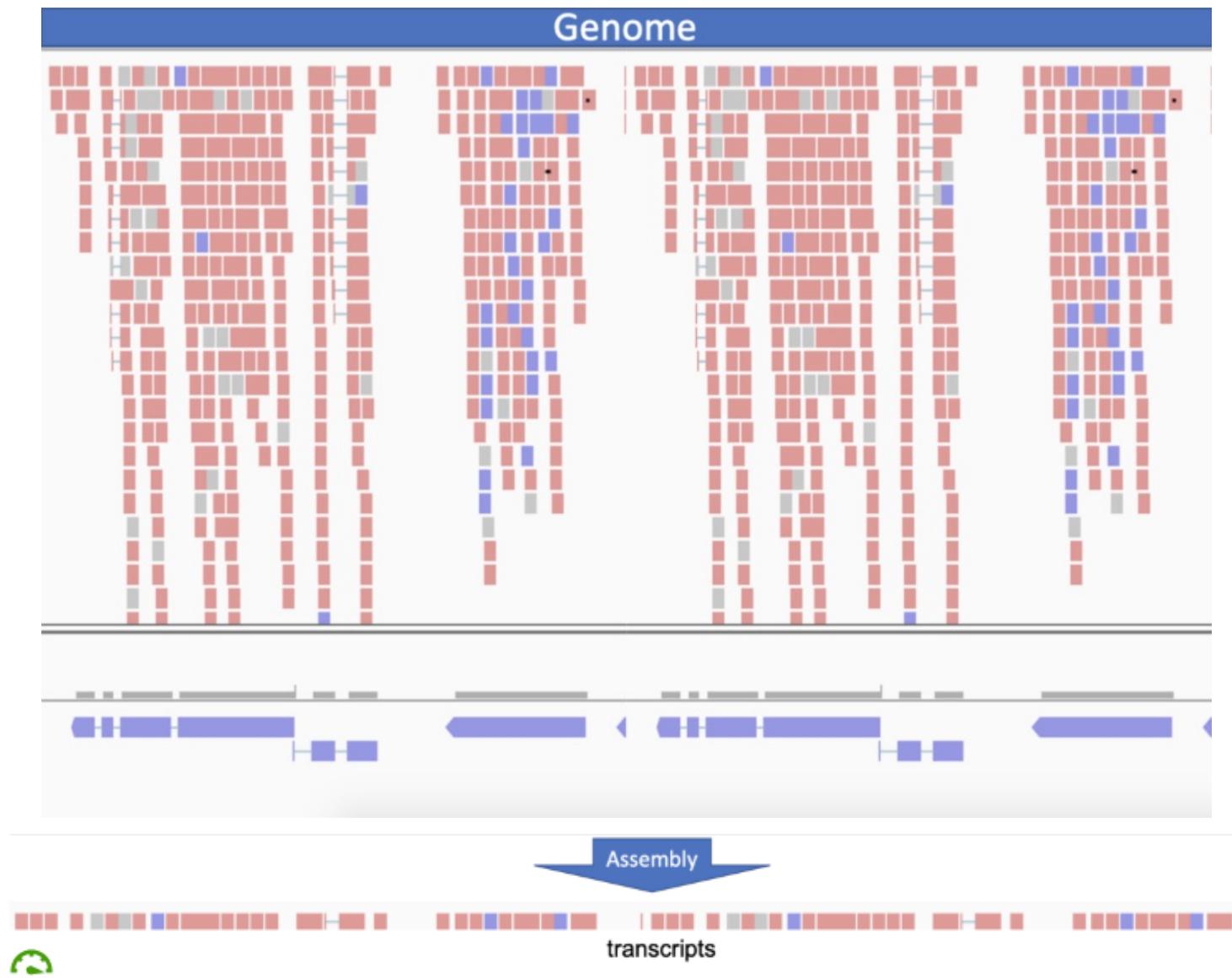
Transcriptome assembly



Transcriptome assembly

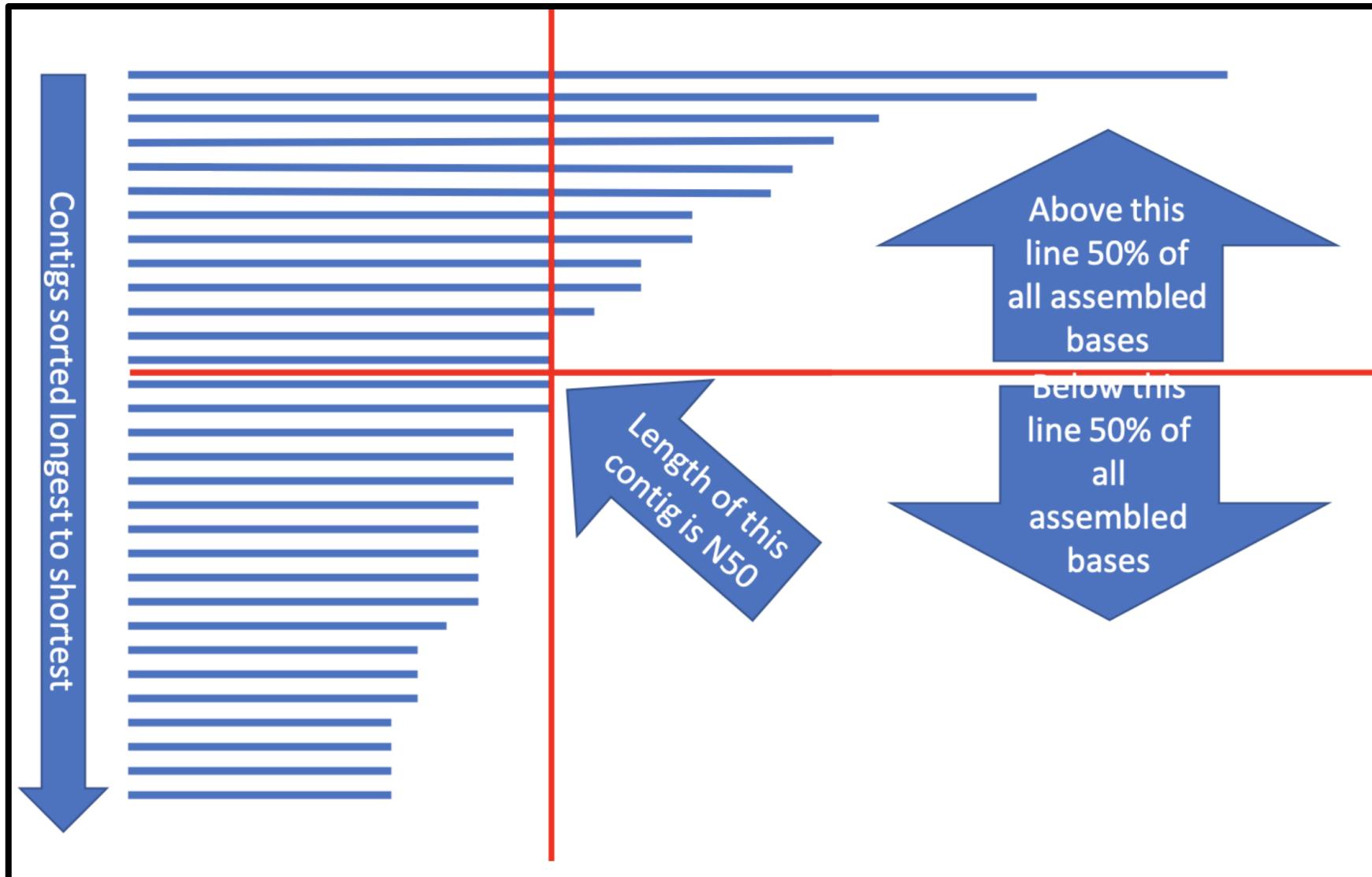


Genome assembly -> transcriptome assembly



How big are the assembled contigs?

N50 indicates the minimal length among the contigs that contain 50 of all assembled bases.



