



CS 68: BIOINFORMATICS

Prof. Sara Mathieson
Swarthmore College
Spring 2018



Outline: Feb 9

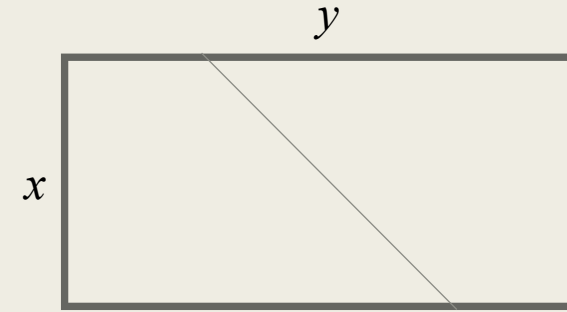
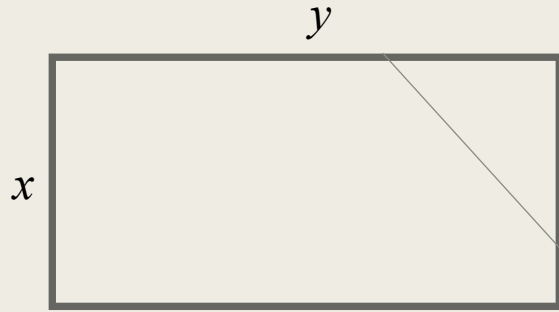
- Finish: local alignment variations
- Begin: Burrows-Wheeler Transform (BWT)
- Application to read mapping

Notes:

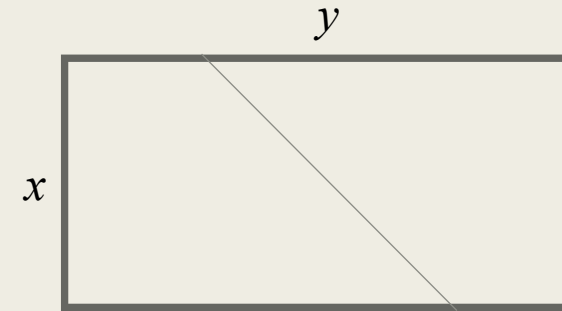
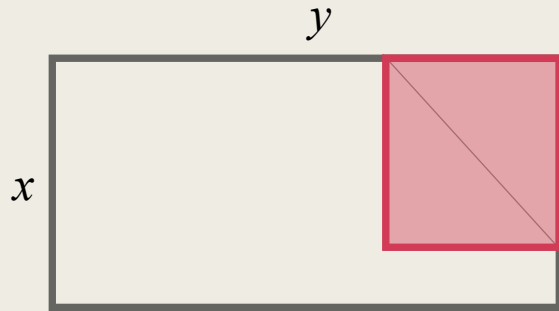
- Reading posted
- Feb 15: Lab 4 (BWT and read mapping)
- Feb 22: practice midterm (I am at SIGCSE)
- Mar 1: midterm 1
- Lab 1 returned on Monday

Local alignment variations

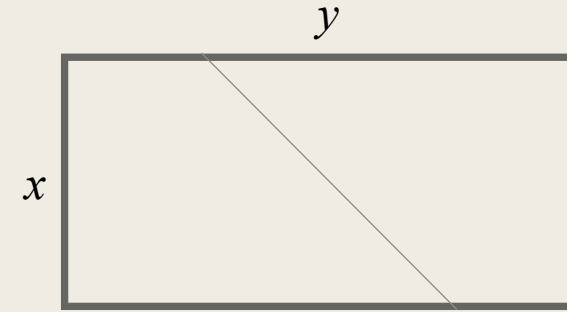
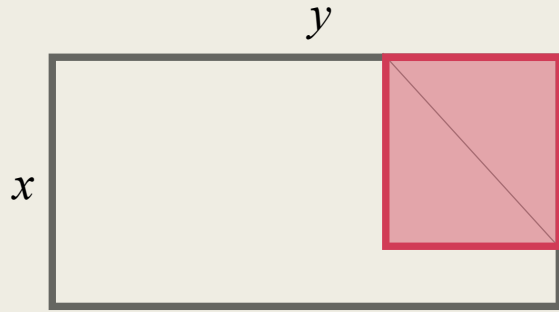
Handout 6: what portion of x aligns to what portion of y ?



