



# Assembly

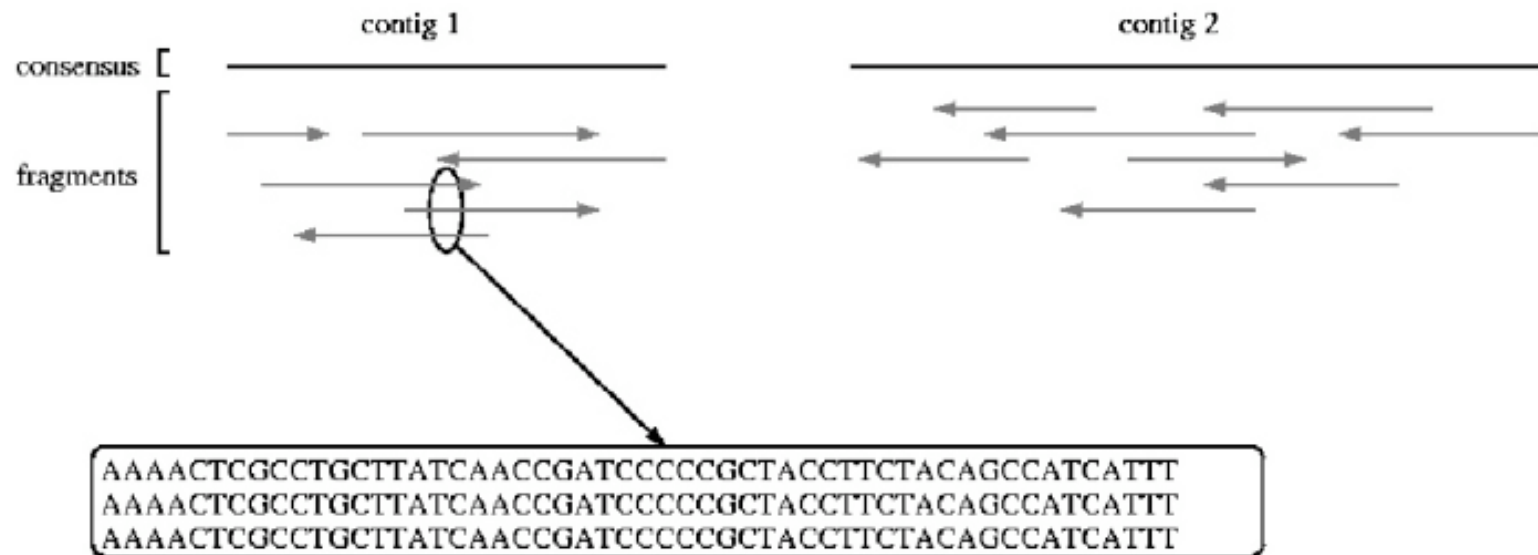


# Outline

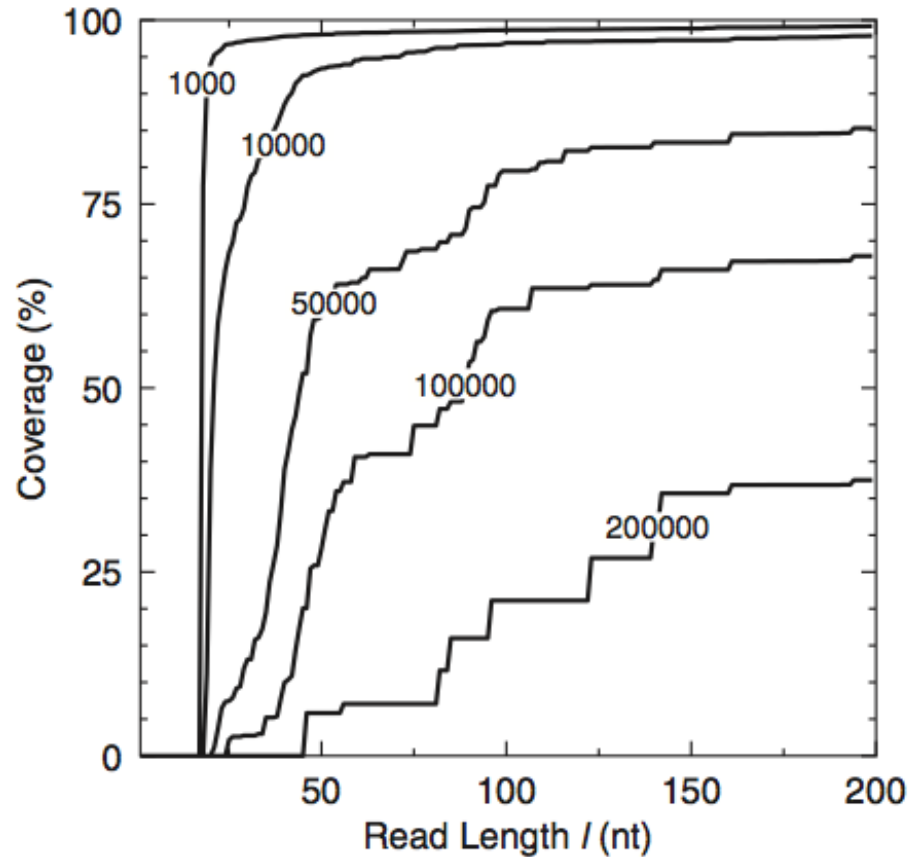
- Random sequencing and coverage
- Assembling by unique overlaps
  - Clustering
  - Memory use
- De Bruijn k-mer graphs
- Heuristics
- Challenges with assembly
  - More data not necessarily better

# Assembly

- Short-read assembly is problematic
- Relies on very deep coverage, ruthless read trimming, paired ends.



# Short read lengths are hard.



**Figure 3.** Percentage of the *E.coli* genome covered by contigs greater than a threshold length as a function of read length.

