

# Exam Review

DISCLAIMER: This is a non-exhaustive list

Also, remember the things I told you Id ask again, from exam 1

Lecture 30

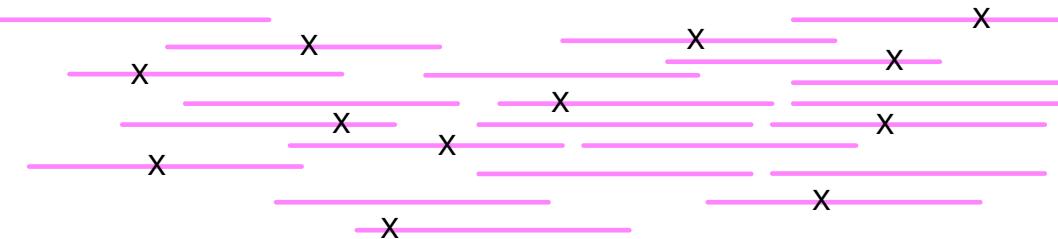
Nov 16, 2016

# Announcements

# Lecture 14

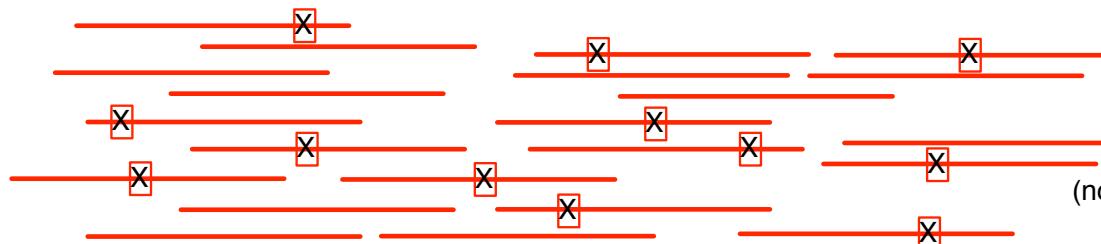
# DIGITAL NORMALIZATION

— - - - - True sequence (unknown)



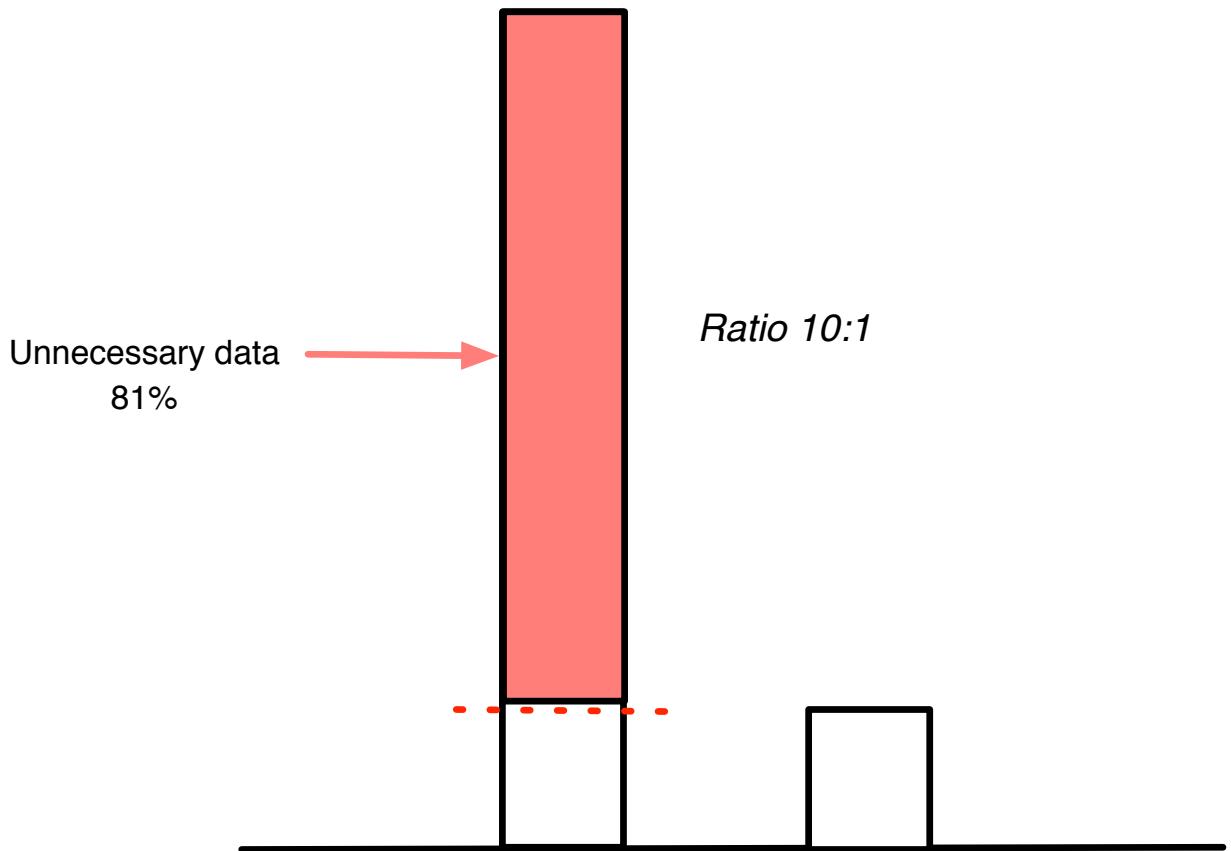
Reads  
(randomly sequenced)

```
for read in dataset:  
    if estimated_coverage(read) < C:  
        accept(read)  
    else:  
        discard(read)
```