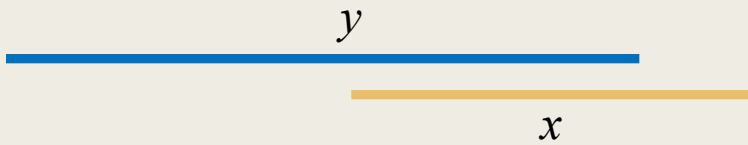
Handout 6: what portion of x aligns to what portion of y ?



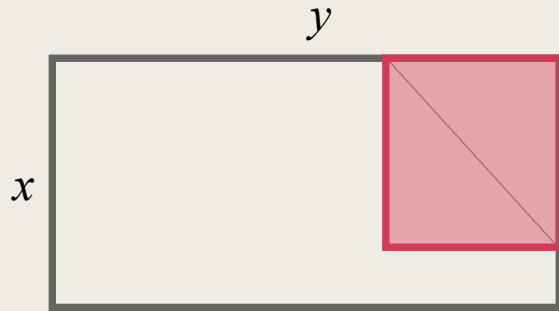
Handout 6: what portion of x aligns to what portion of y ?



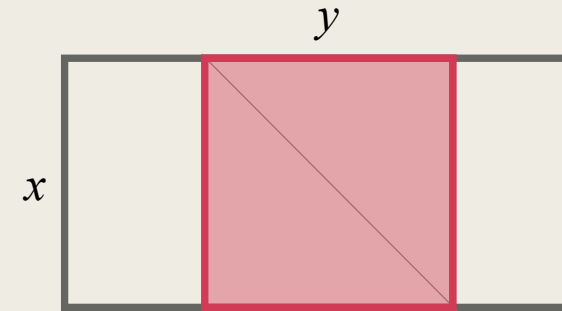
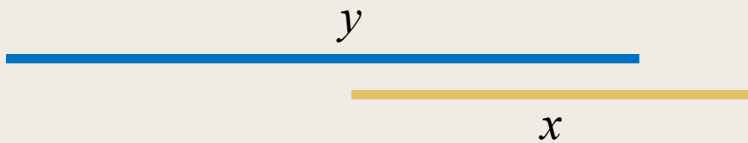
Beginning of x with end of y



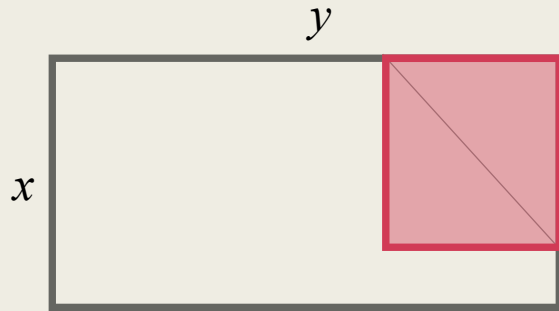
Handout 6: what portion of x aligns to what portion of y ?



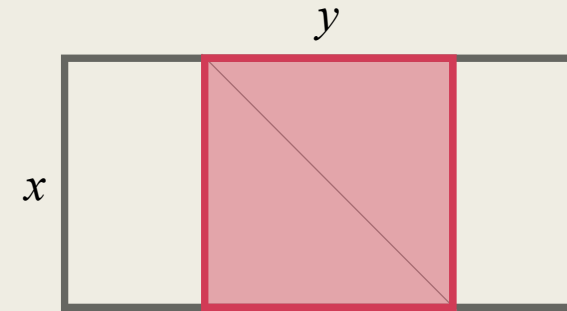
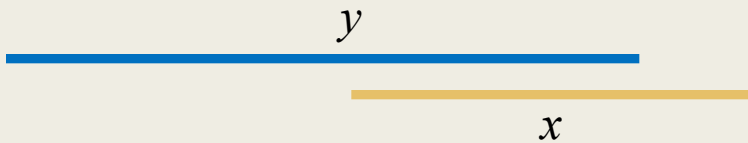
Beginning of x with end of y



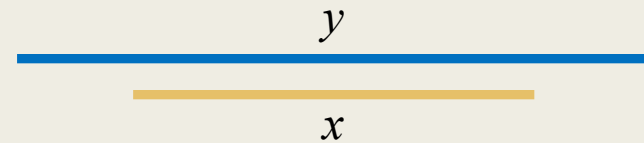
Handout 6: what portion of x aligns to what portion of y ?