Whiteford et al., Nuc.Acid Res, 2005



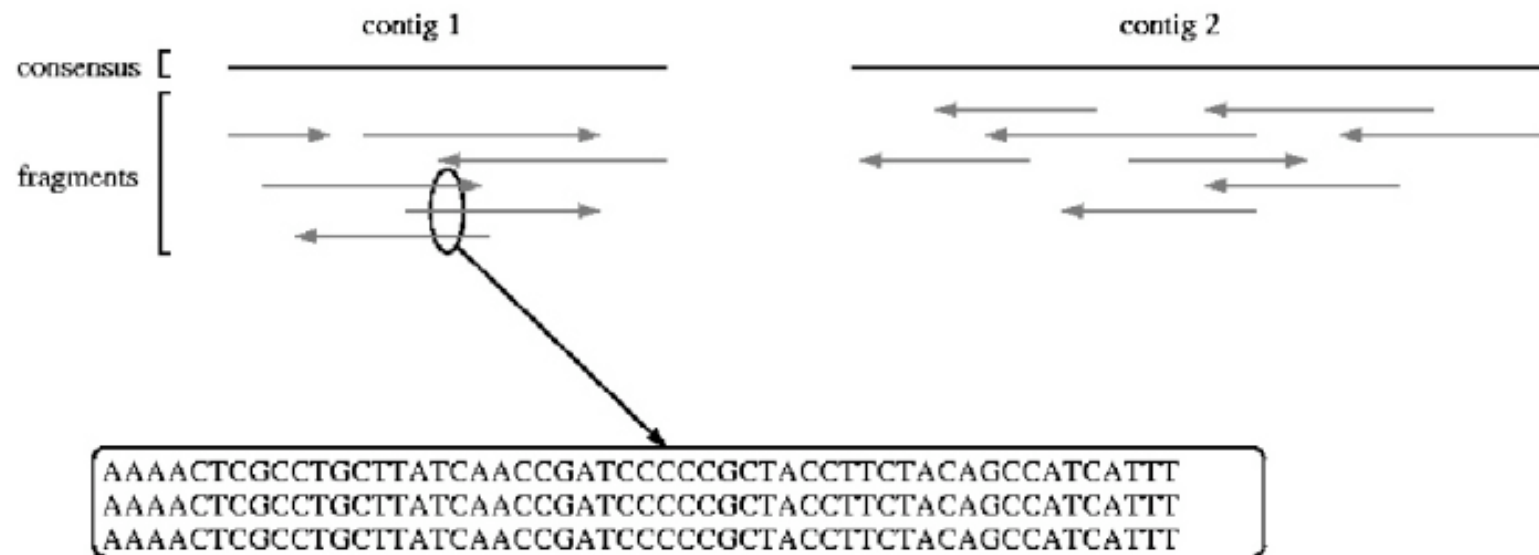
## Two main challenges for *de novo* sequencing, apart from read length

- Repeats.
- Low coverage.

Both introduce breaks in the construction of contigs.

# Repeats

- Overlaps don't place sequences uniquely when there are repeats present.

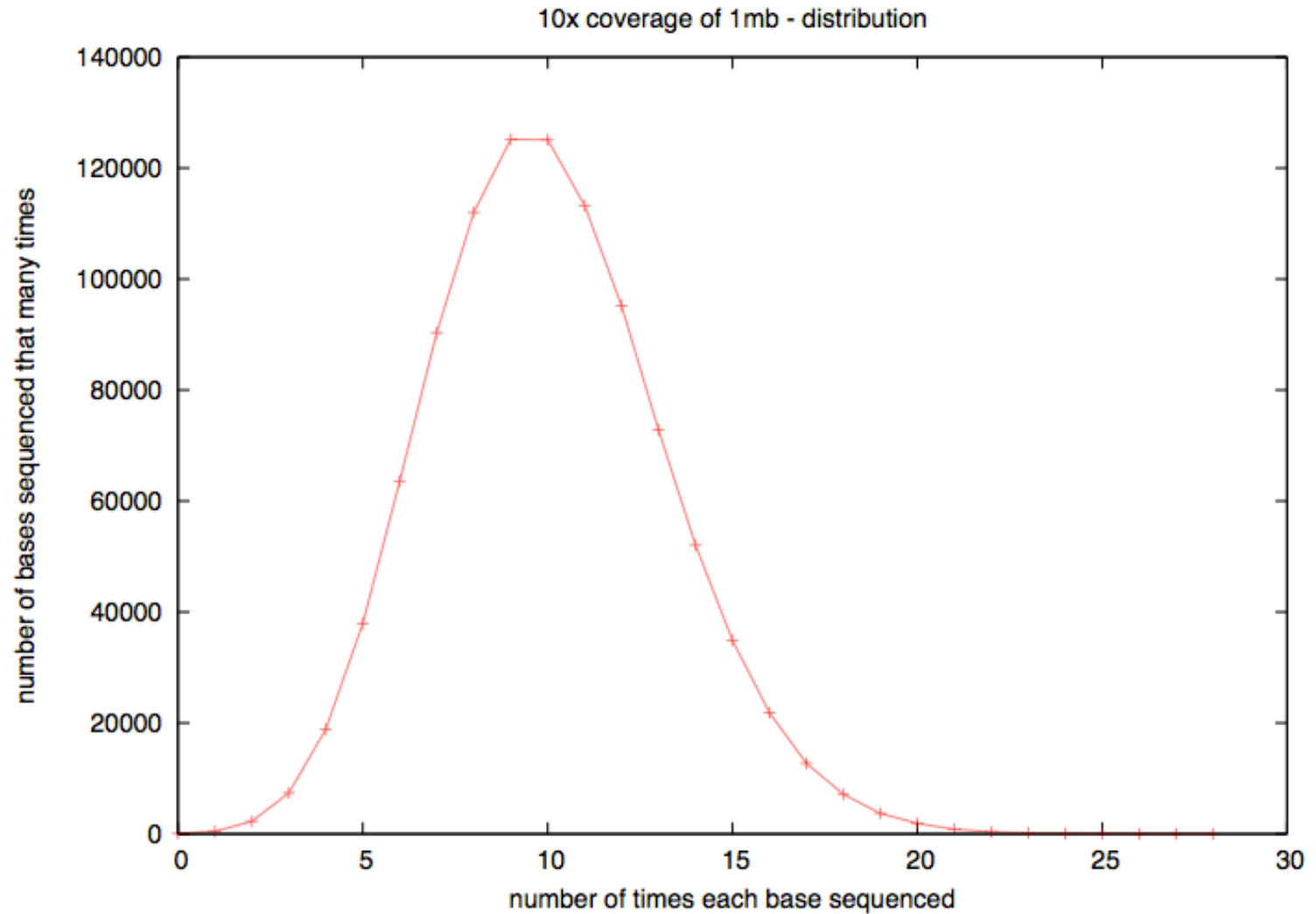




# Coverage

- “1x” doesn’t mean every DNA sequence is read once.
- It means that, if sampling were *systematic*, it would be.
- Sampling isn’t *systematic*, it’s random!

# Coverage varies widely for low x







# Two basic assembly approaches

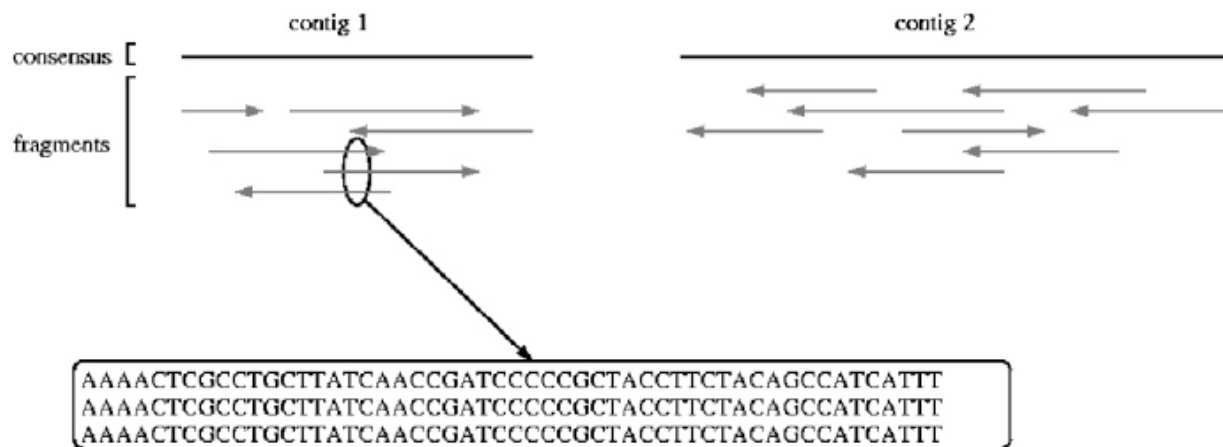
- Overlap/layout/consensus
- De Bruijn k-mer graphs

The former is used for long reads, esp all Sanger-based assemblies. The latter is used because of memory efficiency.