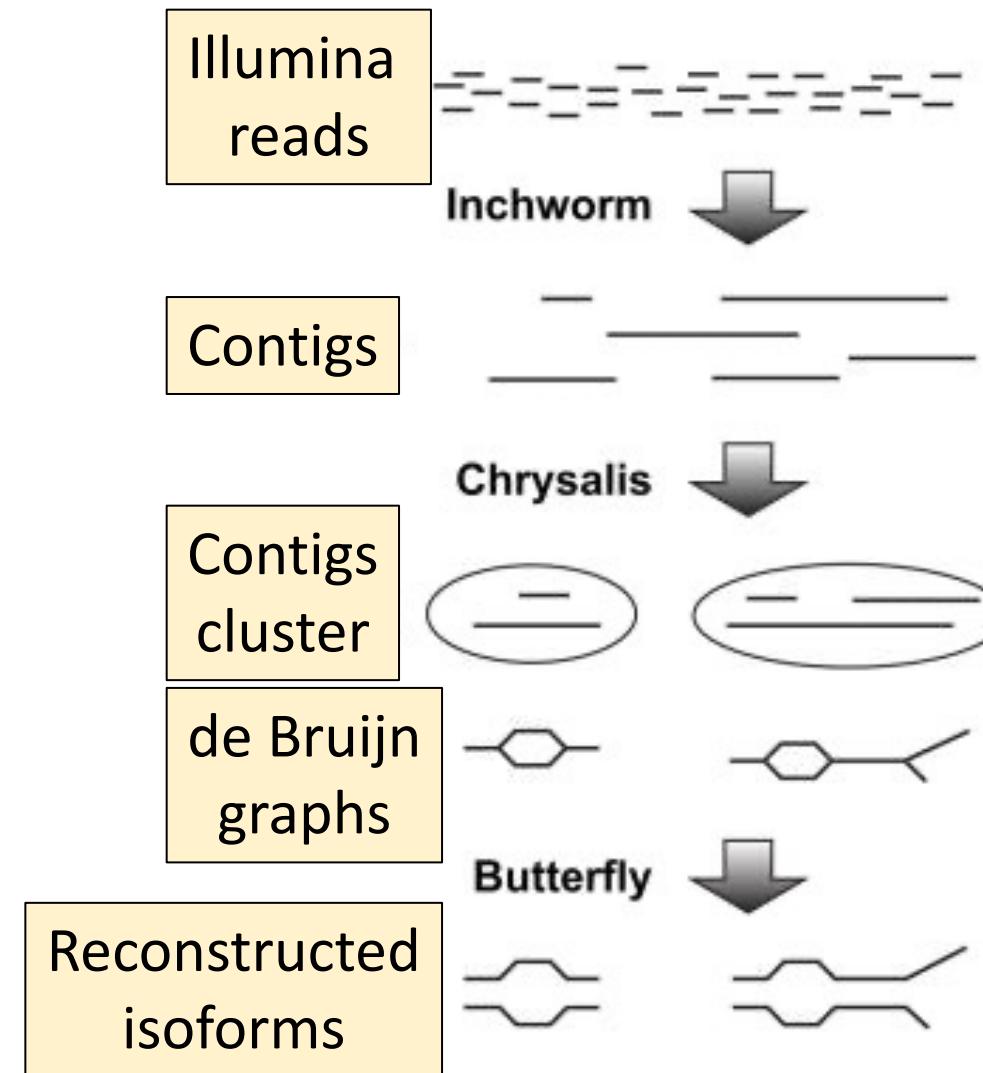
Will N40, N30 be bigger or smaller than N50?

Using Trinity for transcriptome assembly

```
[a.savol@defiance[RNA-Seq]# /usr/local/programs/trinityrnaseq-Trinity-v2.8.4/Trinity
```

```
#####
#
# *Note, a typical Trinity command might be:
#
#     Trinity --seqType fq --max_memory 50G --left reads_1.fq --right reads_2.fq --CPU 6
#
#         (if you have multiple samples, use --samples_file ... see above for details)
#
# and for Genome-guided Trinity, provide a coordinate-sorted bam:
#
#     Trinity --genome_guided_bam rnaseq_alignments.csorth.bam --max_memory 50G
#             --genome_guided_max_intron 10000 --CPU 6
#
# see: /usr/local/programs/trinityrnaseq-Trinity-v2.8.4/sample_data/test_Trinity_Assembly/
#       for sample data and 'runMe.sh' for example Trinity execution
#
#     For more details, visit: http://trinityrnaseq.github.io
#
#####
```

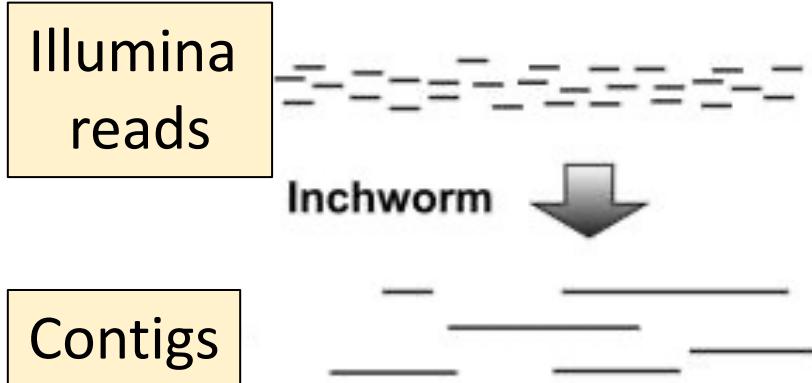
Using Trinity for transcriptome assembly



Trinity also uses kmers for faster contig assembly

```
[a.savol@defiance[trinity_out_dir]# head -n 10 jellyfish.kmers.fa
>526
AAAAAAAAAAAAAAAAAAAAAAA
>1
AGACTGGCCTCTGGTTCAACCCCTGC
>1
TGTCTTATTCTAAGCCACACGCAGA
>105
TGTGGGCTGCGCATTCTCTTCCGGA
>1
AAGACCCCGTGTGTTGATGAAGCAG
```

Using Trinity for transcriptome assembly



```
a.savol@defiance[backup.trinity_out_dir]# vim inchworm.K25.L25.fa
```

```
127 >a3;3 total_counts: 76 Seed: 5 K: 25 length: 46
128 AACAGGCTGCCGGGCTAAGGCCCTAGGATTAGCCCACAGGAGCC
129 >a4;2 total_counts: 938 Seed: 2 K: 25 length: 562
130 TGGTTTCAGTCCTCCTGAGATAATGCCCTGGGTGTGACTTTAGTAGATGCAGA
131 GAGATTGCTCTGCCACCCTCCAGGCTCCACAGCGTGGAACTCTGCATGTTGTTCACT
132 AACACAGGCCAGTCCCAGCACCCCTAGACAGTGAGTCAGTGGAGAGGCGTAGAGAGG
133 CTGAGACTGGCTGCCTACAATCTAACCTCGATAACTTAGGACTAACTCTCCTAGA
134 ATGTTGGTTCTTAAGATAAGGCAGCTCATATAAGTAAGCTATCACACAGGCCAGCTCT
135 GCTCAGAATGGGGGAAGGTAGCACTCTGACCCCTTGCTCTCGAGGGCCTGTCTTCC
136 TCCTCGGCTCTTACTCCTTGAGACTGCACCCCTAACCTCTGTCTCATCTCATTATA
137 GGCCTGGACTCTTGAAGTCTGATATTGATTACAGTGGATGTCAATTGTCTGCTTGCAG
138 CAGGGCACCTGGAACCTATTAGATAAGCAAGAGGTTCCACTGTGAACCCAACGAGCTG
139 CGAGCTGTGACCTTGCTCTCAG
140 >a5;2 total_counts: 79 Seed: 2 K: 25 length: 61
141 TCTTGGGCCCCCAGACACCCCCCTGCCCCCTACACGCCCTGACAGTGTGTACAGCTA
142 C
143 >a6;2 total_counts: 156 Seed: 3 K: 25 length: 89
144 GCCCAATGTTGCCTTGGGAACATTATCTAACCTTATTGAGCAGTTAGGACAATGTCT
145 ACAACTTAAGAGGTAAAGCACCCCTCTAC
146 >a7;2 total_counts: 38 Seed: 2 K: 25 length: 48
147 GGAGGGGCCAGAGGCCATCCAAAAGAGGGGTGTGCACCCTCT
```