Redundant reads  
(not needed for assembly)

# DIGITAL NORMALIZATION



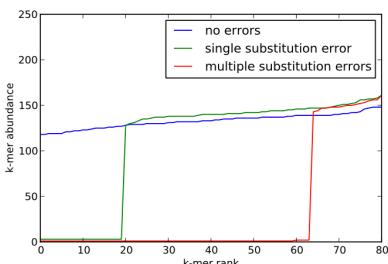
# DIGITAL NORMALIZATION

## Median kmer abundance

0 error: 32,33,34,34,35,36,40

1 error: 1,1,1,32,24,24,40

>1 error: 1,1,1,1,1,32,34



# Lecture 16

# ASSEMBLE A GENOME? GENERAL STRATEGIES

Genome size	Unlimited \$\$	Typical
>10Mb		
10Mb - 100Mb		
> 100 Mb		

# Lecture 18

# ASSEMBLE A GENOME?

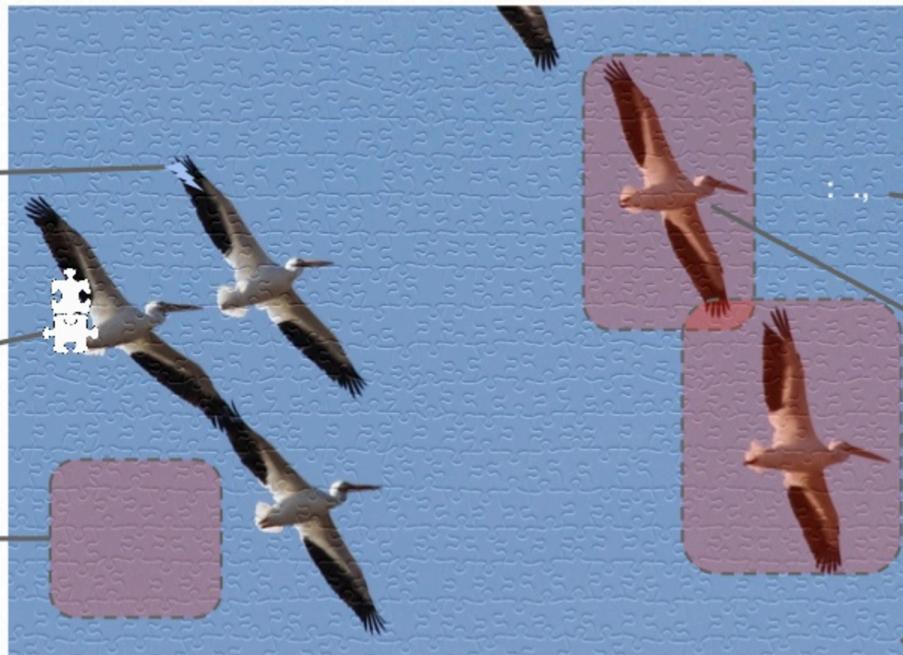
## What makes a jigsaw puzzle hard?

