



CS 68: BIOINFORMATICS

Prof. Sara Mathieson
Swarthmore College
Spring 2018



Outline: Feb 2

- Evaluation of assemblies
- Start: string alignment

Notes:

- Sequence alignment reading posted (from Durbin)

Recap: issues with de Bruijn graph assembly

- 1) Repeats of length $(k-1)$ or longer
- 2) Gaps in coverage
- 3) Differences in coverage
- 4) Sequencing errors

Evaluating Assemblies

Assembly evaluation

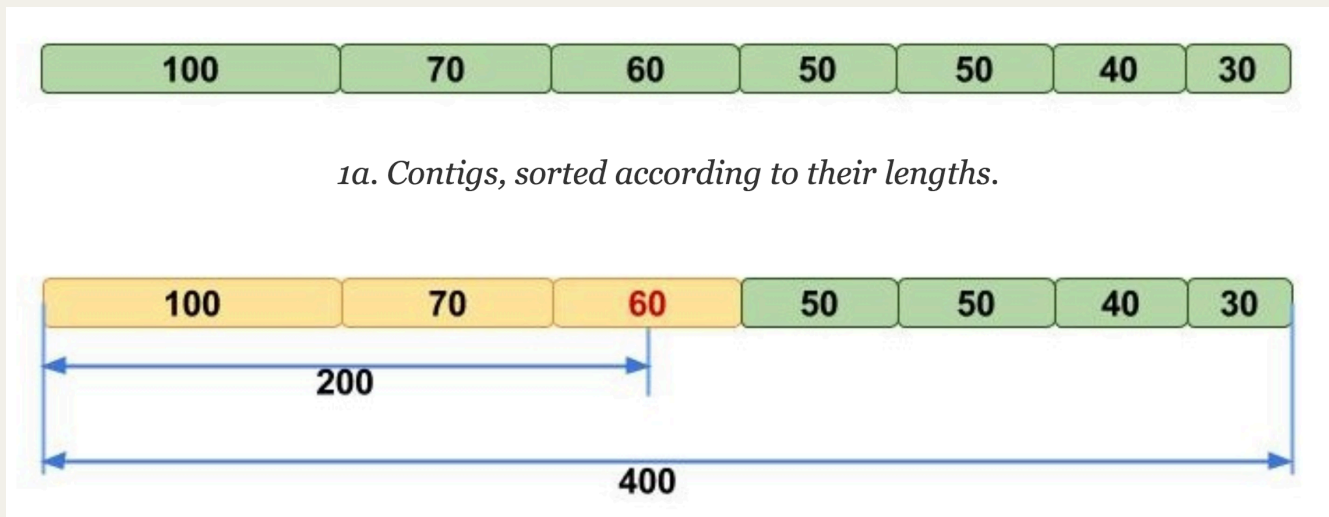
- **N50:** for a set of contigs, N50 is the greatest length such that at least half the bases of the assembly are in a contig with length N50 or longer



1a. Contigs, sorted according to their lengths.

Assembly evaluation

- **N50:** for a set of contigs, N50 is the greatest length such that at least half the bases of the assembly are in a contig with length N50 or longer



Why is N50 a bad evaluation metric?

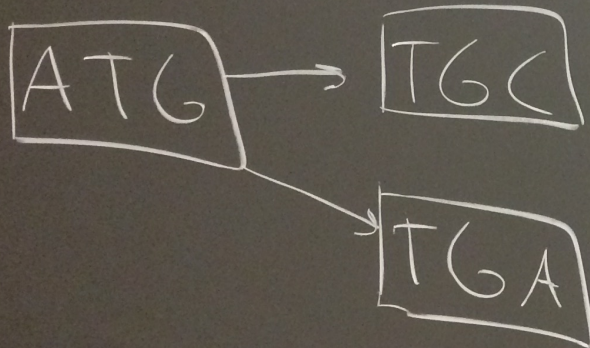
- We could just loop through cycles in our a graph over and over, generating large (incorrect) contigs
- We need a better way to evaluate the quality of assemblies
- Take away: simulated data is every valuable. Take an existing genome, simulate random reads, then try to reconstruct.