Assessing Trinity accuracy

```
a.savol@defiance[test_Trinity_Assembly]# ${TRINITY_HOME}/util/misc/illustrate_ref_comparison.pl  
_indiv_ex_sample_derived/refSeqs.fa trinity_out_dir/Trinity.fasta 90 > assembly_comparison.txt
```

```
1 Loaded 197879 letters in 84 sequences  
2 Searched 29910 bases in 11 sequences  
3 1-2919 [=====] mm9chr5-NM_172722;mm9chr5-231713  
4 1-2919+ -----> TRINITY_DN6_c0_g1_i3 24-2942:5399 (54.07% aln, 100.00% ID)  
5 1-1728+ -----> TRINITY_DN6_c0_g1_i2 24-1751:5528 (31.26% aln, 100.00% ID)  
6 1729-2919+ -----> TRINITY_DN6_c0_g1_i2 1881-3071:5528 (21.54% aln, 100.00% ID)  
7 2649-2919+ -----> TRINITY_DN6_c0_g1_i1 77-347:2804 (9.66% aln, 100.00% ID)  
8  
9 1-1536 [=====] mm9chr6-NM_001083315;mm9chr6-64213  
10 1-1536+ -----> TRINITY_DN0_c0_g1_i1 567-2102:2364 (64.97% aln, 100.00% ID)  
11 1-1056+ -----> TRINITY_DN0_c0_g1_i5 567-1622:2213 (47.72% aln, 100.00% ID)  
12 1-511+ -----> TRINITY_DN0_c0_g1_i2 567-1077:2433 (21.00% aln, 100.00% ID)  
13 23-1536+ -----> TRINITY_DN0_c0_g1_i3 425-1938:2200 (68.82% aln, 100.00% ID)  
14 23-511+ -----> TRINITY_DN0_c0_g1_i4 425-913:2269 (21.55% aln, 100.00% ID)  
15 512-1536+ -----> TRINITY_DN0_c0_g1_i4 983-2007:2269 (45.17% aln, 100.00% ID)  
16 512-1536+ -----> TRINITY_DN0_c0_g1_i2 1147-2171:2433 (42.13% aln, 100.00% ID)  
17 1208-1536+ -----> TRINITY_DN0_c0_g1_i5 1623-1951:2213 (14.87% aln, 100.00% ID)  
18  
19 1-1605 [=====] mm9chr6-NM_022332;mm9chr6-64213  
20 1-1605+ -----> TRINITY_DN0_c0_g1_i2 567-2171:2433 (65.97% aln, 100.00% ID)  
21 1-512+ -----> TRINITY_DN0_c0_g1_i1 567-1078:2364 (21.66% aln, 100.00% ID)  
22 1-512+ -----> TRINITY_DN0_c0_g1_i5 567-1078:2213 (23.14% aln, 100.00% ID)  
23 23-1605+ -----> TRINITY_DN0_c0_g1_i4 425-2007:2269 (69.77% aln, 100.00% ID)  
24 23-512+ -----> TRINITY_DN0_c0_g1_i3 425-914:2200 (22.27% aln, 100.00% ID)  
25 582-1605+ -----> TRINITY_DN0_c0_g1_i3 915-1938:2200 (46.55% aln, 100.00% ID)  
26 582-1605+ -----> TRINITY_DN0_c0_g1_i1 1079-2102:2364 (43.32% aln, 100.00% ID)  
27 582-1125+ -----> TRINITY_DN0_c0_g1_i5 1079-1622:2213 (24.58% aln, 100.00% ID)  
28 1277-1605+ -----> TRINITY_DN0_c0_g1_i5 1623-1951:2213 (14.87% aln, 100.00% ID)  
29  
30 1-1200 [=====] mm9chr9-NM_026942;mm9chr9-69106  
31 1-240+ -----> TRINITY_DN3_c0_g1_i4 66-305:769 (31.21% aln, 100.00% ID)  
32 133-1200+ -----> TRINITY_DN3_c0_g1_i1 871-1938:2619 (40.78% aln, 100.00% ID)  
33 133-240+ -----> TRINITY_DN3_c0_g1_i5 871-978:2899 (3.73% aln, 100.00% ID)  
34 133-240+ -----> TRINITY_DN3_c0_g1_i3 871-978:1442 (7.49% aln, 100.00% ID)  
35 239-1200+ -----> TRINITY_DN3_c0_g1_i2 664-1625:2306 (41.72% aln, 100.00% ID)  
36 241-1200+ -----> TRINITY_DN3_c0_g1_i5 1259-2218:2899 (33.11% aln, 100.00% ID)  
37  
38 1-1104 [=====] mm9chr4-NM_010598;mm9chr4-16498  
39 1-1104+ -----> TRINITY_DN2_c0_g1_i3 324-1427:3739 (29.53% aln, 100.00% ID)  
40 1-77+ ---> TRINITY_DN2_c0_g1_i2 324-400:3697 (2.08% aln, 100.00% ID)  
41 1-77+ ---> TRINITY_DN2_c0_g1_i1 175-251:3548 (2.17% aln, 100.00% ID)  
42 120-1104+ -----> TRINITY_DN2_c0_g1_i2 401-1385:3697 (26.64% aln, 100.00% ID)  
43 120-1104+ -----> TRINITY_DN2_c0_g1_i1 252-1236:3548 (27.76% aln, 100.00% ID)  
44 501-540+ ---->
```

Assessing Trinity accuracy

1 Loaded 197879 letters in 84 sequences

2 Searched 29910 bases in 11 sequences

```

3 1-2919  [=====] mm9chr5-NM_172722;mm9chr5-231713
4 1-2919+ -----> TRINITY_DN6_c0_g1_i3 24-2942:5399 (54.07% aln, 100.00% ID)
5 1-1728+ -----> TRINITY_DN6_c0_g1_i2 24-1751:5528 (31.26% aln, 100.00% ID)
6 1729-2919+ -----> TRINITY_DN6_c0_g1_i2 1881-3071:5528 (21.54% aln, 100.00% ID)
7 2649-2919+ -----> TRINITY_DN6_c0_g1_i1 77-347:2804 (9.66% aln, 100.00% ID)

```

Contig TRINITY_DN6_c0_g1_i3 (length=5399)

Using Trinity for transcriptome assembly

```
a.savol@defiance[test_Trinity_Assembly]# ./runMe.sh
```

Group coding challenge

```
1 #!/usr/bin/env python
2
3 # define a function is_palindrome that takes a string as argument and returns a boolean
4 # that is True if the string is a palindrome, and False if it is not
5
6 def is_palindrome():
7     pass
```

```
In [2]: is_palindrome('deified')
Out[2]: True
```

```
In [3]: is_palindrome('CTCC')
Out[3]: False
```

```
In [4]: is_palindrome('C')
Out[4]: True
```

Group coding challenge: 3 possible solutions

```
1 def is_palindrome(s):  
2     return(s == s[::-1]) # via slicing (simplest)  
3
```

Group coding challenge: 3 possible solutions

```
1 def is_palindrome(s):
2     return(s == s[::-1]) # via slicing (simplest)
3
```

```
4 def is_palindrome(s):
5     all_matches = True
6     for i,c in enumerate(s,1):
7         if c != s[-i]: # walk backward through the string with negative indices
8             all_matches = False
9
10    return(all_matches)
```

Group coding challenge: 3 possible solutions

```
1 def is_palindrome(s):
2     return(s == s[::-1]) # via slicing (simplest)
3
```

```
4 def is_palindrome(s):
5     all_matches = True
6     for i,c in enumerate(s,1):
7         if c != s[-i]: # walk backward through the string with negative indices
8             all_matches = False
9
10    return(all_matches)
```

```
12 def is_palindrome(s):
13     pal = True
14     for fwd,rev in zip(s,s[::-1]):
15         # will evaluate pal irreversibly to False if mismatch
16         pal = pal and fwd == rev
17         print('%s==%s: pal=%r' % (fwd,rev,pal))
18
19     return(pal)
```