# Overlap/layout/consensus

Essentially,

1. Calculate all overlaps
2. Cluster based on overlap.
3. Do a multiple sequence alignment



# K-mers

Essentially, break reads (of any length) down into multiple overlapping words of fixed length  $k$ .

ATGGACCAGATGACAC (k=12) =>

ATGGACCAGATG  
TGGACCAGATGA  
GGACCAGATGAC  
GACCAGATGACA  
ACCAGATGACAC

# K-mers

**Table 1A.** Mean number of false placements of *K*-mers on the genome

<i>K</i>	<i>Escherichia coli</i>	<i>Saccharomyces cerevisiae</i>	<i>Arabidopsis thaliana</i>	<i>Homo sapiens</i>
200	0.063	0.26	0.053	0.18
160	0.068	0.31	0.064	0.49
120	0.074	0.39	0.086	1.7
80	0.082	0.49	0.15	7.2
60	0.088	0.58	0.27	18
50	0.091	0.63	0.39	32
40	0.095	0.69	0.65	78
30	0.11	0.77	1.5	330
20	0.15	1.0	5.7	2100
10	18	63.8	880	40,000

# K-mers

**Table 1B.** Fraction of K-mers having a unique placement on the genome

<i>K</i>	<i>E. coli</i> (%)	<i>S. cerevisiae</i> (%)	<i>A. thaliana</i> (%)	<i>H. sapiens</i> (%)
200	98.5	95.9	97.4	97.6
160	98.3	95.6	97.1	97.2
120	98.2	95.2	96.6	96.6
80	98.0	94.7	95.4	95.2
60	97.8	94.4	94.4	93.1
50	97.7	94.2	93.4	91.2
40	97.6	93.9	92.2	88.3
30	97.4	93.5	90.4	83.4
20	97.0	92.9	86.5	71.8
10	0.0	0.0	0.0	0.0

# Big genomes are problematic

Species	Ploidy	Genome size (kb)	Reference N50 (kb)	Component N50 (kb)	Edge N50 (kb)	Ambiguities per megabase	Coverage (%)	Coverage by perfect edges $\geq 10$ kb (%)
<i>C. jejuni</i>	1	1800	1800	1800	1800	0.0	100.0	100.0
<i>E. coli</i>	1	4600	4600	4600	4600	0.0	100.0	100.0
<i>B. thailandensis</i>	1	6700	3800	1800	890	2.7	99.8	99.5
<i>E. gossypii</i>	1	8700	1500	1500	890	2.6	100.0	99.9
<i>S. cerevisiae</i>	1	12,000	920	810	290	28.7	98.7	94.9
<i>S. pombe</i>	1	13,000	4500	1400	500	19.1	98.8	97.5
<i>P. stipitis</i>	1	15,000	1800	900	700	8.6	97.9	96.3
<i>C. neoformans</i>	1	19,000	1400	810	770	4.5	96.4	93.4
<i>Y. lipolytica</i>	1	21,000	3600	2200	290	6.2	99.1	98.6
<i>Neurospora crassa</i>	1	39,000	660	640	90	17.4	97.0	92.5
<i>H. sapiens</i> region	2	10,000	10,000	490	2	68.2	97.3	0.2

