

Space-Optimal Alignments

Wednesday, February 10, 2021 10:46 AM

Previously,

Alignment models: { Global (NW algo)

Local (SW algo)

semi-global (Gotoh's algo)

Gap models:

{ Linear gap (g)

{ Affine gap (b+g)

Complexities:

Runtime = $O(mn)$

→ Space = $O(mn)$

"Reads" ~ 100 bp

S_1 : $m \sim 100$
 S_2 : $n \sim 100$

space for DP table $\propto 10^4$ cells

$\begin{matrix} \% & D \\ \% & I \end{matrix}$ ⇒ each cell ≈ 12 bytes
 Space (DP table) $\propto 120$ KB

Gene S ~ 10^4 bp

Space $\propto (10^4 \times 10^4) \times 12$

$\approx 12 \times 10^8$

≈ 1200 MB

≈ 1.2 GB

Genomes ~

Covid (viral)

30K

900M

~ 1B cells

12 GB

10^9 (human)

10^{18} cells

↓
Exabytes

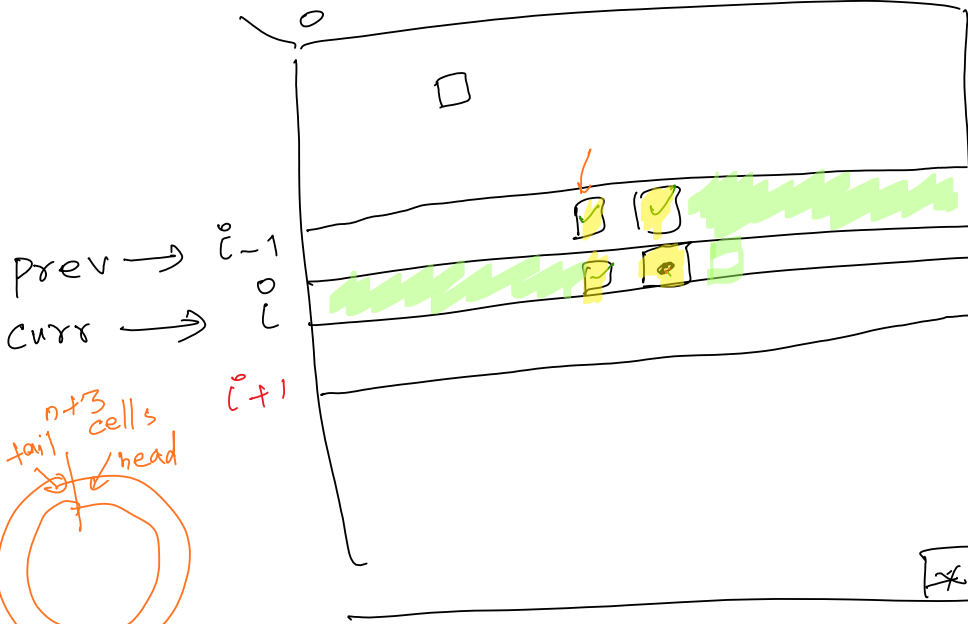
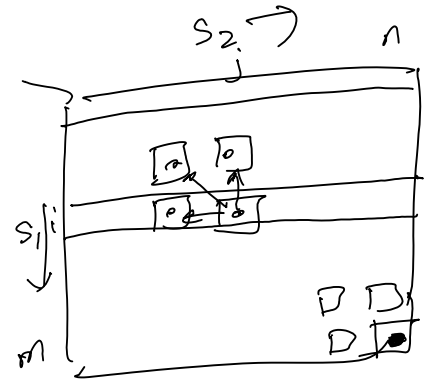
Space-optimal alignments

Wednesday, February 10, 2021 10:46 AM

Q) Can we store only part of the DP table?

Approach:

a) Do we need the full table for computing the optimal score?



Space cost

$$= O(m + n + 2n)$$

↑ ↑ ↑
 s_1 s_2 2 rows (curr, prev)

Linear

(opt score)