Using symlinks

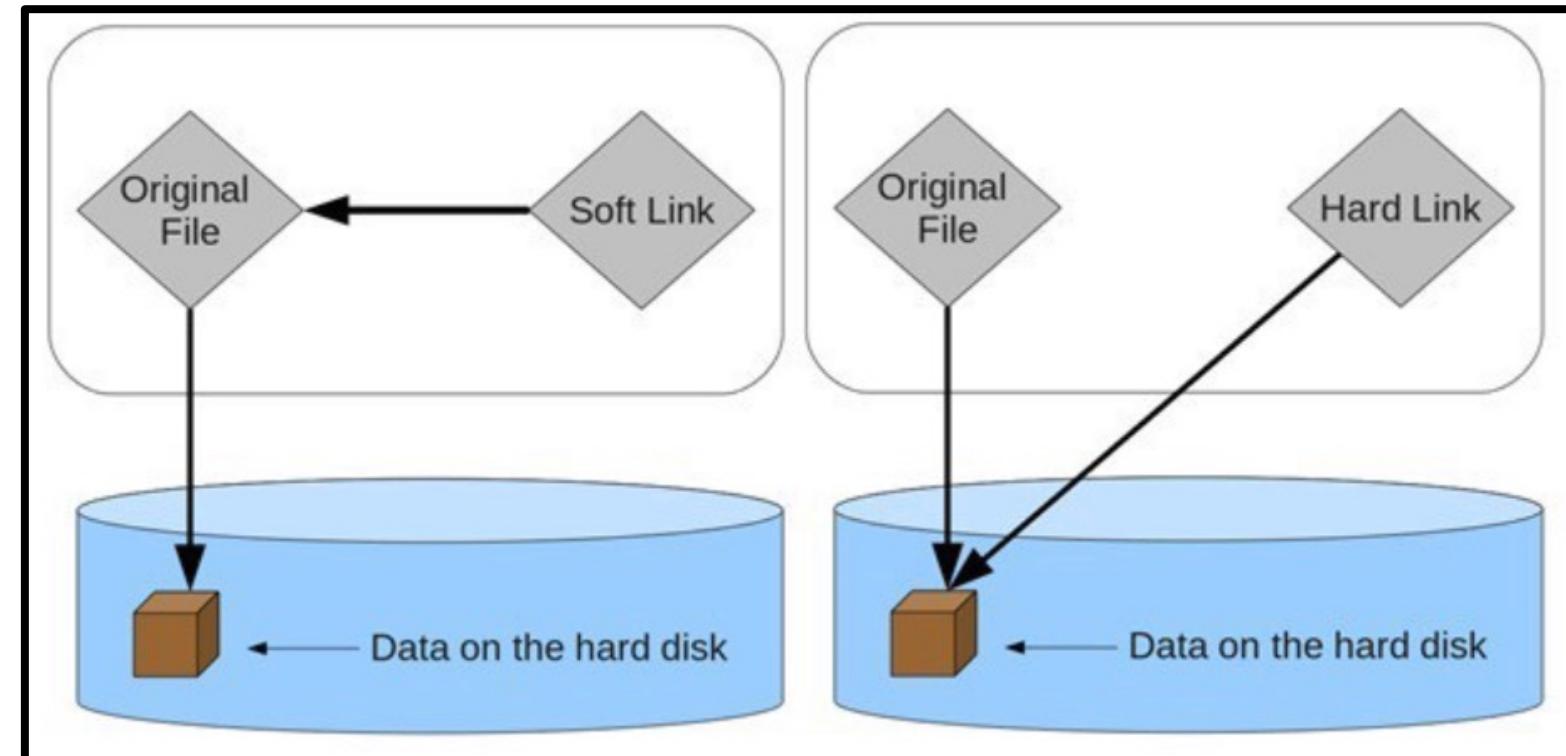
A symbolic link (aka *symlink* or *soft link*) is a special type of file that points to another file or directory.

Symlinks enable multiple file names to point to the same file

Symlinks help avoid unnecessary file copying.

```
a.savol@defiance[RNA-Seq]# ln -s ../../module-06-a-savol/RNA-Seq/Paired Paired
a.savol@defiance[RNA-Seq]# ll
total 0
drwxr-xr-x. 2 a.savol BioInfo 74 Mar 20 17:36 .
drwxr-xr-x. 4 a.savol BioInfo 84 Mar 20 17:18 ..
lrwxrwxrwx. 1 a.savol BioInfo 34 Mar 20 17:20 bam -> ../../module06-a-savol/RNA-Seq/bam
lrwxrwxrwx. 1 a.savol BioInfo 38 Mar 20 17:36 Paired -> ../../module-06-a-savol/RNA-Seq/Paired
lrwxrwxrwx. 1 a.savol BioInfo 34 Mar 20 17:20 sam -> ../../module06-a-savol/RNA-Seq/sam
lrwxrwxrwx. 1 a.savol BioInfo 39 Mar 20 17:20 Unpaired -> ../../module06-a-savol/RNA-Seq/Unpaired
```

	Soft link (symbolic link)	Hard link
Command	<code>ln -s <TARGET> <LINK_NAME></code>	<code>ln <TARGET> <LINK_NAME></code>
Use case	Access a file from multiple locations (or even multiple computers)	Access same data and be robust to name changes
Consequences of renamed/moved/deleted Original file	Broken link	Nothing
Consequences of deleted link	None	None



Merging bam files

```
1 #!/usr/bin/env bash
2 # mergeAll.sh
3 # Usage: bash scripts/mergeAll.sh 1>results/logs/mergeAll.log 2>results/logs/mergeAll.err
4
5 # List the files with ls and redirect (>) the output to bamIn.txt
6 ls data/bam/Aip*.sorted.bam > data/bam/bamIn.txt
7
8 # Merge the files using bamIn.txt with the -b option
9 samtools merge -b data/bam/bamIn.txt data/bam/AipAll.bam
```

Running Transcriptome Assembly with Trinity

```
1 #!/usr/bin/env bash
2 # runTrinity.sh
3 # Usage: bash scripts/runTrinity.sh 1>results/logs/trinity_guided.log 2>results/logs/trinity_guided.err
4
5 Trinity \
6 --genome_guided_bam data/bam/AipAll.bam \
7 --genome_guided_max_intron 10000 \
8 --max_memory 10G --CPU 4 \
9 --output results/trinity_guided
```

Creating a pipeline

```
1 #!/usr/bin/bash
2 #SBATCH --partition=short          # choose from debug, express, or short
3 #SBATCH --job-name=trinity
4 #SBATCH --time=20:00:00            # the code pieces should run in far less than 1 hour
5 #SBATCH -N 1                      # nodes requested
6 #SBATCH -n 4                      # task per node requested
7 #SBATCH --mem=10Gb
8 #SBATCH --exclusive
9 #SBATCH --output="batch-%x-%j.output"    # where to direct standard output; will be batch-jobname-jobID.output
10 #SBATCH --mail-type=ALL
11 #SBATCH --mail-user=<your user ID>@northeastern.edu # Update to your user name!
12
13 # Usage: sbatch sbatch_transcriptome.sh
14 # Assumes input data is in /home/$USER/AiptasiaRNASeq/data/
15
16 echo "Starting our analysis $(date)"
17
18 echo "Loading our BINF6308 Anaconda environment, which includes Trinity."
19 module load anaconda3/2021.11
20 source activate BINF-12-2021
21 echo "Loading samtools."
22 module load samtools/1.10
23
24 echo "Make directory for data files"
25 mkdir -p data/
26
27 # part of a bigger sbatch script (e.g., #SBATCH lines above)
28 echo "Moving trimmed FASTQ data to the working directory"
29 cp -r /home/$USER/AiptasiaRNASeq/data/trimmed data/trimmed
30 echo "Moving alignment files to working directory"
31 cp -r /home/$USER/AiptasiaRNASeq/data/bam data/bam
32 cp -r /home/$USER/AiptasiaRNASeq/data/sam data/sam
33
34 echo "Make directory for log files"
35 mkdir -p results/logs/
```

Creating a pipeline

```
37 echo "Starting Guided Assembly $(date)"
38 echo "Merge all BAM alignment files $(date)"
39 bash scripts/mergeAll.sh 1>results/logs/$SLURM_JOB_NAME-$SLURM_JOB_ID-mergeAll.log 2>results/logs/$SLURM_JOB_NAME-$SLURM_JOB_ID-mergeAll.err
40
41 echo "Assemble the Guided Transcriptome $(date)"
42 bash scripts/runTrinity.sh 1>results/logs/$SLURM_JOB_NAME-$SLURM_JOB_ID-runTrinity.log 2>results/logs/$SLURM_JOB_NAME-$SLURM_JOB_ID-runTrinity.err
43
44 echo "Analyze the Guided Transcriptome $(date)"
45 bash scripts/analyzeTrinity.sh 1>results/$SLURM_JOB_NAME-$SLURM_JOB_ID-trinity_guided_stats.txt 2>results/logs/$SLURM_JOB_NAME-$SLURM_JOB_ID-analyzeTrinity.err
46
47 echo "Guided Assembly complete $(date)"
48
49
50 echo "Starting De Novo Assembly $(date)"
51 echo "Assemble the De Novo Transcriptome $(date)"
52 bash scripts/trinityDeNovo.sh 1>results/logs/$SLURM_JOB_NAME-$SLURM_JOB_ID-trinityDeNovo.log 2>results/logs/$SLURM_JOB_NAME-$SLURM_JOB_ID-trinityDeNovo.err
53
54 echo "Analyze the De Novo Transcriptome $(date)"
55 bash scripts/analyzeTrinityDeNovo.sh 1>results/$SLURM_JOB_NAME-$SLURM_JOB_ID-trinity_de_novo_stats.txt 2>results/logs/$SLURM_JOB_NAME-$SLURM_JOB_ID-analyzeTrinityDeNovo.err
56
57 echo "De Novo Assembly complete $(date)"
58
59 echo "Moving key files back to /home"
60 cp -r results/trinity_guided /home/$USER/AiptasiaRNASeq/data/trinity_guided
61 cp -r results/trinity_de_novo /home/$USER/AiptasiaRNASeq/data/trinity_de_novo
62 cp results/trinity*stats.txt /home/$USER/AiptasiaRNASeq/data/
63
64 echo "Assemblies complete $(date)"
```