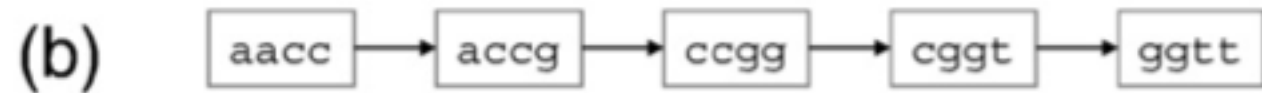
# K-mer graphs

(a) aaccgg



# K-mer graphs - overlaps

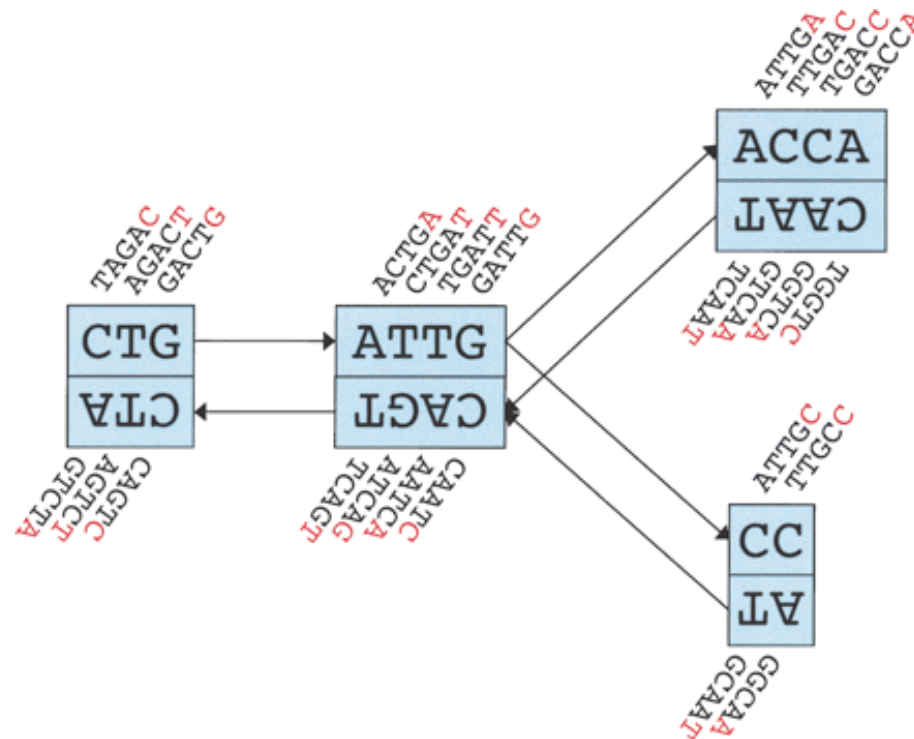
(a) aaccgg  
ccggtt





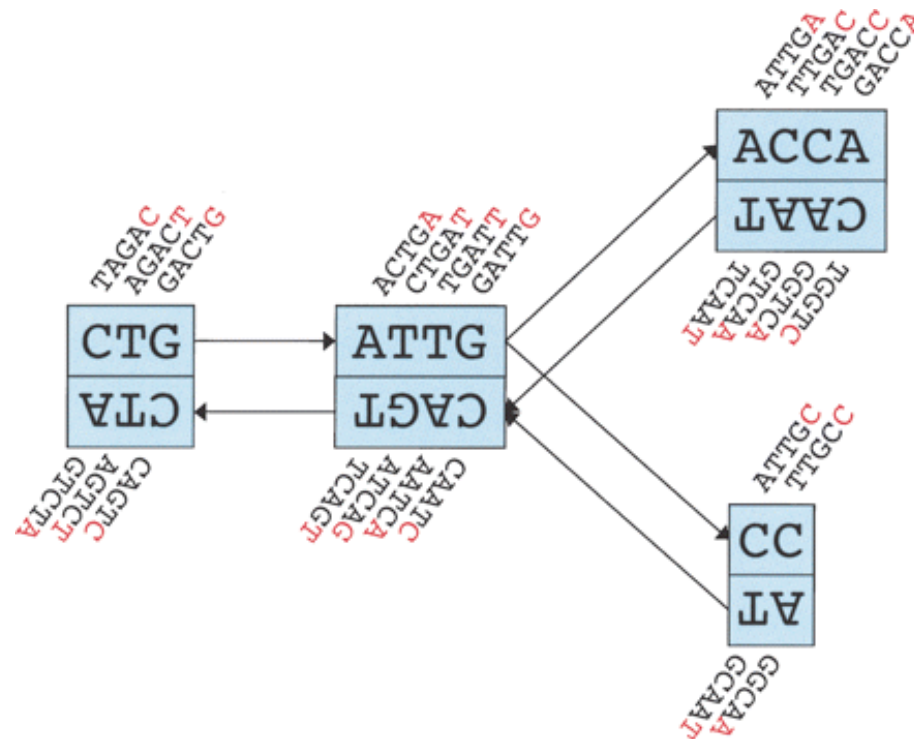
# K-mer graphs - vlvet

Each rectangle represents k-mer from overlapping reads (at angles).



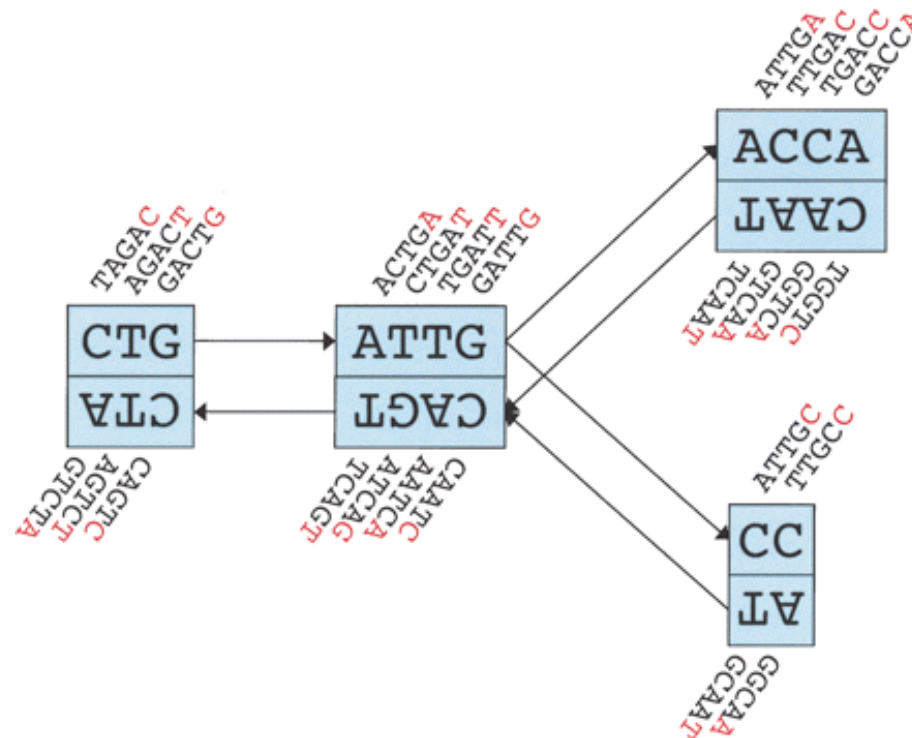
# K-mer graphs

Connectivity shows *possible* contiguous sequences.

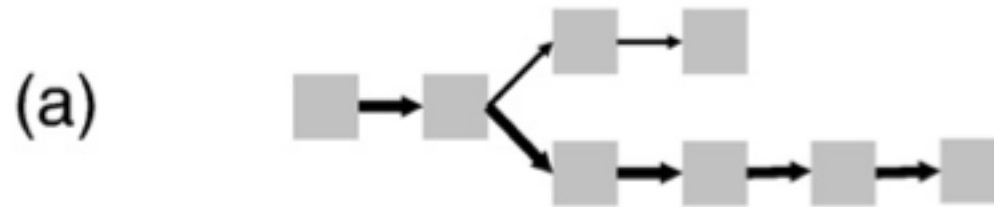


# K-mer graphs

Goal is to construct a *single* path through this graph.



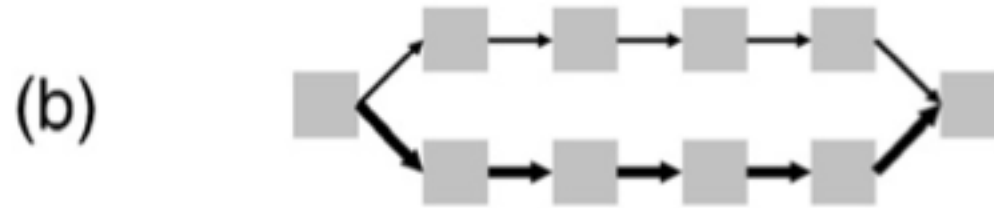
# K-mer graph complexity - spur



(Short) dead-end in graph.

Can be caused by error at the end of some overlapping reads

# K-mer graph complexity - bubble



Multiple parallel paths that diverge and join.

Caused by sequencing error and true polymorphism / polyploidy in sample.

# K-mer graph complexity – “frayed rope”



Converging, then diverging paths.

Caused by repetitive sequences.

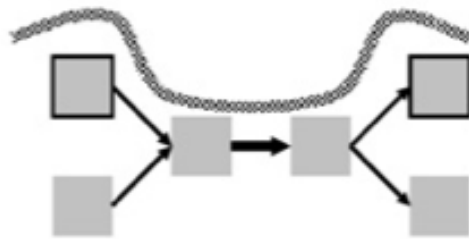


# Resolving graph complexity

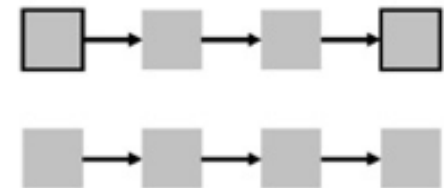
- Primarily heuristic approaches.
- Detecting complex graph structures can generally not be done efficiently.
- Much of the divergence in functionality of new assemblers comes from this.
- Three examples:

# Read threading

(before)



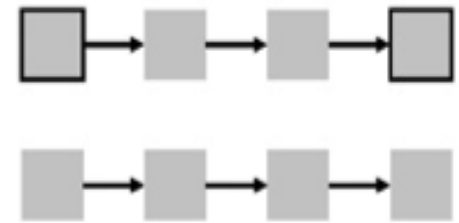
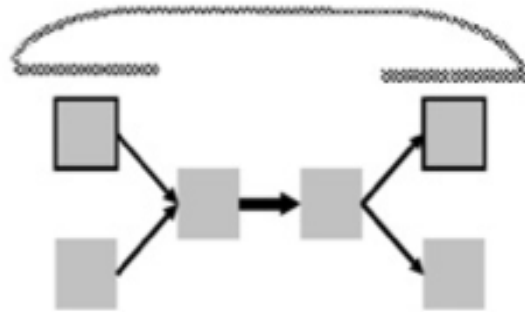
(after)



Single read spans k-mer graph => extract the single-read path.