No box

Frayed pieces

Missing pieces

Repetitive regions



Lots of pieces

Dirty pieces

Multiple copies

No corners  
(circular genomes)

# **OTHER ASSEMBLY TERMS**

Unitig

Contig

scaffold

# Lecture 19

# First law of assembly

If a suffix of read A is similar to a prefix of read B...

TCTATATCTCGGCTCTAGG  
| | | | | | |  
TATCTCGACTCTAGGCC

...then A and B might *overlap* in the genome

TCTATATCTCGGCTCTAGG  
GGCGTCTATATCTCGGCTCTAGGCCCTCATT TTTT  
TATCTCGACTCTAGGCC

# Second law of assembly

More coverage leads to more and longer overlaps

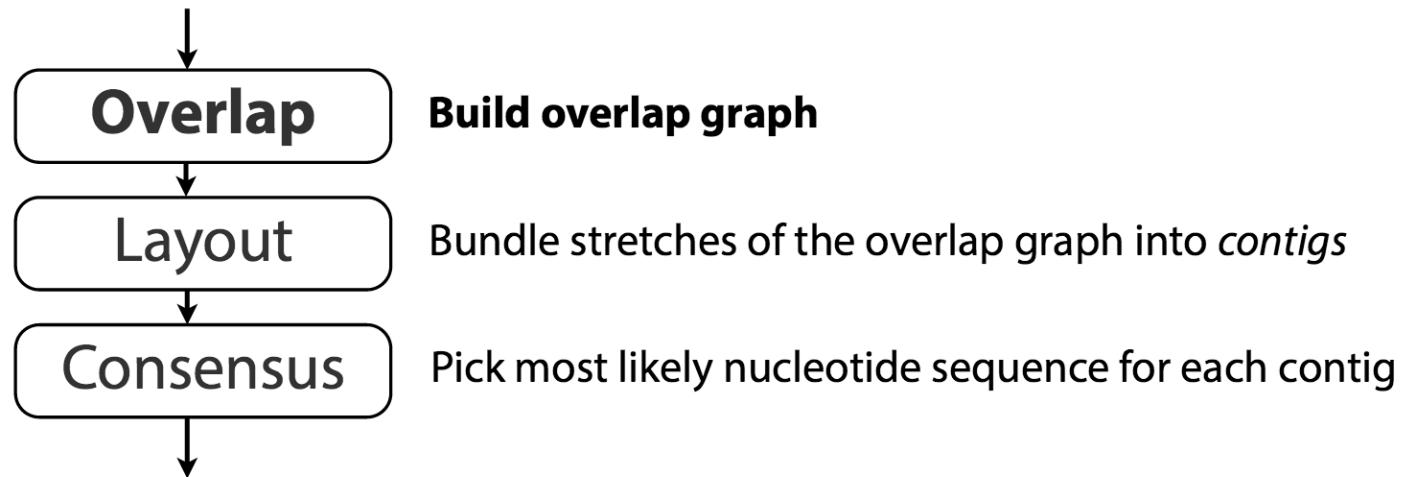
CTAGGCCCTCAATTTT  
CTCGGCTCTAGGCCCTCATT  
TCTATATCTCGGCTCTAGG  
GGCGTCGATATCT  
GGCGTCTATATCTCGGCTCTAGGCCCTCATT  
CTAGGCCCTCAATTTT  
GGCTCTAGGCCCTCATT  
CTCGGCTCTAGGCCCTCATT  
TATCTCGACTCTAGGCCCTCA  
TCTATATCTCGGCTCTAGG  
GGCGTCTATATCTCG  
GGCGTCTATATCT  
less coverage  
more coverage