Space-optimal alignments

Wednesday, February 10, 2021 10:46 AM

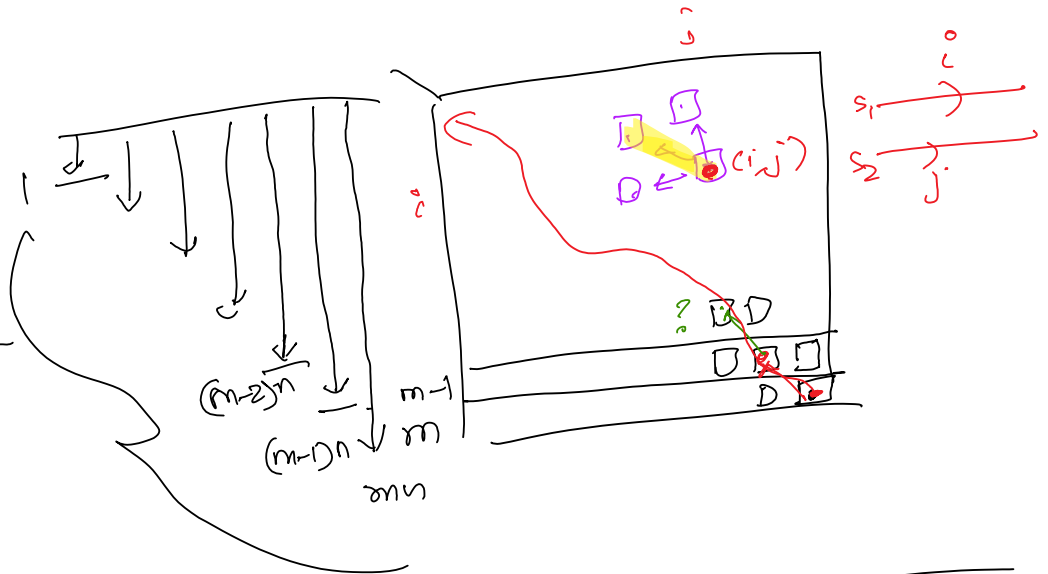
Q) what if we also need to output/reconstruct the optimal path?