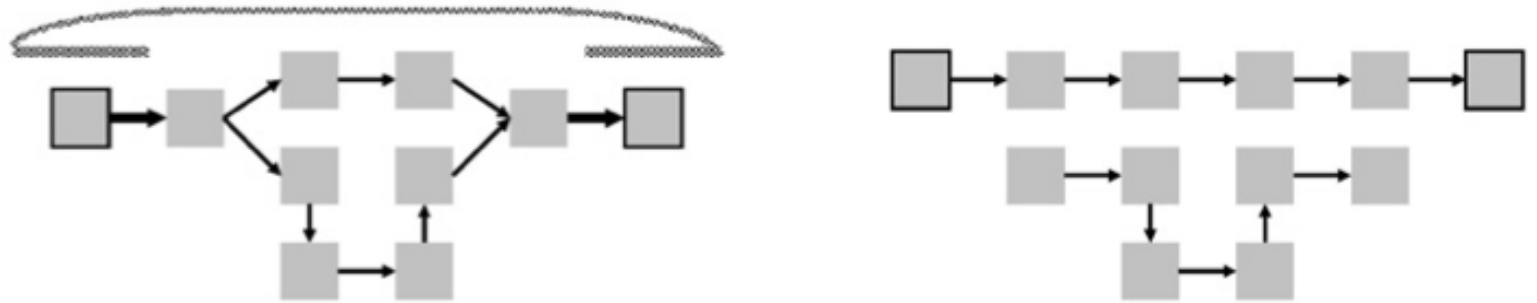


# Mate threading



Resolve “frayed-rope” pattern caused by repeats, by separating paths based on mate-pair reads.

# Path following

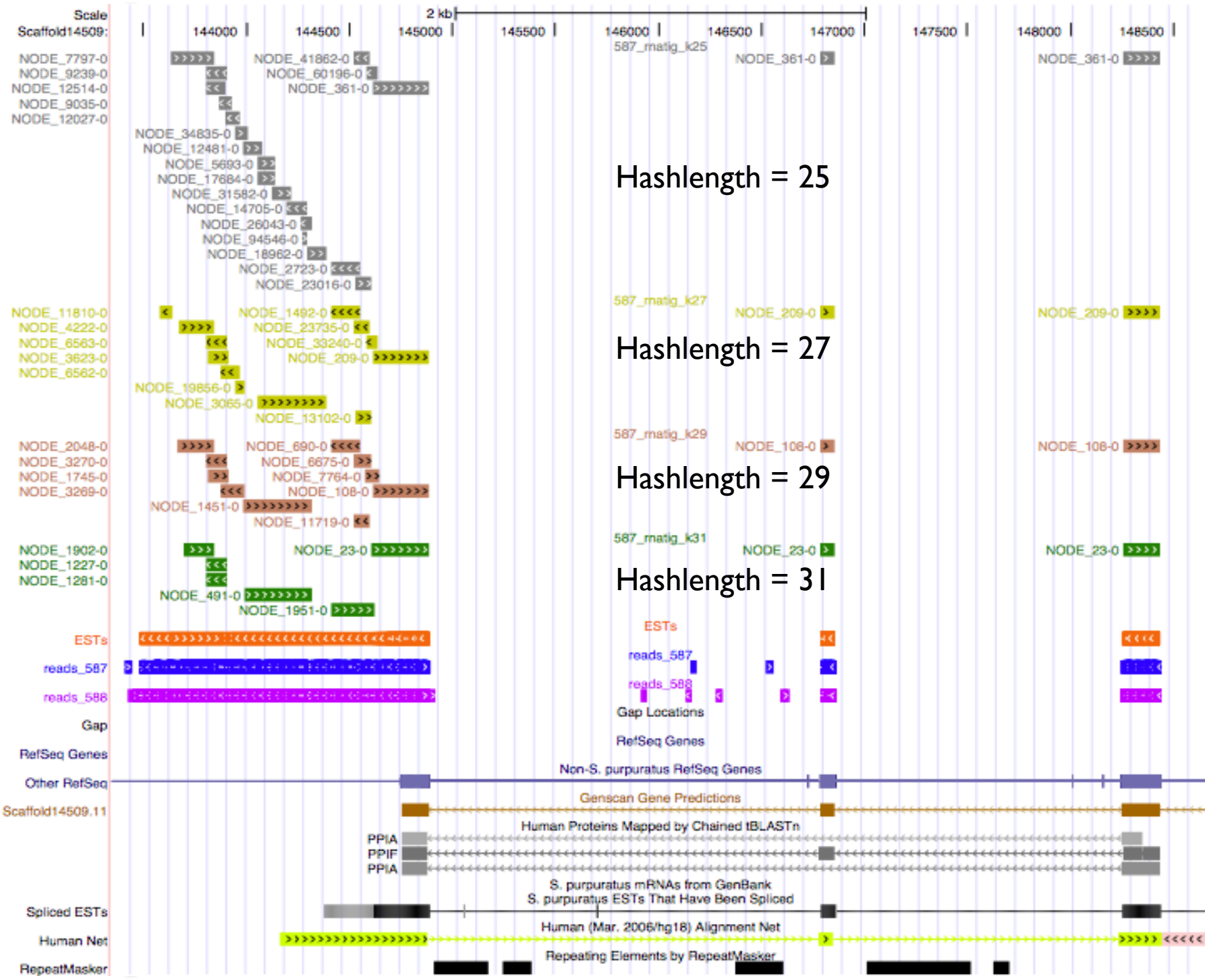


Reject inconsistent paths based on mate-pair reads and insert size.



# More assembly problems

- RNAseq has variation in copy number; naïve assemblers can treat this as repetitive and eliminate it.
- Assembly requires gobs of memory (4 lanes, 60m reads => ~ 150gb RAM)
- How to evaluate?
  - Lamprey RNAseq: 25% of reads map to assembled contigs.