# **ASSEMBLY**

- OLC Assembly: Characteristics

# ASSEMBLY

- OLC Assembly



# Lecture 20

# **ASSEMBLY – DE BRUIJN**

Hamiltonian Path Problem

Eulerian Path Problem

# Lecture 21

# **Genome Assembly**

**Lecture 21**  
**Oct 21, 2016**

# **ASSEMBLY – DE BRUIJN**

GATTACAGTTCA

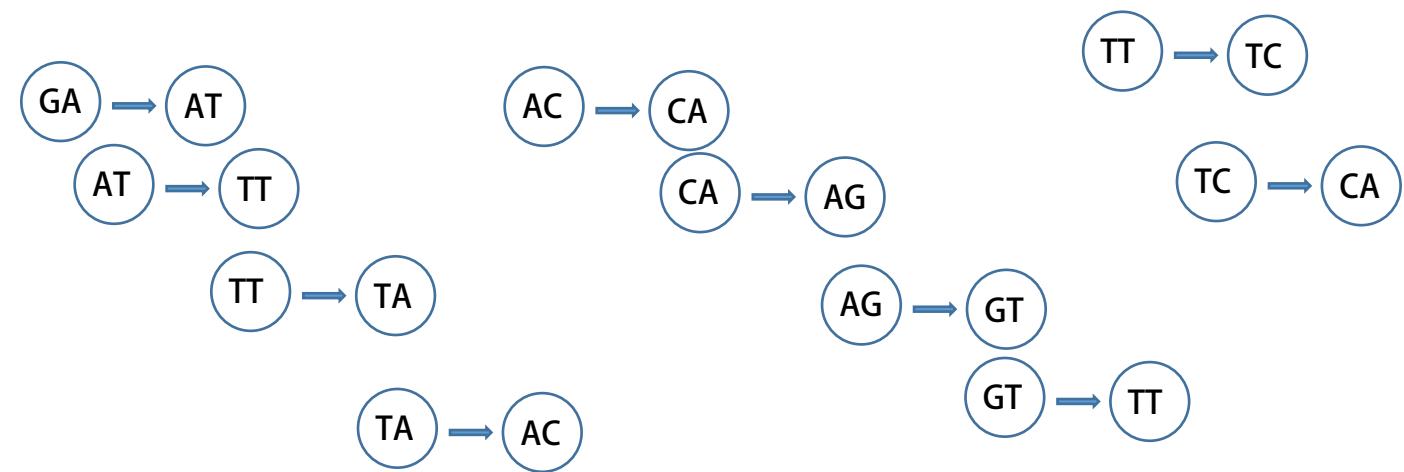
# ASSEMBLY – DE BRUIJN

GATTAC  
GAT  
ATT  
TTA  
TAC

ACAGTTCA  
ACA  
CAG  
AGT  
GTT  
TTC  
TCA

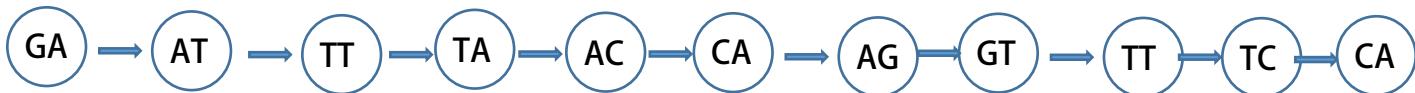
# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



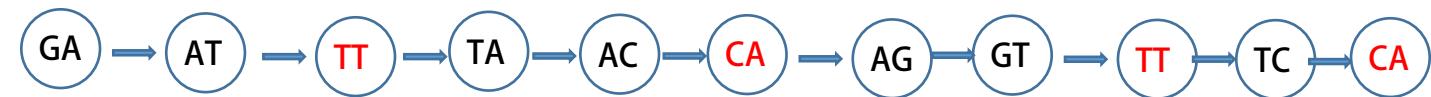
# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