brute force

Time cost

$$\propto n(1+2+\dots+m)$$

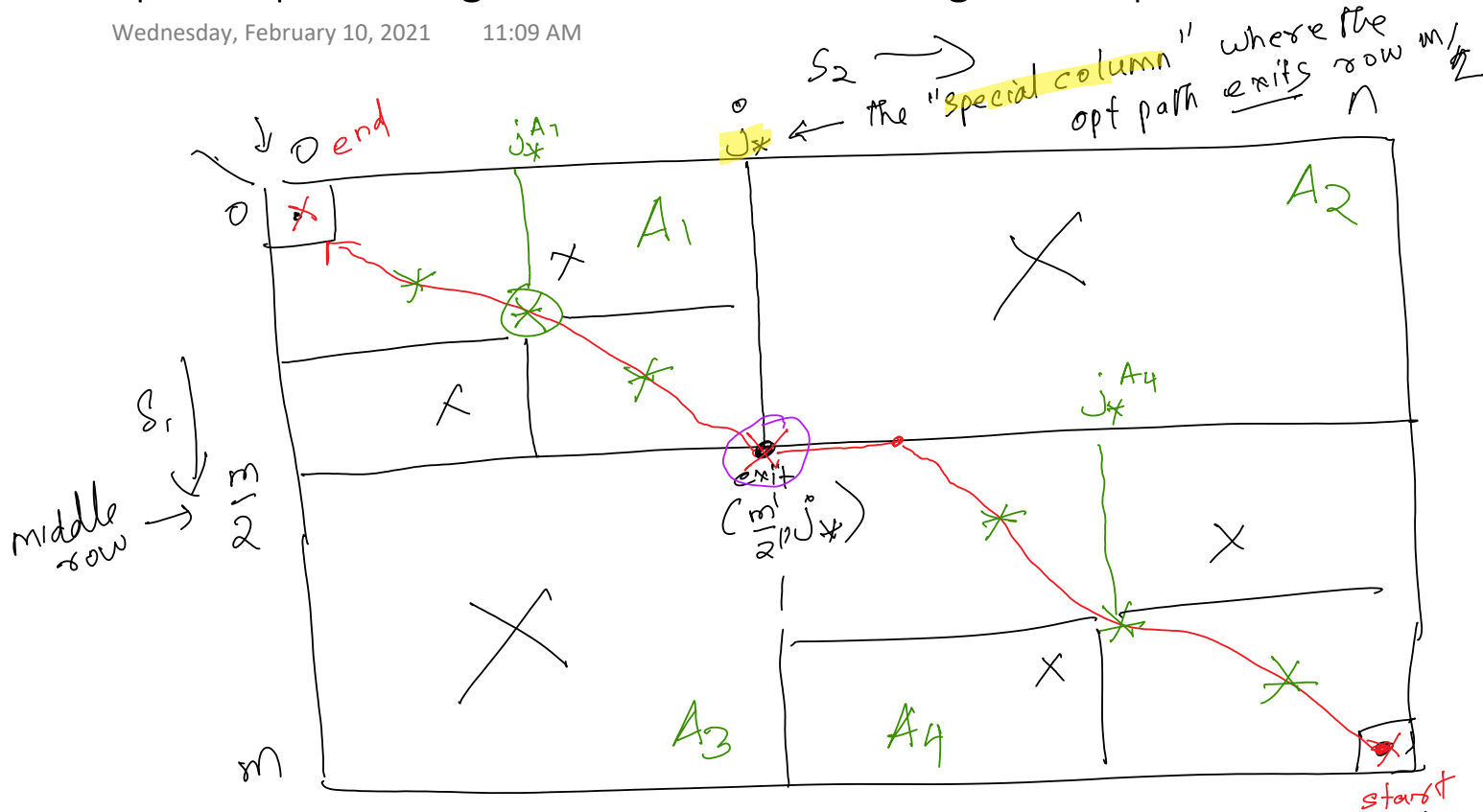
$$\propto n(m^2)$$



Goal: To compute the opt path in $O(m+n)$ space and in $O(mn)$ time (space-optimal)

Space-optimal alignments: The Hirschberg Technique

Wednesday, February 10, 2021 11:09 AM



$$A_1 + A_2 + A_3 + A_4 = (m+1)(n+1)$$

oracle $\Rightarrow j^*$

Recursive step

1

2

3

4

...

cells computed

mn

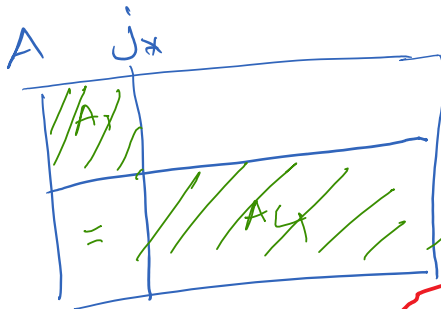
$+ \frac{mn}{2}$

$+ \frac{mn}{4}$

\dots

$\frac{mn}{8}$

$$= A_1 + A_4$$



$$mn \left(1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots \right)$$

$$\Rightarrow 2mn = O(mn) \text{ time}$$

If we could find j^* then we could do the optimal reface in $O(mn)$ time.

Determining the oracle (for the special column)

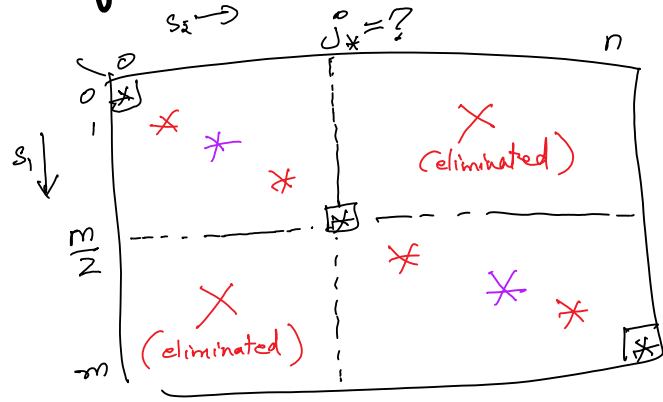
Friday, February 12, 2021 10:20 AM

Q) How to determine J_x^* (i.e. design of the oracle)?

Key Observation:

If you split an optimal path \mathcal{P} into disjoint subpaths $\{P_1, P_2, \dots, P_k\}$ then each subpath P_i has to be optimal for that sub problem.

$$\text{opt path } \mathcal{P} \hat{=} P_1 \cdot P_2 \cdot \dots \cdot P_i \cdot \dots \cdot P_k$$