Assessing Trinity output

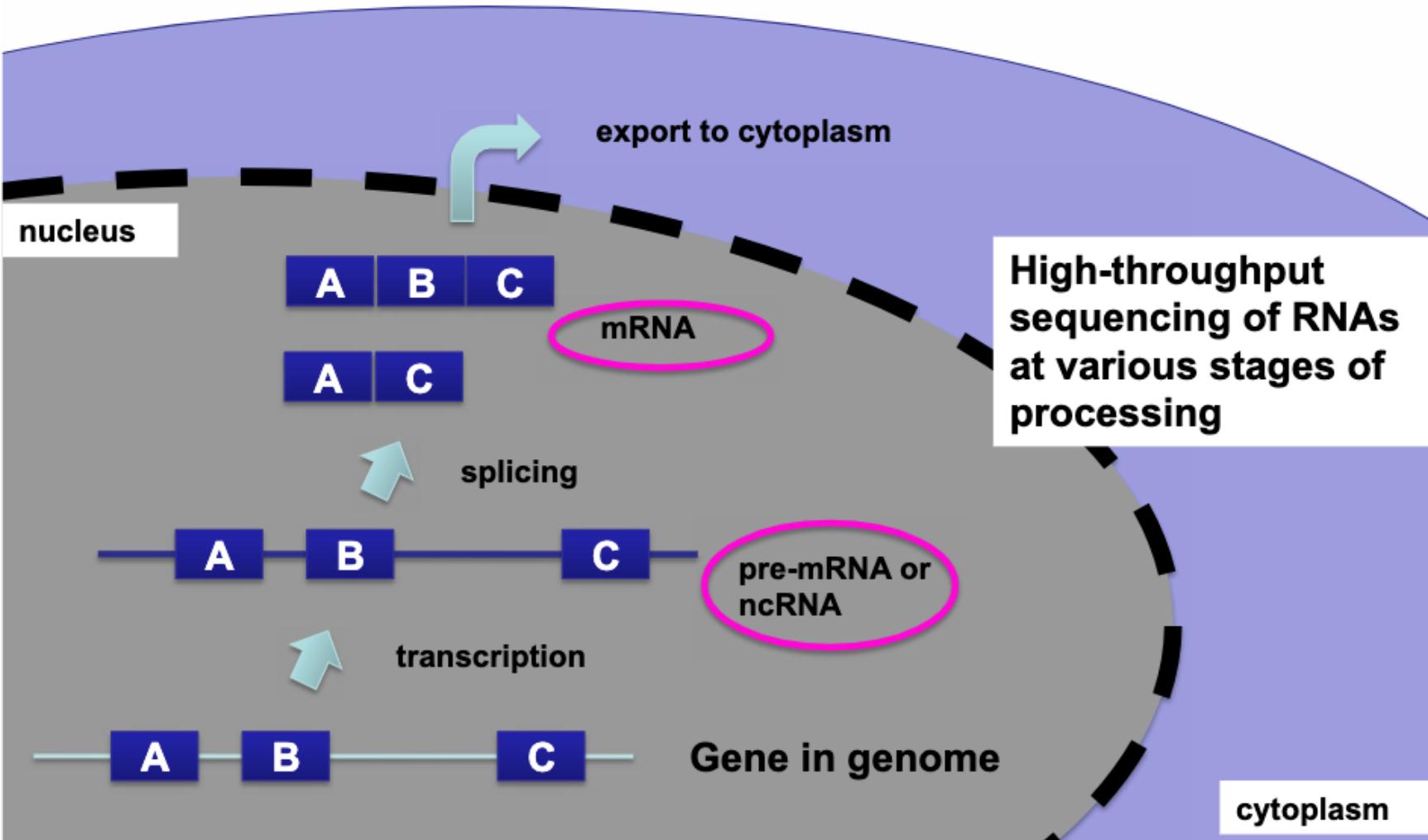
```
1 ##
2 ##
3 ## #####
4 ## ## Counts of transcripts, etc.
5 ## #####
6 ## Total trinity 'genes': 26028
7 ## Total trinity transcripts: 27595
8 ## Percent GC: 38.49
9 ##
10 ## #####
11 ## Stats based on ALL transcript contigs:
12 ## #####
13 ##
14 ## Contig N10: 1917
15 ## Contig N20: 1292
16 ## Contig N30: 952
17 ## Contig N40: 727
18 ## Contig N50: 572
19 ##
20 ## Median contig length: 356
21 ## Average contig: 503.55
22 ## Total assembled bases: 13895363
23 ##
24 ##
25 ## #####
26 ## ## Stats based on ONLY LONGEST ISOFORM per 'GENE':
27 ## #####
28 ##
29 ## Contig N10: 1599
30 ## Contig N20: 1103
31 ## Contig N30: 820
32 ## Contig N40: 639
33 ## Contig N50: 516
34 ##
35 ## Median contig length: 347
36 ## Average contig: 470.82
37 ## Total assembled bases: 12254396
```

```
1 #!/usr/bin/env bash
2 # analyzeTrinity.sh
3 # Usage: bash scripts/analyzeTrinity.sh 1>results/trinity_
4 # guided_stats.txt 2>results/logs/analyzeTrinity.err
5 TrinityStats.pl results/trinity_guided/Trinity-GG.fasta
```

Lecture 7: Outline

- Transcriptome Assembly
 - Assembly v. alignment
 - Reference-guided v. *de novo*
 - Index-based alignment
 - Using Trinity (demo)
- Links in Unix
- Using IGV (round 1)

Review of gene expression



Courtesy of Cole Trapnell. Used with permission.

Slide courtesy Cole Trapnel

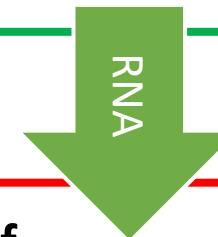
How is RNA-seq different from DNA-seq? What are we looking for?

The message:

'It was the best of times, it was the vorst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of Light, it was the season of Darkness ...'

The features:

worst>vorst



1 belief	2 age
1 best	2 epoch
1 darkness	2 season
1 foolishness	2 times
1 incredulity	8 it
1 light	8 of
1 wisdom	8 the
1 vorst	8 was

"expression"

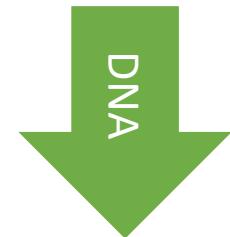
How is RNA-seq different from DNA-seq? What are we looking for?

The features:

worst>vorst

The question:

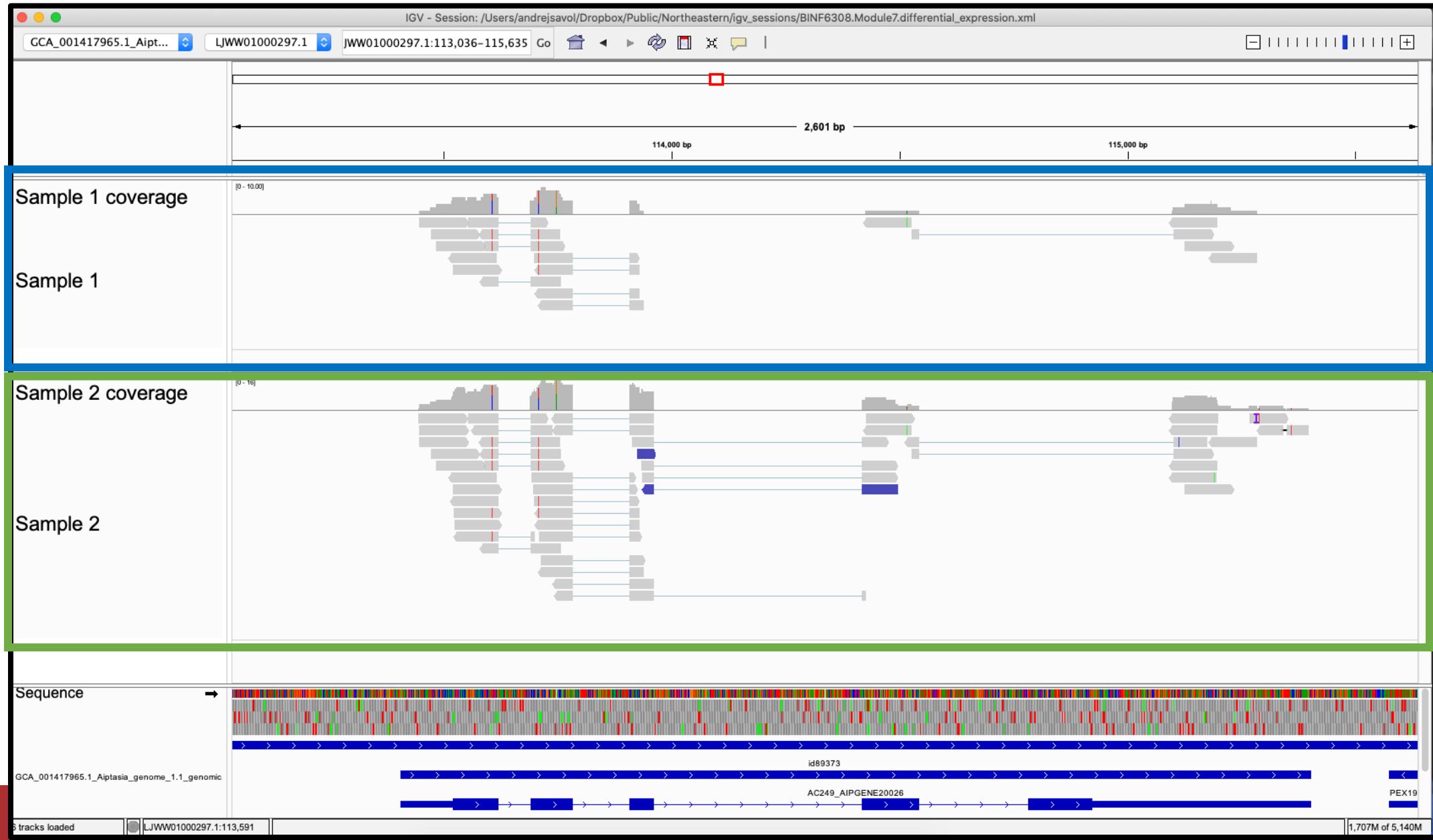
Which words are misspelled?



Which words occur more often than expected?

Calculating gene expression: a visual introduction

RNA-seq data from two samples



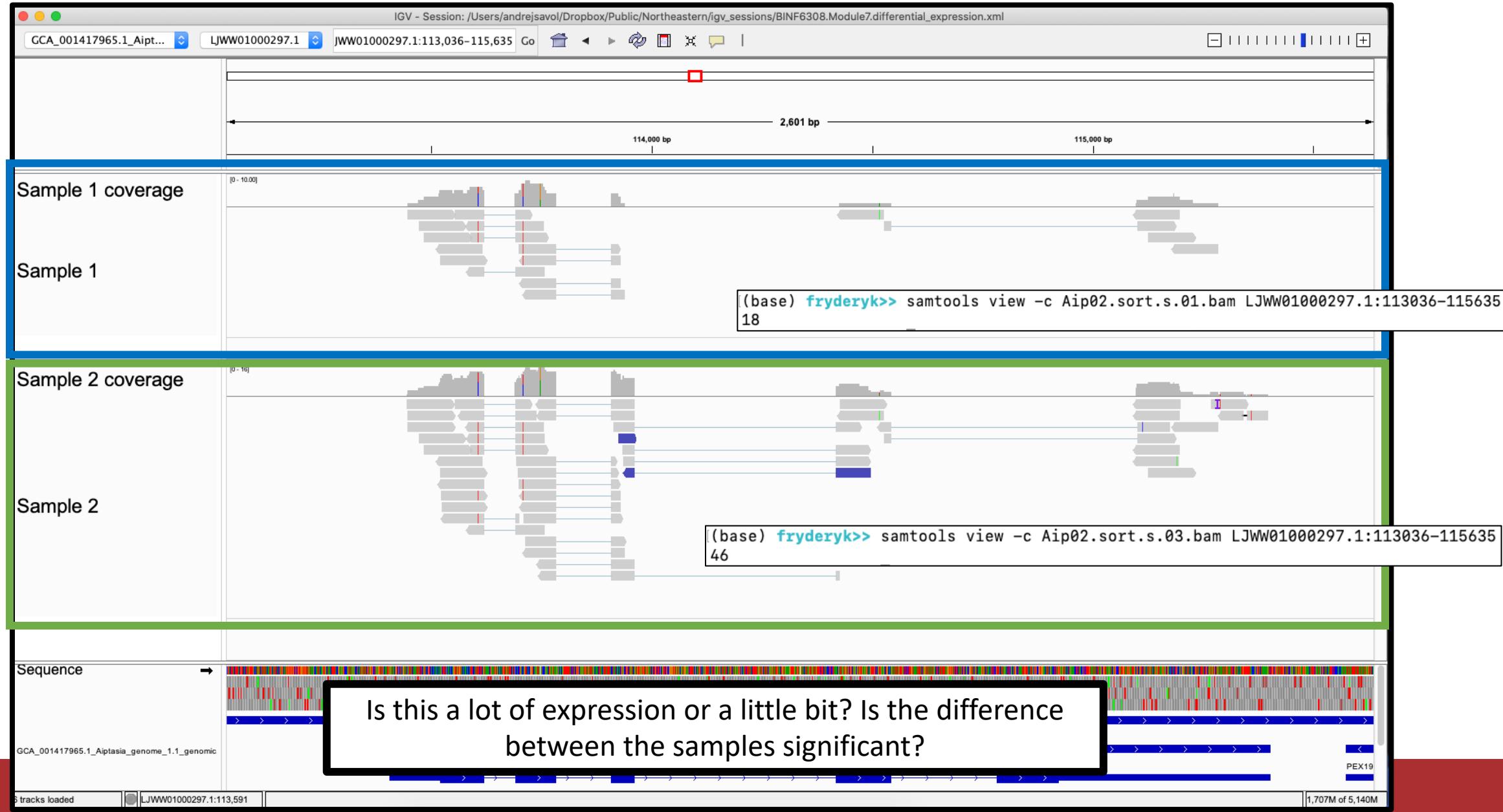
Calculating gene expression: a visual introduction

RNA-seq data from two samples



Calculating gene expression: a visual introduction

RNA-seq data from two samples



Why not just count reads per gene?



Some requirements for a unit of gene expression?

Normalized by gene length (so we can compare between genes)

Normalized by sequencing depth (so we can compare between samples)

$$\text{RPKM of a gene: RPKM} = \frac{\text{Number of reads mapped to gene} \times 10^3 \times 10^6}{\text{Total number of mapped reads} \times \text{Gene length in bp}}$$

$$\text{FPKM of a gene: FPKM} = \frac{\text{Number of fragments mapped to gene} \times 10^3 \times 10^6}{\text{Total number of mapped reads} \times \text{Gene length in bp}}$$

$$\text{TPM of a gene: TPM} = A \times \frac{1}{\sum(A)} \times 10^6 \text{ Where } A = \frac{\text{Total reads mapped to gene} \times 10^3}{\text{Gene length in bp}}$$

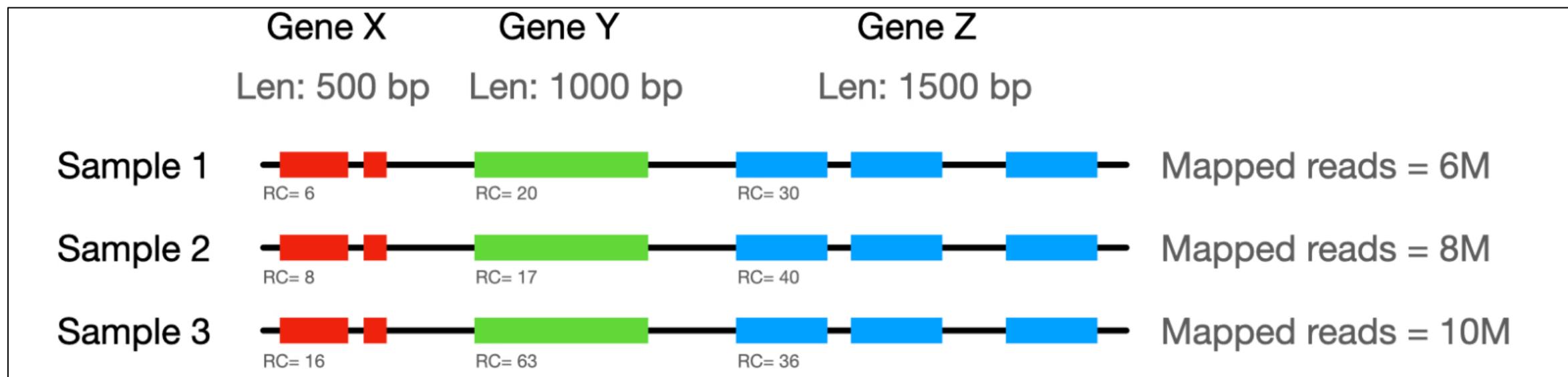
Toy example to compare RPKM and TPM

$$\text{RPKM of a gene: } \text{RPKM} = \frac{\text{Number of reads mapped to gene} \times 10^3 \times 10^6}{\text{Total number of mapped reads} \times \text{Gene length in bp}}$$