$k=4$

repeat
length
 $k-1$

ATGC

ATGA



① $\{100, 70, 60, 50, 50, 40, 30\}$
 $\underbrace{\hspace{1.5cm}}_{230}$

total 400

$$N_{50} = 60$$

② $\{10000, 150, 30, 20\}$

★

$$N_{50} = 10000$$

③ $\{100, 100, 100, 100, 100, 100\}$

$$N_{50} = 100$$

④ $\{1000, 250, 250, 250, 250\}$

$$N_{50} = 1000$$

⑤ $\{1000, 250, 250, 250, 250, 5\}$

$$\boxed{N50=250}$$

⑥ - easy to compute
- no ground truth required

⑦ - assembling sequence
that doesn't come from
the reference

⑧

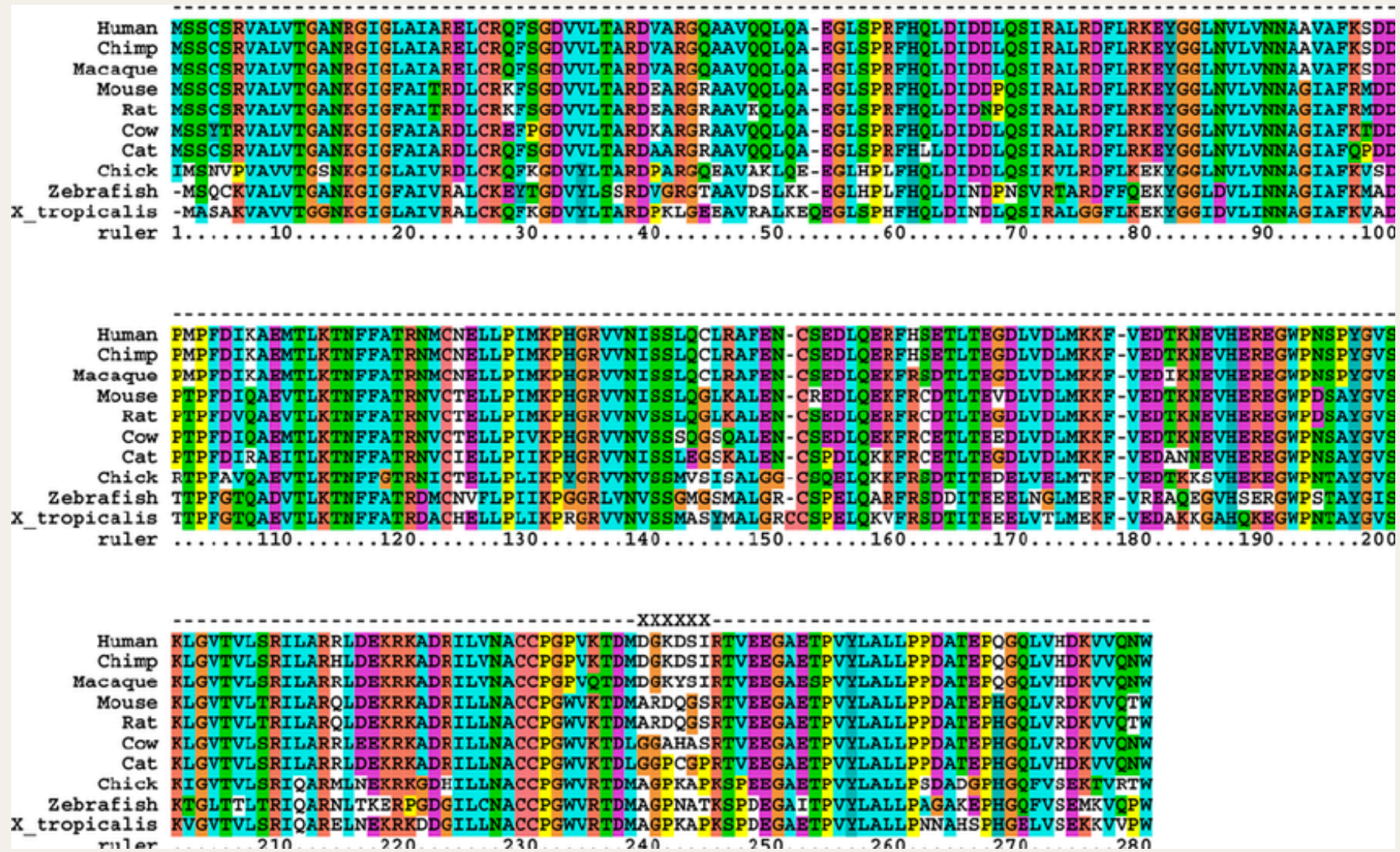
→ sequence alignment.