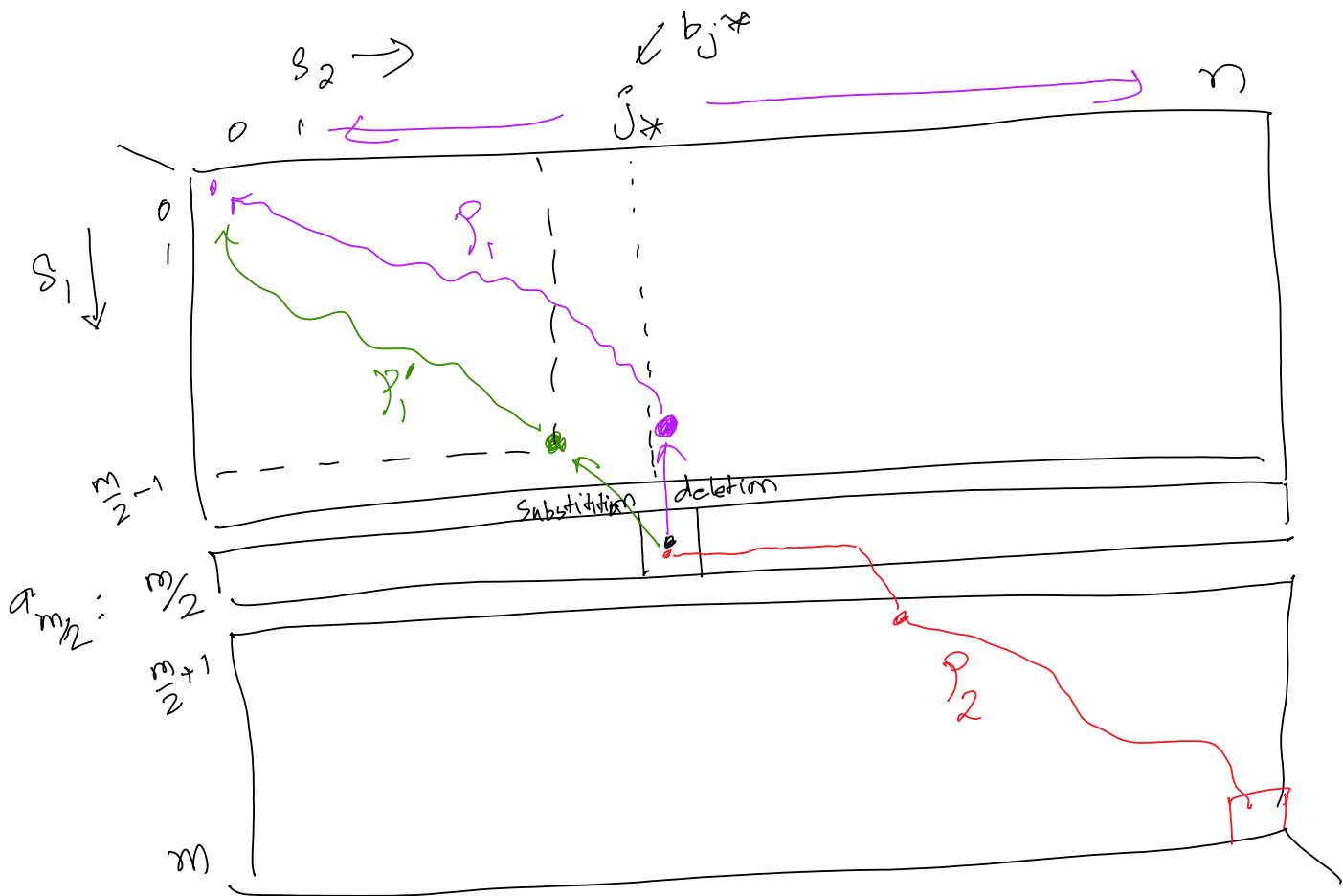
An example:

let the (unknown) optimal alignment be:

	a_1	a_2	a_3	a_4	a_5		a_6	a_7	a_8	a_9	
s_1 :	a	a	c	g	c	-	-	t	a	c	g
s_2 :	a	t	c	g	c	a	a	t	a	-	g
	b_1	b_2	b_3	b_4	b_5		b_6	b_7	b_8	b_9	b_{10}
P :	P_1					J_x^*	P_2				
	$\left(\begin{array}{l} s_1[a_1 \dots a_5] \\ s_2[b_1 \dots b_5] \end{array} \right)$					\uparrow	$\left(\begin{array}{l} s_1[a_6 \dots a_9] \\ s_2[b_6 \dots b_{10}] \end{array} \right)$				

Oracle for j^*

Friday, February 12, 2021 10:49 AM



Our optimal path

$$P = \left\{ \begin{array}{l} P_1 \cdot \binom{a_{m/2}}{-} \cdot P_2 \\ P_1' \cdot \binom{a_{m/2}}{b_{j^*}} \cdot P_2 \end{array} \right\} \text{ (or)}$$

$$\text{Score}(P) = \max \left\{ \begin{array}{l} \text{score}(P_1) + g + \text{score}(P_2) \\ \text{score}(P_1') + g(a_{m/2}, b_{j^*}) + \text{score}(P_2) \end{array} \right\}$$

Hirschberg technique (1975)

Oracle for j^*

Friday, February 12, 2021 11:32 AM

At row $m/2$:

Evaluate all possible j and figure out j^*

for $j = 0$ to n {

$$\text{score}_1 = \text{score}(s_1[1.. \frac{m}{2}], s_2[1.. j]) + \delta(a_{\frac{m}{2}}, b_j) + \text{score}(s_1[\frac{m}{2}+1.. m], s_2[j+1.. n])$$

$$\text{score}_2 = \text{score}(s_1[1.. \frac{m}{2}-1], s_2[1.. j]) + g + \text{score}(s_1[\frac{m}{2}+1.. m], s_2[j+1.. n])$$

$$\text{score}_j = \max \left\{ \begin{array}{l} \text{score}_1 \\ \text{score}_2 \end{array} \right\}$$

$O(mn)$

$j^* = \text{Arg Max}_j (\text{score}_j) \quad 0 \leq j \leq n$