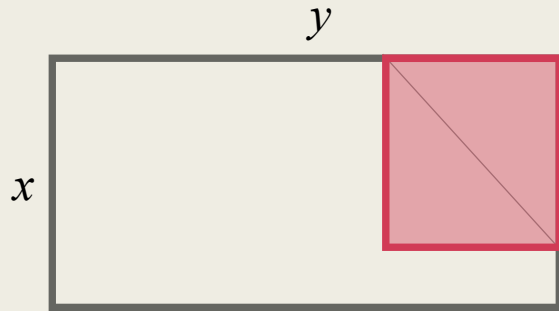
Beginning of x with end of y



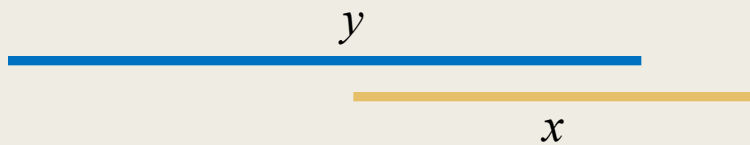
x aligns with the middle of y



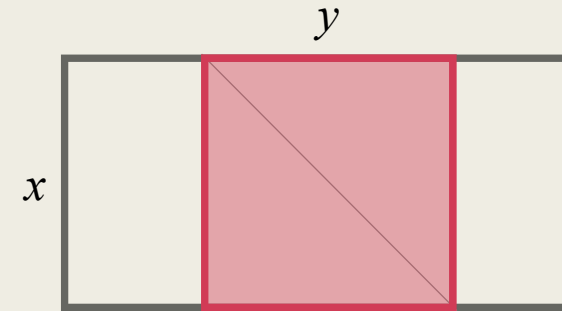
Handout 6: what portion of x aligns to what portion of y ?



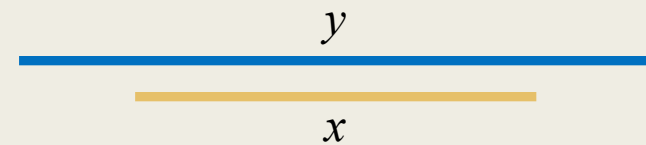
Beginning of x with end of y



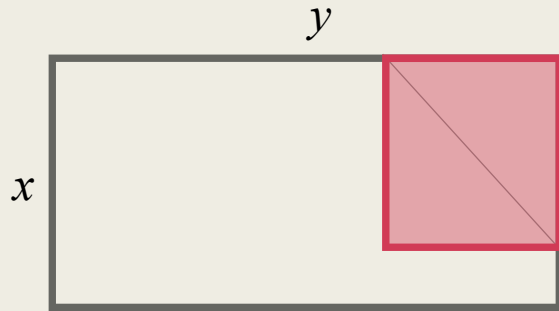
DP algorithm modifications:



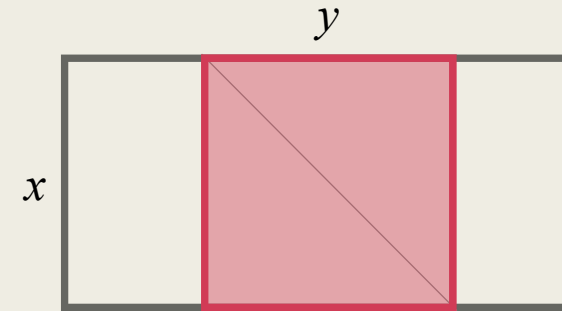
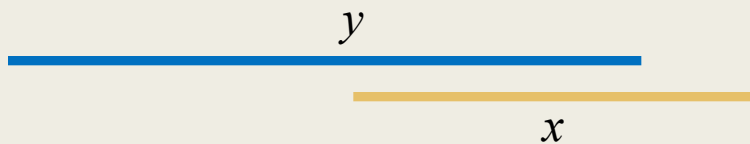
x aligns with the middle of y



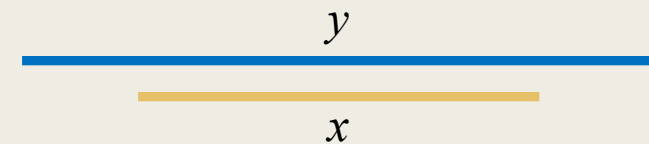
Handout 6: what portion of x aligns to what portion of y ?



