

Lecture 2: Sequence Alignment



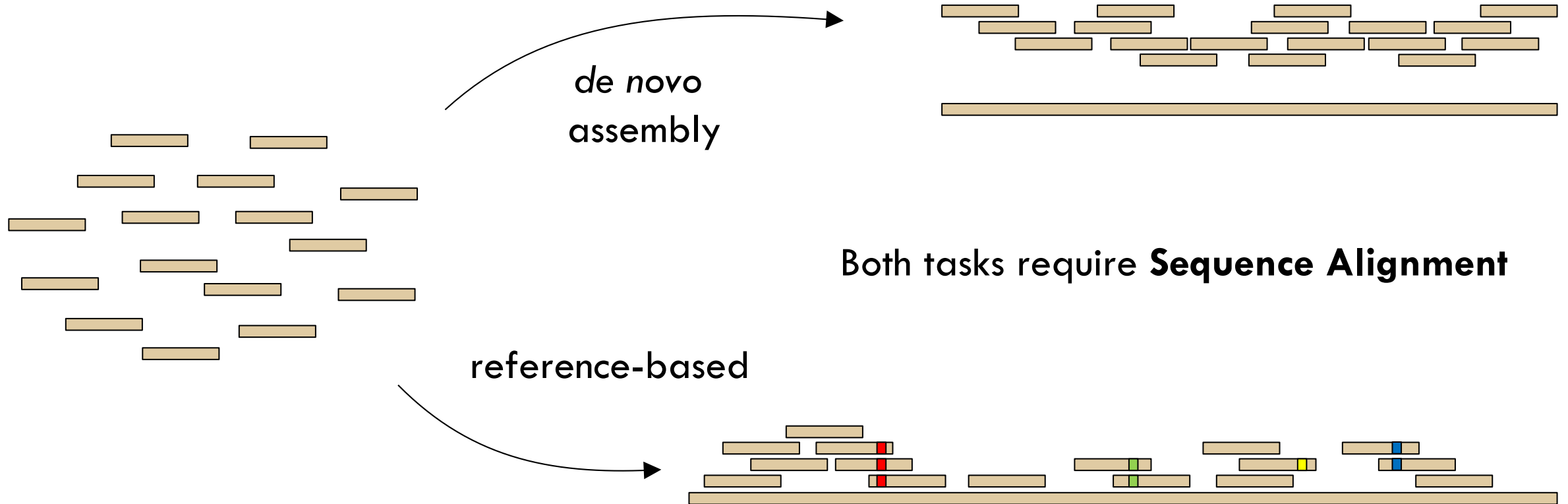
ECE 365 - Data Science and Genomics

Announcements

- Lab 1 (Introduction to Genomics) released
 - ▣ due March 11 at 11:59pm
 - ▣ Submit via gradescope

From last lecture: processing shotgun sequencing data

- 1 If a genome from that species has never been sequenced before



- 2 If a genome from that species has already been sequenced

Sequence Alignment

$x =$ **AGGCTATCACCTGACCTCCAGGCCGATGCCC**

$y =$ **TAGCTAGCACGACCGCGGTCGATTGCCCCGAC**

“Definition”: Given two sequences x and y , place gaps (‘-’) in them so that the resulting sequences “line up well”

-AGGCTA**T****CACCTGACCT****T****CCAGG****C****CGA--TGCCC---**
TAG-CTA**G****CAC--GACC****G****C--GG****T****CGATTGCCCCGAC**

What is a good alignment?

x = AGGCTAGTT

y = AGCGAAGTTT

AGGCTAGTT-
AGCGAAGTTT

6 matches, 3 mismatches, 1 gap

AGGCTA-GTT-
AG-CGAAGTTT

7 matches, 1 mismatch, 3 gaps

AGGC-TA-GTT-
AG-CG-AAGTTT

7 matches, 0 mismatch, 5 gaps

Scoring function

- We will score points for matches, and penalize mismatches and gaps:

mismatch: CTATCAC
 CTAGCAC

gap: CTATCAC
 CTA-CAC

- Scoring function: match: $+m$
 mismatch: $-s$
 gap: $-d$

$$\text{Score} = (\# \text{ matches}) \times m - (\# \text{ mismatches}) \times s - (\# \text{ gaps}) \times d$$

- Let's look at some examples!