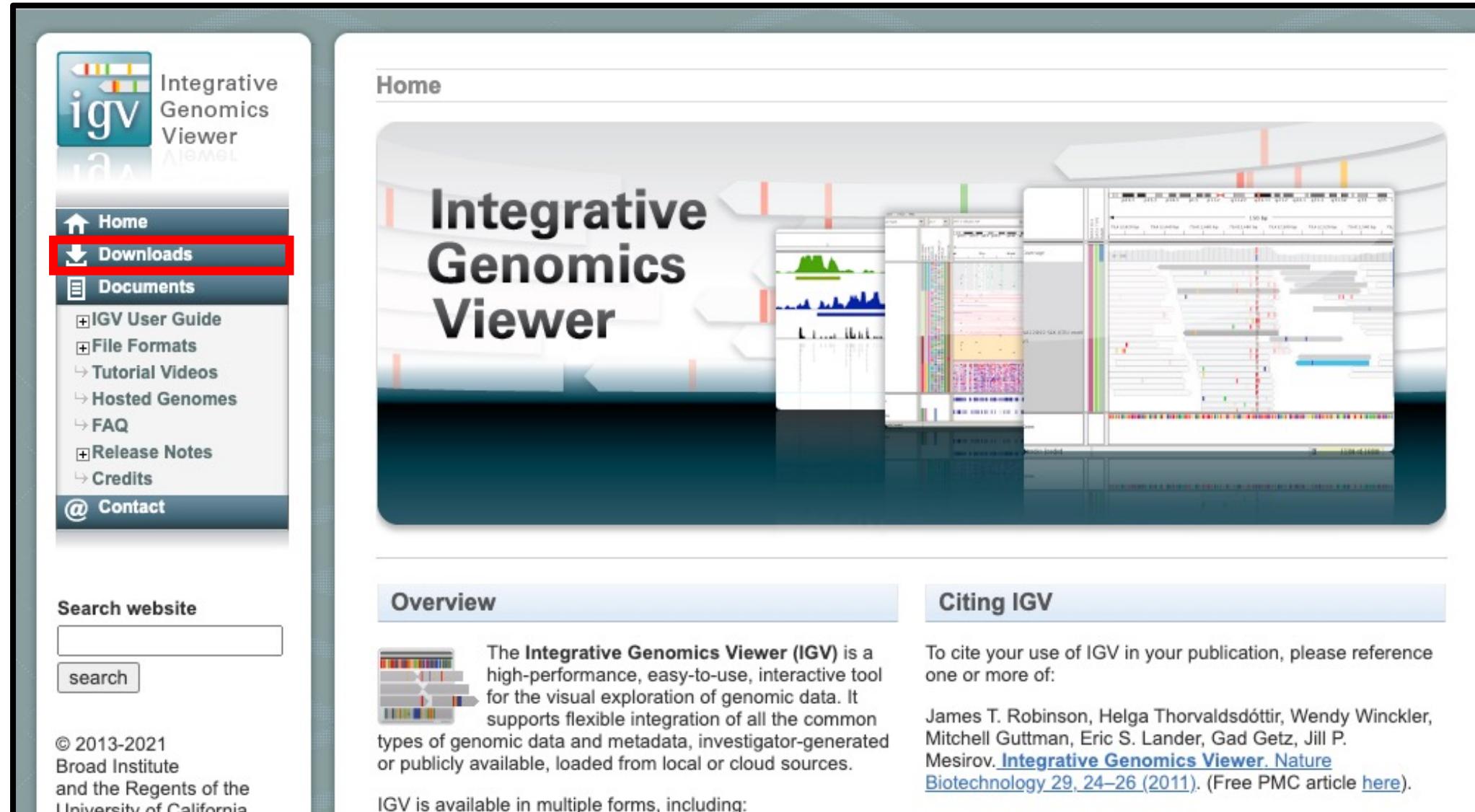
Reads
Total mapped \times Gene length
(in millions) \times In kilobases

$$\text{TPM of a gene: } \text{TPM} = A \times \frac{1}{\sum(A)} \times 10^6 \text{ Where } A = \frac{\text{Total reads mapped to gene} \times 10^3}{\text{Gene length in bp}}$$

A = $\frac{\text{Reads}}{\text{Gene length}} \text{ In kilobases}$



Introduction to IGV (Integrated Genomics Viewer)



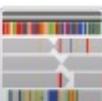
The screenshot shows the official website for the Integrative Genomics Viewer (IGV). On the left, a sidebar contains a logo, a navigation menu with items like Home, Downloads (which is highlighted with a red box), Documents, and various guides and videos, and a search bar. The main content area features a large title "Integrative Genomics Viewer" and a preview of the software's graphical user interface, which displays genomic tracks and data. Below the preview, there are sections for "Overview" and "Citing IGV", along with copyright information for the Broad Institute.

Home

Integrative Genomics Viewer

IGV is available in multiple forms, including:

Overview

 The Integrative Genomics Viewer (IGV) is a high-performance, easy-to-use, interactive tool for the visual exploration of genomic data. It supports flexible integration of all the common types of genomic data and metadata, investigator-generated or publicly available, loaded from local or cloud sources.