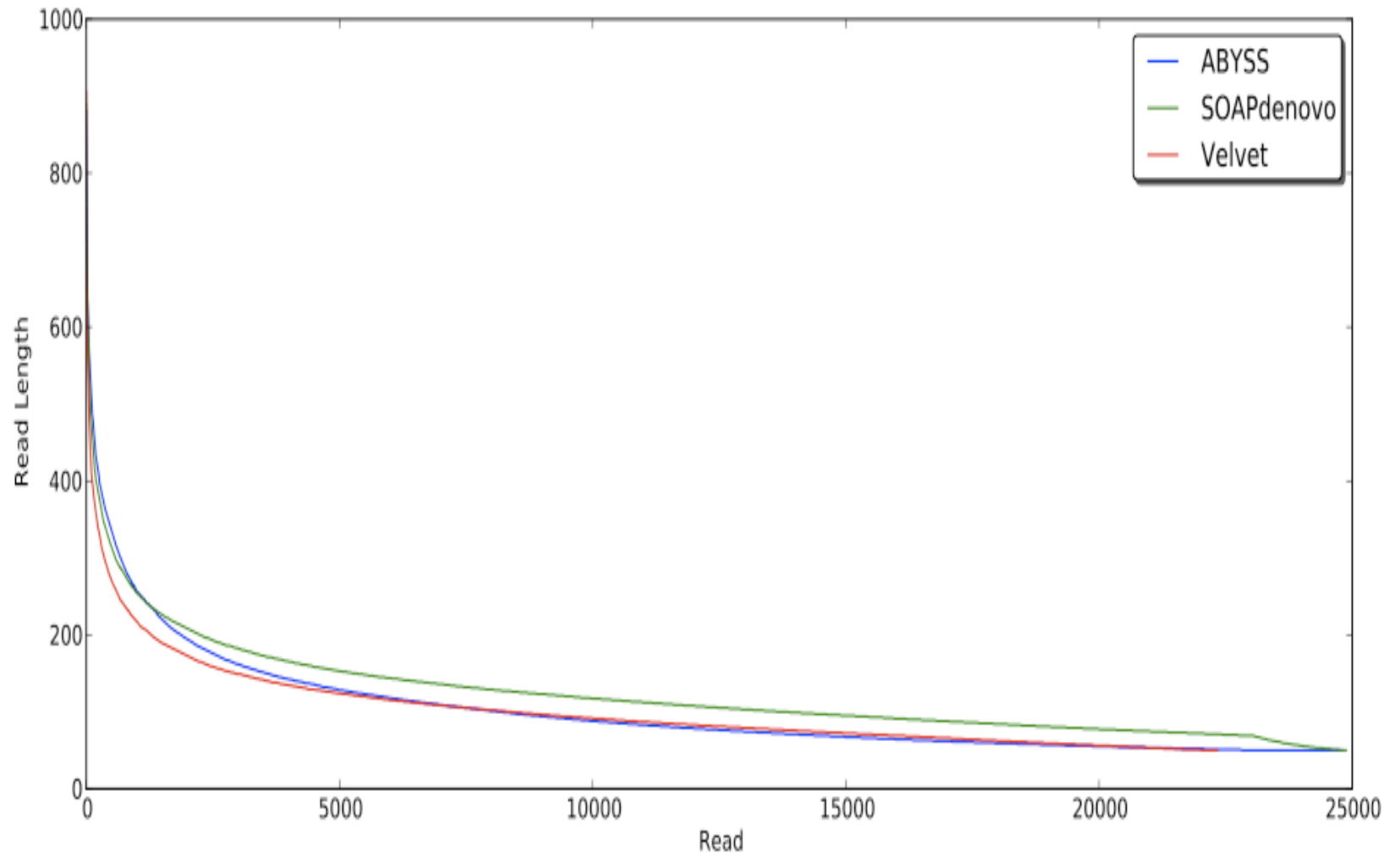
Hashlength = 25

Hashlength = 27

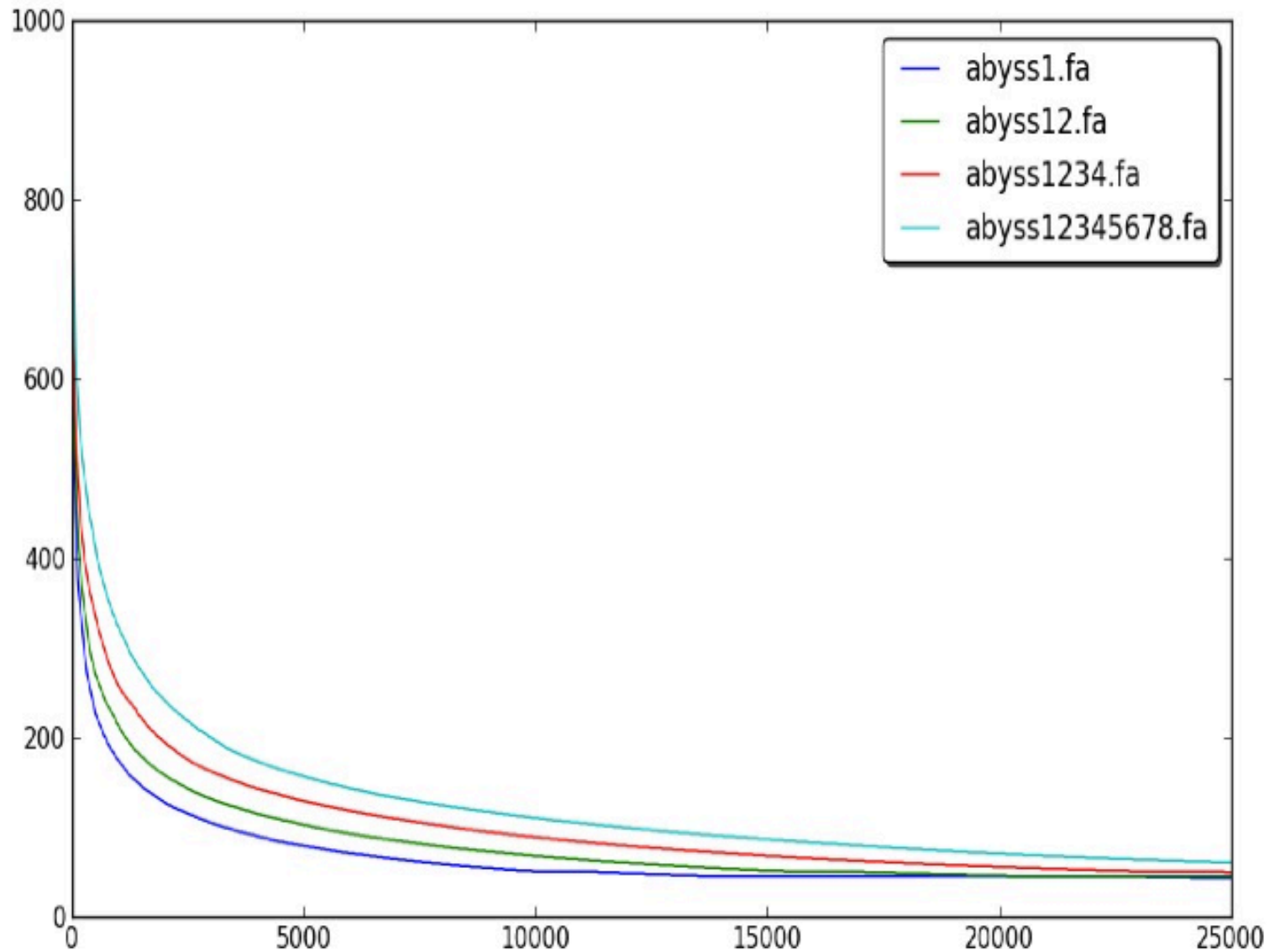
Hashlength = 29

Hashlength = 31

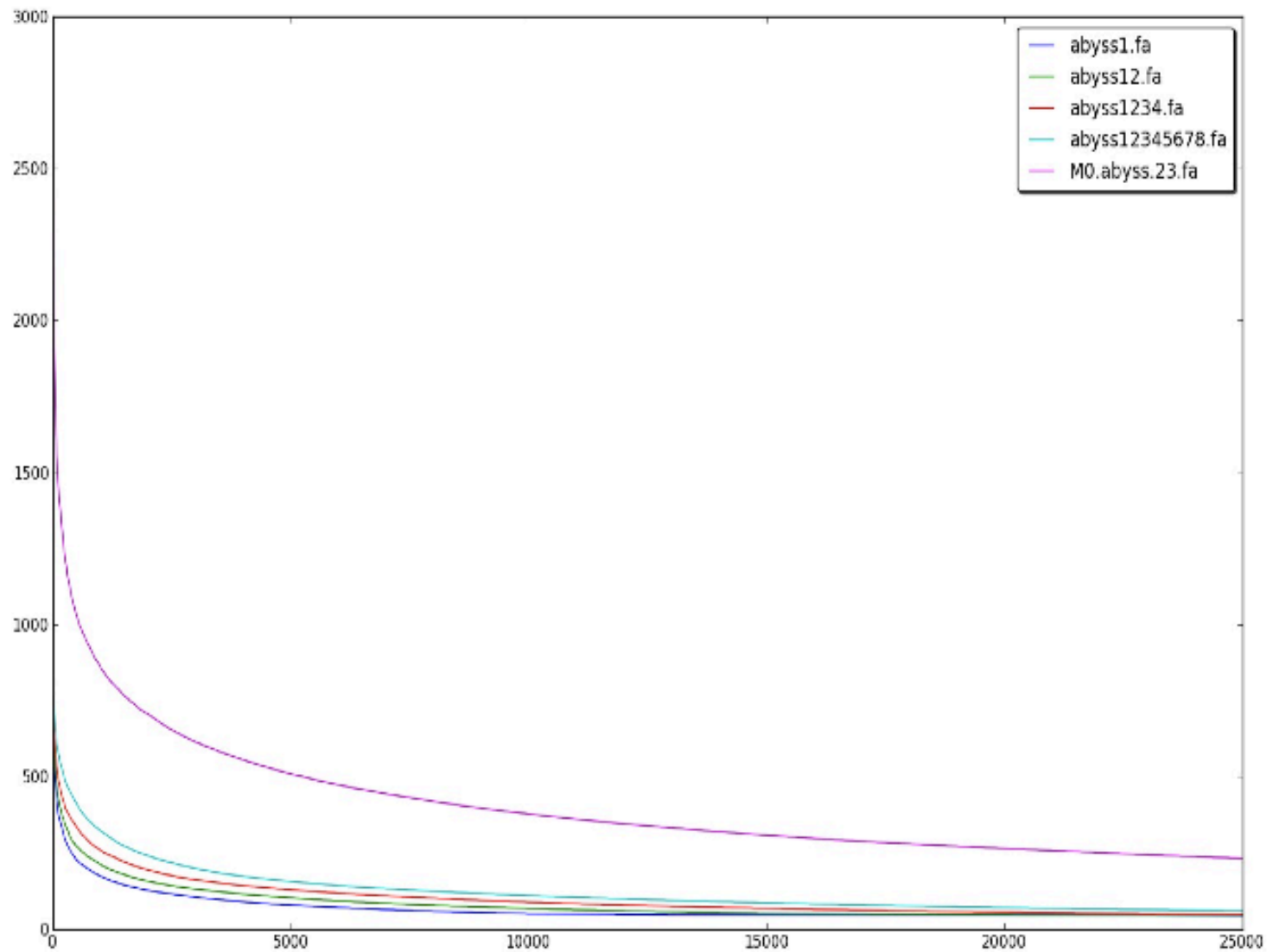
# Trying out different assemblers



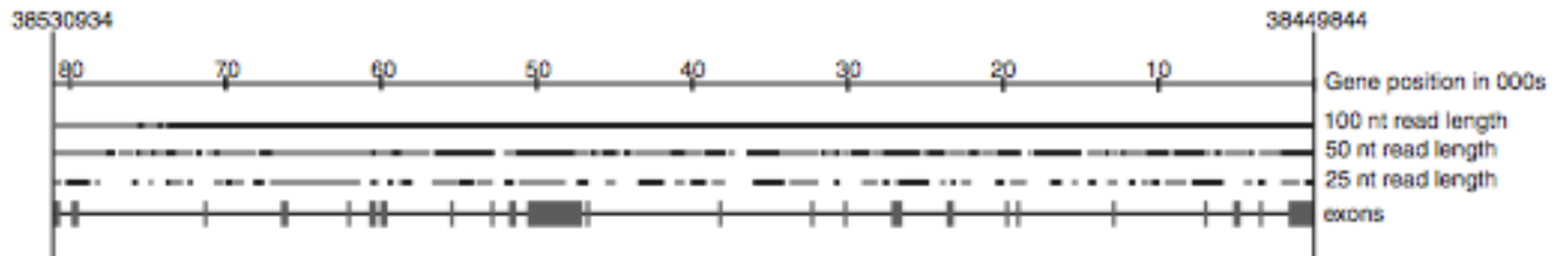
# Adding more data starts to saturate



# Longer reads help a lot... add 75bp



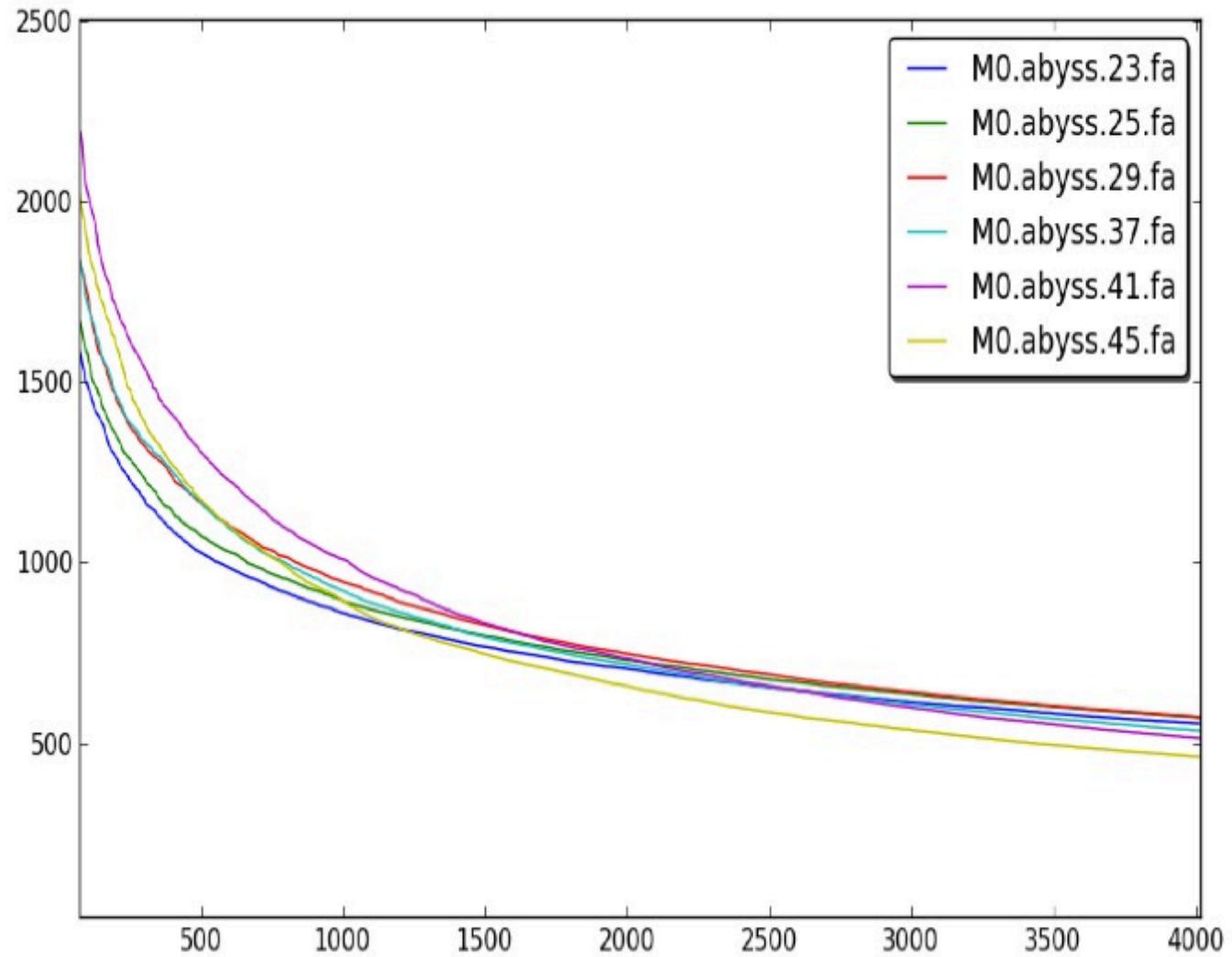