Finding best **global** alignment

- We use an algorithm based on **dynamic programming**

match:

mismatch:

gap:

+1

-1

-1

	x=	G	C	A	T	T	C
y=	0						
G							
A							
T							
T							
A							
C							

Finding best **global** alignment

- We use an algorithm based on *dynamic programming*

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

		G	C	A	T	T	C
	0						
G							
A							
T							
T							
A							
C							

score of best alignment
between GCATT and GA

Finding best **global** alignment

- We use an algorithm based on *dynamic programming*

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

		G	C	A	T	T	C
	0						
G							
A							
T							
T							
A							
C							

score of best alignment
between GCA and GATTA

Finding best **global** alignment

- We use an algorithm based on *dynamic programming*

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

		G	C	A	T	T	C
	0	-1	-2	-3	-4	-5	-6
G	-1	1	0				
A	-2						
T	-3						
T	-4						
A	-5						
C	-6						

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + m, & \text{(if } x_i = y_j) \\ \underline{H_{i-1,j-1}} - s, & \text{(if } x_i \neq y_j) \\ \underline{H_{i-1,j}} - d \\ \underline{H_{i,j-1}} - d \end{cases}$$

Finding best **global** alignment

- We use an algorithm based on *dynamic programming*

G C A T T - C
G - A T T A C

5 - 0 - 2 = 3

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

		G	C	A	T	T	C
	0	-1	-2	-3	-4	-5	-6
G	-1	1	0	-1	-2	-3	-4
A	-2	0	0	1	0	-1	-2
T	-3	-1	-1	0	2	1	0
T	-4	-2	-2	-1	1	3	2
A	-5	-3	-3	-1	0	2	2
C	-6	-4	-2	-2	-1	1	3

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + m, & \text{(if } x_i = y_j) \\ H_{i-1,j-1} - s, & \text{(if } x_i \neq y_j) \\ H_{i-1,j} - d \\ H_{i,j-1} - d \end{cases}$$

Finding best **global** alignment

- We use an algorithm based on *dynamic programming*

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

		G	C	A	T	T	C
	0	-1	-2	-3	-4	-5	-6
G	-1	1	0	-1	-2	-3	-4
A	-2	0	0	1	0	-1	-2
T	-3	-1	-1	0	2	1	0
T	-4	-2	-2	-1	1	3	2
A	-5	-3	-3	-1	0	2	2
C	-6	-4	-2	-2	-1	1	3

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + m, & \text{(if } x_i = y_j) \\ H_{i-1,j-1} - s, & \text{(if } x_i \neq y_j) \\ H_{i-1,j} - d \\ H_{i,j-1} - d \end{cases}$$

GCATT-C

G-ATTAC

Finding best **global** alignment

- We use an algorithm based on *dynamic programming*

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

		G	C	A	T	T	C
	0	-1	-2	-3	-4	-5	-6
G	-1	1	0	-1	-2	-3	-4
A	-2	0	0	1	0	-1	-2
T	-3	-1	-1	0	2	1	0
T	-4	-2	-2	-1	1	3	2
A	-5	-3	-3	-1	0	2	2
C	-6	-4	-2	-2	-1	1	3

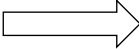
Needleman-Wunsch
algorithm

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + m, \\ \quad (\text{if } x_i = y_j) \\ H_{i-1,j-1} - s, \\ \quad (\text{if } x_i \neq y_j) \\ H_{i-1,j} - d \\ H_{i,j-1} - d \end{cases}$$

GCATT-C
G-ATTAC

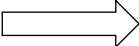
Global vs Local Alignment

- Global alignment: Align x and y fully

x =	AGGCTAGTT		AGGC-TA-GTT-
y =	AGCGAAGTTT		AG-CG-AAGTTT

Global vs Local Alignment

- Global alignment: Align x and y fully

x = **AGGCTAGTT**  **AGGC-TA-GTT-**
y = **AGCGAAGTTT** **AG-CG-AAGTTT**

- Local alignment: Align a substring of x to a substring of y

x = **ACCTACTCGCAATATGCTAGCAGCAACGTA**
y = **GTCGGGATCGCAATCTGCAGTCCGCTTACC**

- **ACCTACTCGCAATATGCTAG-CAGCAACGTA**
| | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | |
- **GTCGGGATCGCAATCTGC-AGTCCGCTTACC-**

Finding best **local** alignment

- We can adapt the previous algorithm to find local alignments

Finding best **local** alignment

- We can adapt the previous algorithm to find local alignments

match:
mismatch:
gap:

+3

-3

-2

		T	G	T	T	A	C	G	G
	0	0	0	0	0	0	0	0	0
G	0	0	3						
G	0								
T	0								
T	0								
G	0								
A	0								
C	0								
T	0								
A	0								

Smith-Waterman algorithm

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + m, & (\text{if } x_i = y_j) \\ H_{i-1,j-1} - s, & (\text{if } x_i \neq y_j) \\ H_{i-1,j} - d \\ H_{i,j-1} - d \\ \mathbf{0} \end{cases}$$

Finding best **local** alignment

- We can adapt the previous algorithm to find local alignments

match:
mismatch:
gap:

+3

-3

-2

		T	G	T	T	A	C	G	G
	0	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3	3
G	0	0	3	1	0	0	0	3	6
T	0	3	1	6	4	2	0	1	4
T	0	3	1	4	9	7	5	3	2
G	0	1	6	4	7	6	4	8	6
A	0	0	4	3	5	10	8	6	5
C	0	0	2	1	3	8	13	11	9
T	0	3	1	5	4	6	11	10	8
A	0	1	0	3	2	7	9	8	7

Smith-Waterman
algorithm

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + m, \\ \quad (\text{if } x_i = y_j) \\ H_{i-1,j-1} - s, \\ \quad (\text{if } x_i \neq y_j) \\ H_{i-1,j} - d \\ H_{i,j-1} - d \\ 0 \end{cases}$$

GTT - AC

GTTGAC

Finding best **local** alignment

- We can adapt the previous algorithm to find local alignments

match: +3
 mismatch: $s = +3$
 gap: $d = +2$

		T	G	T	T	A	C	G	G
	0	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3	3
G	0	0	3	1	0	0	0	3	6
T	0	3	1	6	4	2	0	1	4
T	0	3	1	4	9	7	5	3	2
G	0	1	6	4	7	6	4	8	6
A	0	0	4	3	5	10	8	6	5
C	0	0	2	1	3	8	13	11	9
T	0	3	1	5	4	6	11	10	8
A	0	1	0	3	2	7	9	8	7

Smith-Waterman
algorithm

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + m, & (\text{if } x_i = y_j) \\ H_{i-1,j-1} - s, & (\text{if } x_i \neq y_j) \\ H_{i-1,j} - d \\ H_{i,j-1} - d \\ 0 \end{cases}$$

GTT-AC
GTTGAC

Alignment problem variations

Global alignment

GCATT-C
G-ATTAC

		G	C	A	T	T	C
	0	-1	-2	-3	-4	-5	-6
G	-1	1	0	-1	-2	-3	-4
A	-2	0	0	1	0	-1	-2
T	-3	-1	-1	0	2	1	0
T	-4	-2	-2	-1	1	3	2
A	-5	-3	-3	-1	0	2	2
C	-6	-4	-2	-2	-1	1	3

Local alignment

T GTT-AC GG
G GTTGAC TA

		T	G	T	T	A	C	G	G
	0	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3	3
G	0	0	3	1	0	0	0	3	6
T	0	3	1	6	4	2	0	1	4
T	0	3	1	4	9	7	5	3	2
G	0	1	6	4	7	6	4	8	6
A	0	0	4	3	5	10	8	6	5
C	0	0	2	1	3	8	13	11	9
T	0	3	1	5	4	6	11	10	8
A	0	1	0	3	2	7	9	8	7

Alignment problem variations

Global alignment

GCATT-C
G-ATTAC

		G	C	A	T	T	C
	0	-1	-2	-3	-4	-5	-6
G	-1	1	0	-1	-2	-3	-4
A	-2	0	0	1	0	-1	-2
T	-3	-1	-1	0	2	1	0
T	-4	-2	-2	-1	1	3	2
A	-5	-3	-3	-1	0	2	2
C	-6	-4	-2	-2	-1	1	3

Alignment to reference

GCCTT-C → short
TCT G-ACCTC GCGT
long

[illegible]

Local alignment

T GTT-AC GG
G GTTGAC TA

		T	G	T	T	A	C	G	G
	0	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3	3
G	0	0	3	1	0	0	0	3	6
T	0	3	1	6	4	2	0	1	4
T	0	3	1	4	9	7	5	3	2
G	0	1	6	4	7	6	4	8	6
A	0	0	4	3	5	10	8	6	5
C	0	0	2	1	3	8	13	11	9
T	0	3	1	5	4	6	11	10	8
A	0	1	0	3	2	7	9	8	7

Alignment problem variations

Global alignment

		G	C	A	T	T	C
	0	-1	-2	-3	-4	-5	-6
G	-1	1	0	-1	-2	-3	-4
A	-2	0	0	1	0	-1	-2
T	-3	-1	-1	0	2	1	0
T	-4	-2	-2	-1	1	3	2
A	-5	-3	-3	-1	0	2	2
C	-6	-4	-2	-2	-1	1	3

GCATT-C
G-ATTAC

Local alignment

		T	G	T	T	A	C	G	G
	0	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3	3
G	0	0	3	1	0	0	0	3	6
T	0	3	1	6	4	2	0	1	4
T	0	3	1	4	9	7	5	3	2
G	0	1	6	4	7	6	4	8	6
A	0	0	4	3	5	10	8	6	5
C	0	0	2	1	3	8	13	11	9
T	0	3	1	5	4	6	11	10	8
A	0	1	0	3	2	7	9	8	7

T GTT-AC GG
G GTTGAC TA

Alignment to reference

[illegible]

GCCT-C
TCT G-CTAC GCGT

Overlap alignment

GCCTTCA
TTCT G-CTTCA

[illegible]