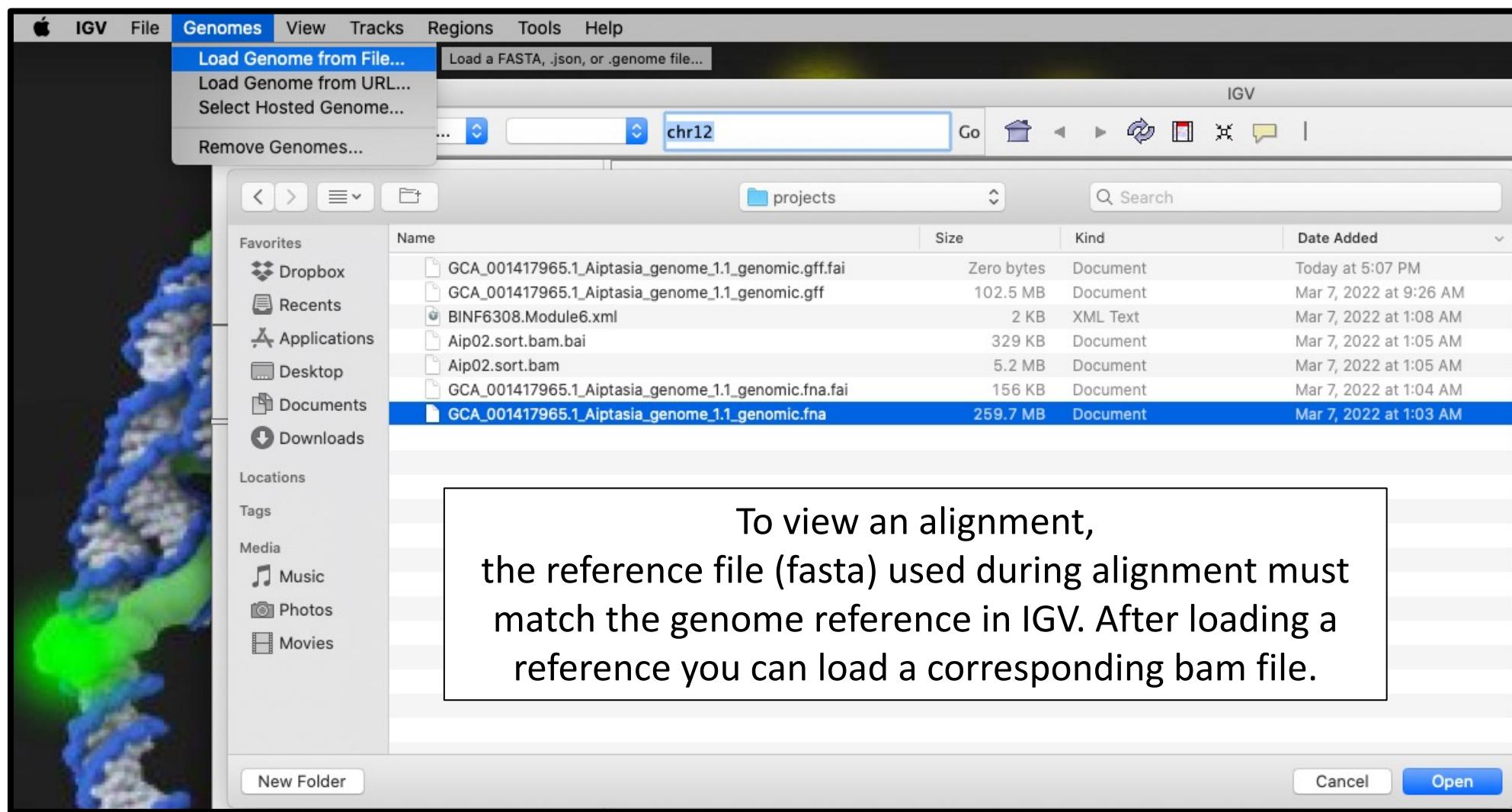
Citing IGV

To cite your use of IGV in your publication, please reference one or more of:

James T. Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, Jill P. Mesirov. [Integrative Genomics Viewer](#). *Nature Biotechnology* 29, 24–26 (2011). (Free PMC article [here](#)).

© 2013-2021
Broad Institute
and the Regents of the
University of California

Generating a *.genome file for IGV



To view an alignment in IGV:

1. Load your reference genome as a fasta file:

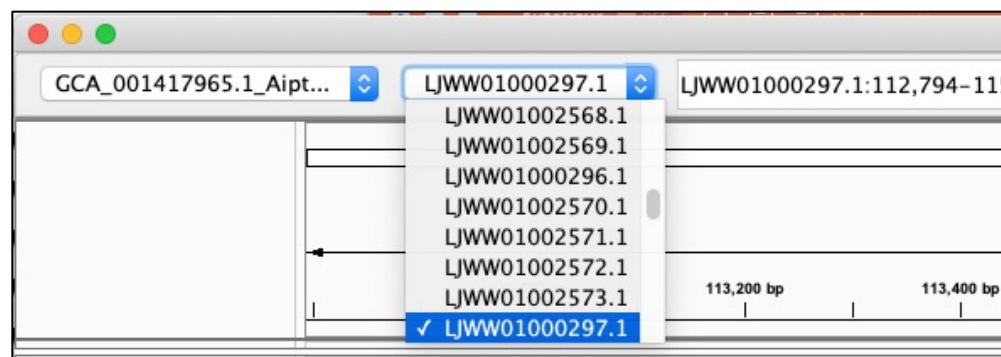
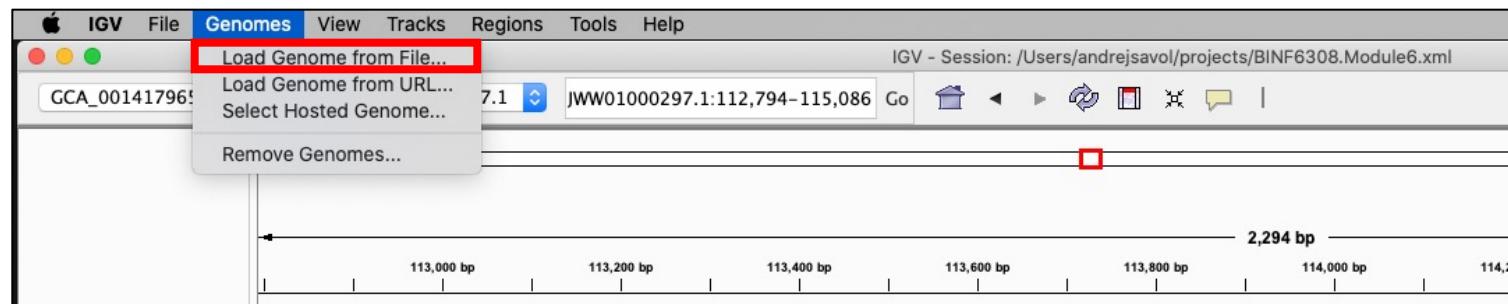
GCA_001417965.1_Aiptasia_genome_1.1_genomic.fna

2. Load your new bam file:

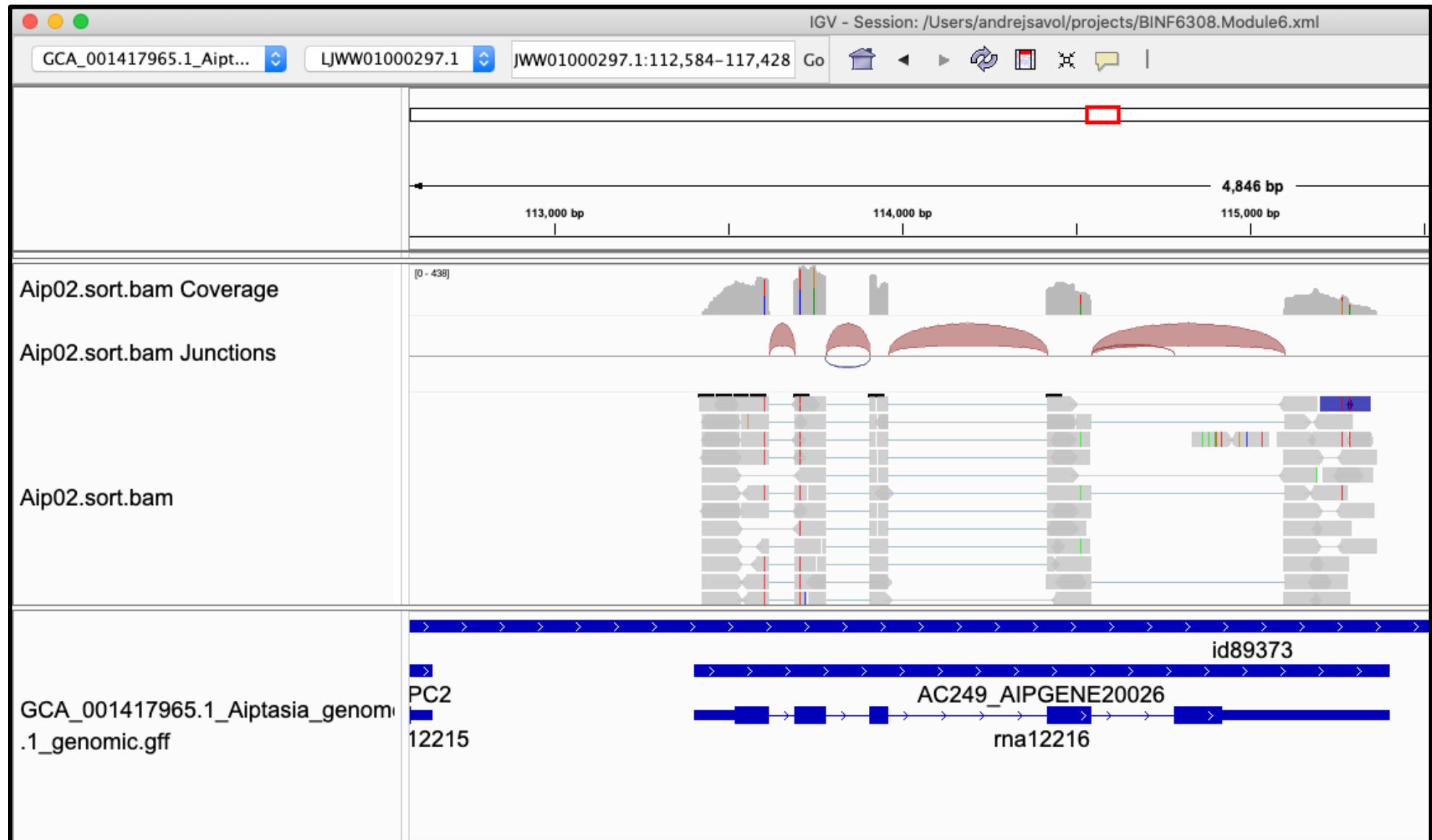
Aip02.sort.bam

3. Navigate to a chromosome that has a decent number of reads (via the dropdown menu):

LJWW01000297.1



To view an alignment in IGV:



Which track is only typical in RNA (versus DNA)?