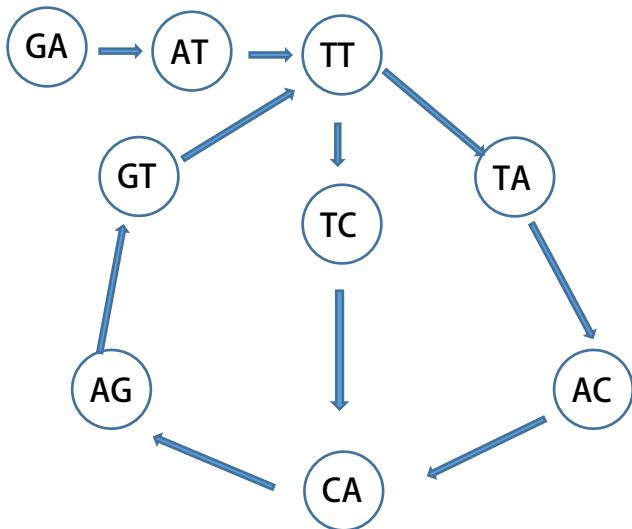
# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



# Lecture 23

# EVALUATING GENOME ASSEMBLIES

# Lecture 24

# **TRANSCRIPTOME ASSEMBLY**

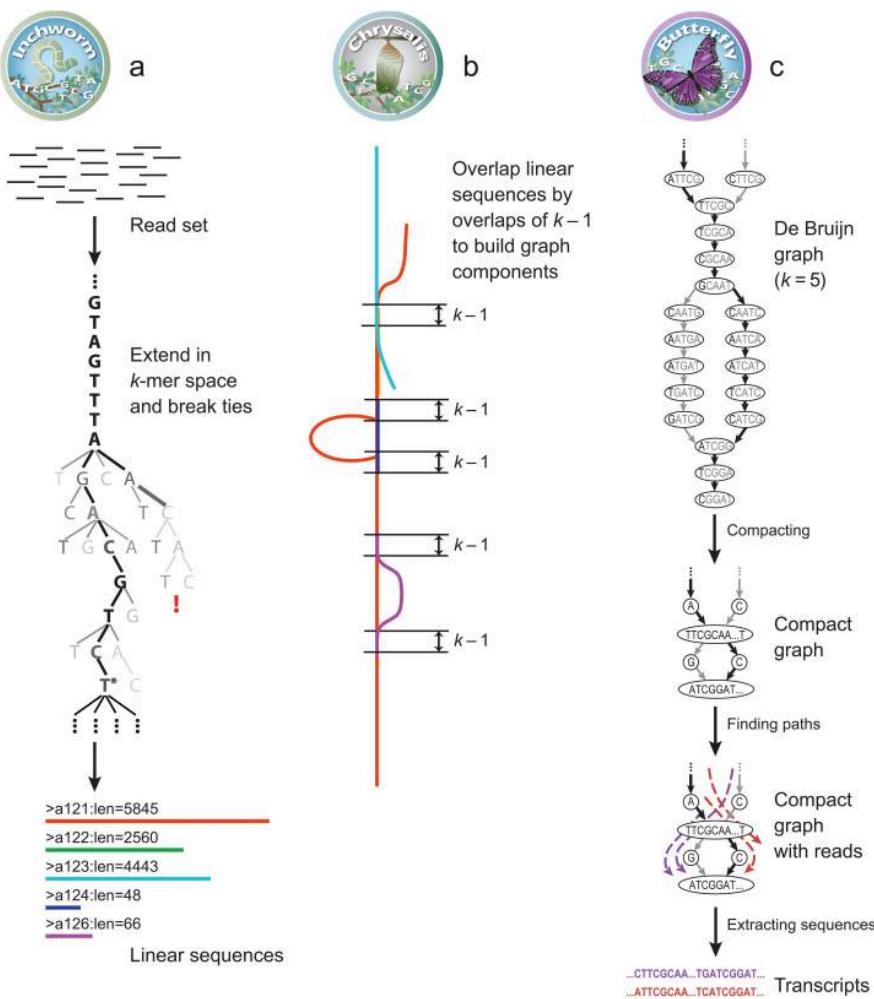
Why easier than genome assembly?