Beginning of x with end of y



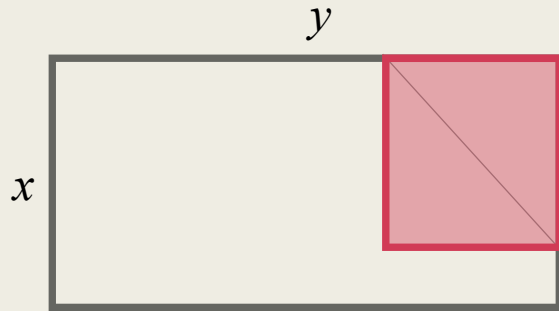
x aligns with the middle of y



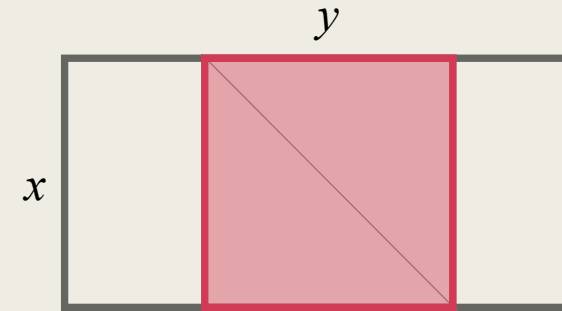
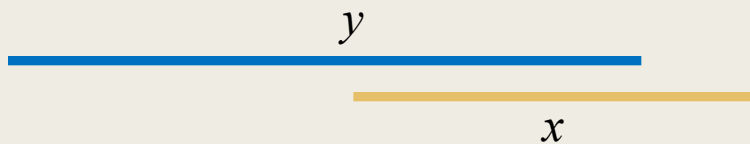
DP algorithm modifications:

1) Initialization: 0^{th} row and 0^{th} column with 0's to not penalize leading/trailing gaps

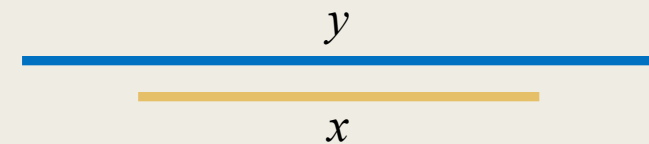
Handout 6: what portion of x aligns to what portion of y?



Beginning of x with end of y



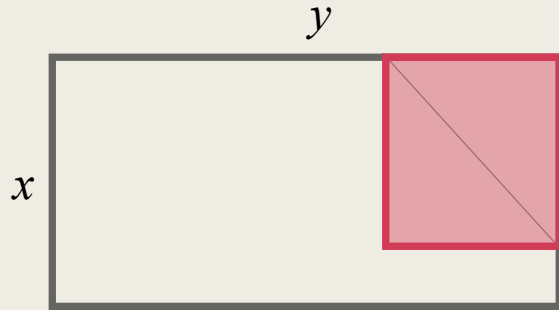
x aligns with the middle of y



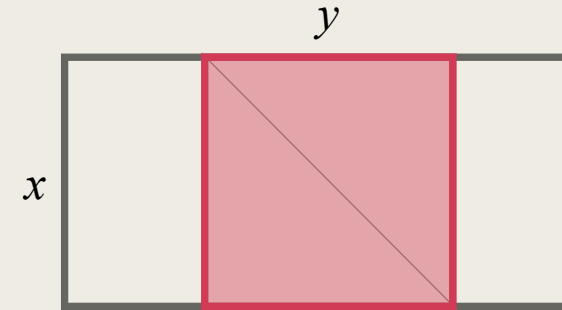
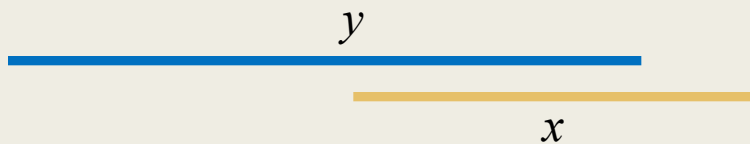
DP algorithm modifications:

- 1) **Initialization:** 0^{th} row and 0^{th} column with 0's to not penalize leading/trailing gaps
- 2) **Recursion:** to fill in the rest of the table, use global alignment (i.e. don't restart at 0)

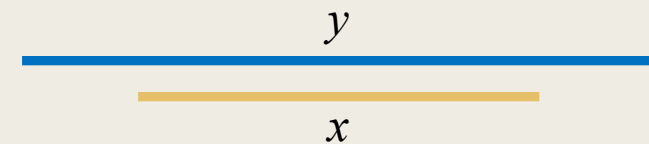
Handout 6: what portion of x aligns to what portion of y ?