Sequence Alginment

Next Topic: sequence alignment

- Goal: given two sequences, what is the best match or “alignment” between them?
- Global alignment: align the entire sequences start to finish
- Local alignment: find portions of the two sequences with high similarity
- Homologous: sequences that are similar due to descent from a common ancestor
- Usually we are aligning homologous sequences (not sequences from completely different regions of the genome)

Example alignments: human, chimp, macaque + other species



Why sequence alignment?

- Understand evolutionary relationships between different species
- In particular: understanding fast-evolving bacterial and viral strains is important for health
- Understand protein function
- Understand diversity at the species level (important for diseases with a genetic component)

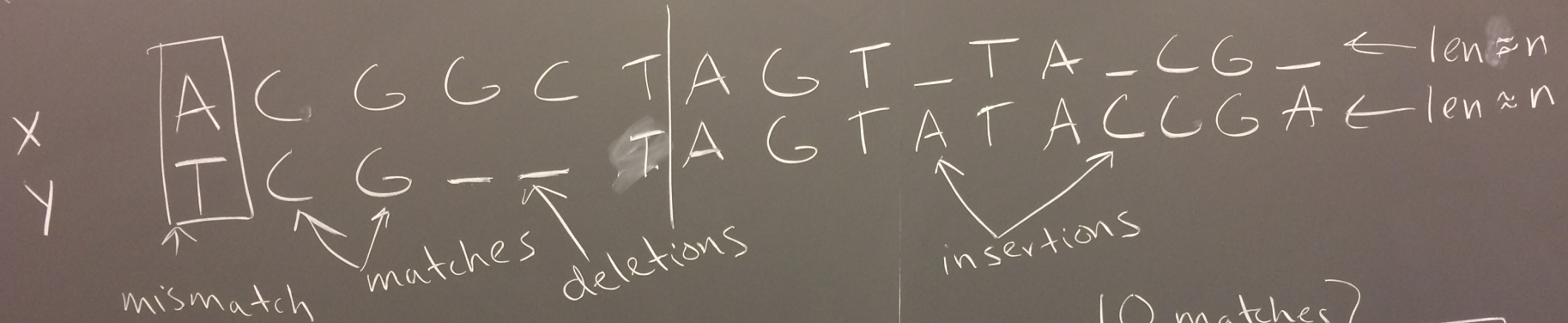
Example

■ ACGGCTAGTTACG

■ TCGTAGTATACCGA

- How should we “line them up” to get the best overlap?

Sequence Alignment



How do we score?

match: +1
 mismatch: -1
 gap: -2

$$S(x, y) = \left. \begin{array}{l} 10 \text{ matches} \\ -1 \text{ mismatch} \\ -10 \text{ gaps} \end{array} \right\} \boxed{-1}$$

want the best score
(highest)

Naive

$\binom{2n}{n}$ # bases

n max choices for gaps $\approx O(2^n)$

X —————

Y —————

Better way?

X = $\overset{x_1}{A} \overset{x_2}{A} \mid \overset{x_3}{A} \overset{x_4}{C}$

Y = $\overset{y_1}{A} \mid \overset{y_2}{G} \overset{y_3}{C}$

3 ways to end

① C
..... C

② C
C —

③ C —
..... C

| | | | | | |
|---|---|----|-----|----|----|
| | | 0 | 1 | 2 | 3 |
| i | | - | A | G | C |
| 0 | - | 0 | -2 | -4 | -6 |
| 1 | A | -2 | (1) | | |
| 2 | A | -4 | | | |
| 3 | A | -6 | | | |
| 4 | C | -8 | | | |

AAAC

$S(i, j)$ = best alignment score from $X_i \rightarrow Y_j$

A, AG, AGC
 $S(0,1)$ $S(0,2)$ $S(0,3)$

A, A-
 A, -2 -2 \rightarrow -4, -A
 | -4