# **TRANSCRIPTOME ASSEMBLY**

Why harder than genome assembly?

# TRANSCRIPTOME ASSEMBLY

Trinity



# Lecture 26

# TRANSCRIPTOME ASSEMBLY – EVALUATION

Error type	Transcripts	Assembly	Read evidence
Family collapse	geneAA geneAB geneAC n=3	n=1	 <p>bases in reads ATCGGAATCGGTT ATAGGTATTGGTA ATAGGGATCGGTG</p>
Chimerism	geneC geneB n=2	n=1	 <p>coverage</p>
Unsupported insertion	n=1	n=1	no reads align to insertion
Incompleteness	n=1	n=1	read pairs align off end of contig
Fragmentation	n=1	n=4	bridging read pairs
Local misassembly	n=1	n=1	read pairs in wrong orientation
Redundancy	n=1	n=3	all reads assign to best contig

# TRANSCRIPTOME ASSEMBLY – EVALUATION – BUSCO

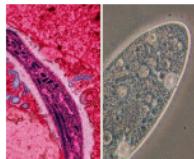
Datasets (Beta versions, updated sets and additional lineages coming soon)



Bacteria sets



Eukaryota sets



Protists sets



Metazoa sets



Fungi sets



Plants set

[Download all datasets](#)

Image credits

# Lecture 27

# How to assess “abundance”

RPKM — Reads per kilobase per million mapped reads

FPKM — Fragments per kilobase per million mapped reads

Don't use these measures, TPM measures the  
“same thing”, but in a better way.

TPM — Transcripts per million

Useful for visualization / assessment etc.

(Estimated) Number of Reads

These are what are used (after normalization)  
for differential expression. Why can't we use TPM?

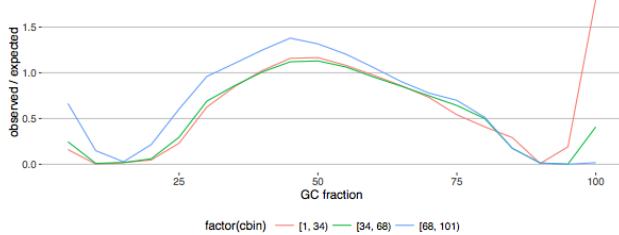
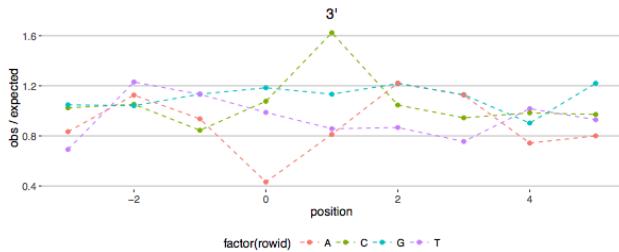
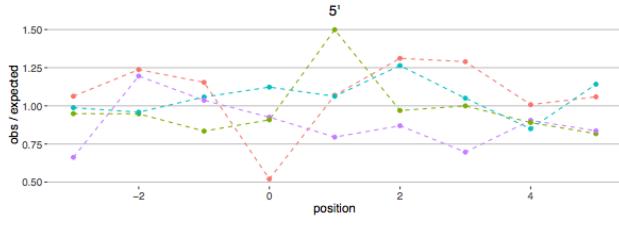
# Biases abound in RNA-seq data

Biases in prep & sequencing can have a significant effect on the fragments we see.

Fragment gc-bias<sup>1</sup>—  
The GC-content of the fragment affects the likelihood of sequencing

Sequence-specific bias<sup>2</sup>—  
sequences surrounding fragment affect the likelihood of sequencing

Positional bias<sup>2</sup>—  
fragments sequenced non-uniformly across the body of a transcript



<sup>1</sup>:Love, Michael I., John B. Hogenesch, and Rafael A. Irizarry. "Modeling of RNA-seq fragment sequence bias reduces systematic errors in transcript abundance estimation." bioRxiv (2015): 025767.

<sup>2</sup>:Roberts, Adam, et al. "Improving RNA-Seq expression estimates by correcting for fragment bias." Genome biology 12.3 (2011): 1.