Beginning of x with end of y



x aligns with the middle of y



DP algorithm modifications:

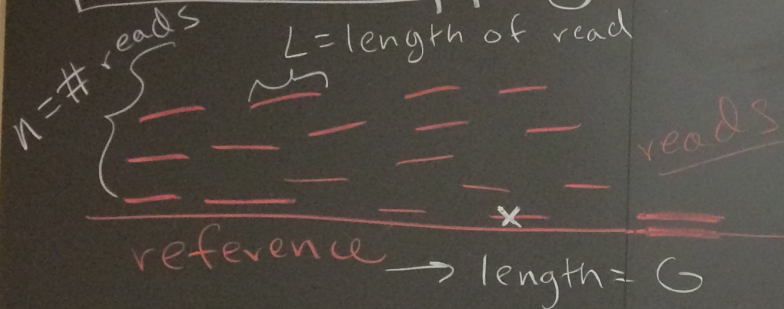
- 1) **Initialization:** 0^{th} row and 0^{th} column with 0's to not penalize leading/trailing gaps
- 2) **Recursion:** to fill in the rest of the table, use global alignment (i.e. don't restart at 0)
- 3) **Traceback:** start at the maximum value along the last row or last column

Next topic: read mapping with BWT

What if y is very long and x is very short?

- Exactly the case if **y = entire genome** and **x = single read**
- Reading mapping: Given an already assembled reference genome and reads from a newly sequence individual of the same species, what is the position of each read?
- Three options:
 - *1) Read aligns perfectly with reference*
 - *2) Read aligns with a few differences (representing population-level variation)*
 - *3) Read does not align at all (insertion in the newly sequenced individual)*

Read Mapping



$O(G \cdot L)$

insertion

runtime \approx billions

$O(L \cdot G \cdot n)$

50-100

billions

BWT

$S = \text{banana}\$$

- special char
- alphabetically first

$BWT(S) = CCCC AAAA AAB B$

$S = A B C A B C$

$6 \cdot A + 12 \cdot B + 6$

$\pi(S)$	rank	$\pi^{\text{sorted}}(S)$	BW(S)
banana\$	5	\$banana	a
anana\$b	4	a\$banan	n
nana\$ba	7	ana\$ban	n
ana\$ban	3	anana\$b	b
na\$bana	6	banana\$	\$
a\$banan	2	na\$bana	a
<u>\$banana</u>	1	nana\$ba	a
all cyclic permutations of S		<div style="display: flex; justify-content: space-around; align-items: center;"> F (first) L (last) </div>	<div style="display: flex; justify-content: space-around; align-items: center;"> repeat anana </div>

B

G.A

<div style="display: flex; justify-content: space-around;"> F L </div>	backtrace	reconstruct s
$\$ \rightarrow a_1$	$\$ \rightarrow a_1$	$a\$$
$a_1 \rightarrow n_1$	$a_1 \rightarrow n_1$	$na\$$
$a_2 \rightarrow n_2$	$n_1 \rightarrow a_2$	$ana\$$
$a_3 \rightarrow b_1$	$a_2 \rightarrow n_2$	$nana\$$
$b_1 \rightarrow \$$	$n_2 \rightarrow a_3$	$anana\$$
$n_1 \rightarrow a_2$	$a_3 \rightarrow b_1$	$banana\$$
$n_2 \rightarrow a_3$	$b_1 \rightarrow \$$	
	<u><u>STOP</u></u>	

Claim If F has k copies
of char $c : c_1, c_2, \dots, c_k$,
the their order is preserved
in L

$$\sigma_1 < \sigma_2 < \dots < \sigma_k \Leftrightarrow \sigma_1 c_1 < \sigma_2 c_2 < \dots < \sigma_k c_k$$