...this matches simulations



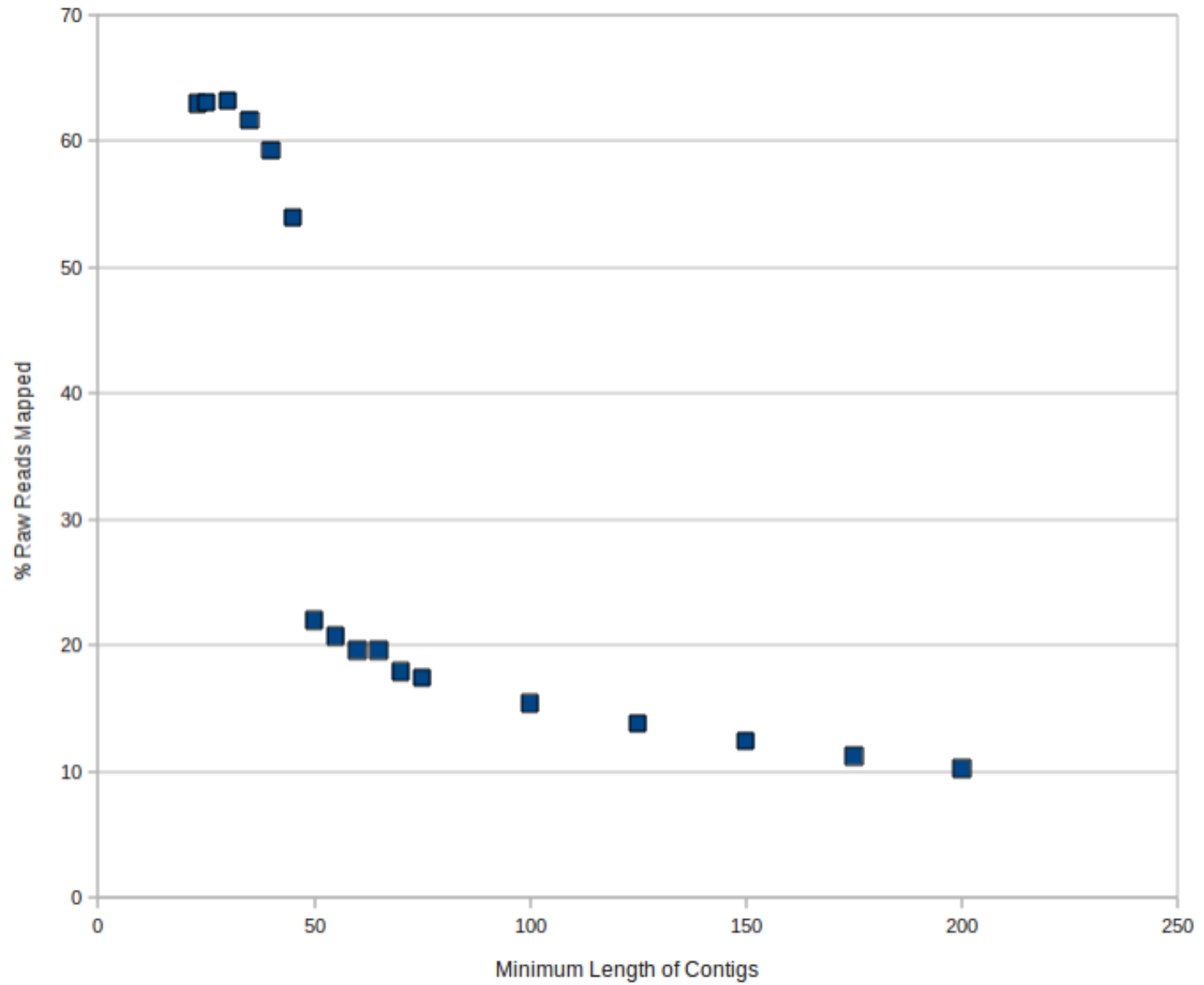
Whiteford et al., Nuc.Acid Res, 2005



# Different k values matter a bit



## How useful is assembly, anyway?





## RNAseq issues

- Most assemblers expect  $\sim 1:1$  ratio between different portions of the genome.
- Therefore they treat high-coverage k-mers as likely repeats.
- In RNAseq, genes are differentially expressed.



# Practical issues

- Do you have enough memory?
- Trim vs use quality scores?
- When is your assembly as good as it gets?
- Paired-end vs longer reads?
  
- More data is not *necessarily* better, if it introduces more errors.

# Today's schedule

- Short-read assembly (Jason; morning)
- Using UCSC genome browser to visualize things (Likit; afternoon)

*Reminder:* terminate your instances, please



- BBQ and fire.