

# From Orthogonal Vectors to Longest Common Subsequence

Definition: Longest Common Subsequence Problem (LCS)

Input: strings  $x, y$  of length  $|x|=n \geq |y|=m$

Task: Compute longest string  $z$  that is a subsequence of both  $x$  and  $y$

Example: (We write  $LCS(x, y) = |z|$ )

A G C A T  
 G A C

A G C A T  
 G A C

Dynamic Programming algorithm:

Idea: Build a table

	$x[1]$	...	$x[n]$
$y[1]$			
$\vdots$			
$y[n]$			

$$T[i, j] = LCS[x[1..i], y[1..j]]$$

$$\Rightarrow T[n, m] = LCS(x, y)$$

Determine new entry by computation on "surrounding" entries

$$T[0, j] = 0 \text{ and } T[i, 0] = 0$$

How to determine  $T[i, j]$  for  $1 \leq i \leq n$  and  $1 \leq j \leq m$

~~$x[i] = y[j]$~~

"delete in x"      "delete in y"

$$T[i, j] = \max \{ T[i-1, j], T[i, j-1] \}$$

If  $x[i] = y[j]$ :

$$T[i, j] = \max \{ T[i, j], T[i-1, j-1] + 1 \}$$

match

Time Complexity:  $O(n^3)$

Fastest Known Algorithm:  $O(n^2 / \log^2 n)$  [Masek, Paterson '80]

Theorem: Assuming OVH, ~~the~~ LCS has no algorithm with running time  $O(n^{2-\delta})$  (for constant  $\delta > 0$ ).

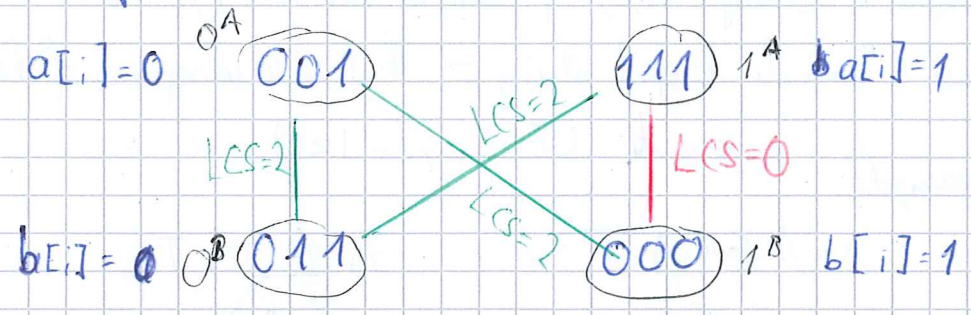
[Bringmann/Kimnour '15, Alford/Bachurs/V. Williams '15]



Given OV-instance  $A, B \in \{0,1\}^d$   $|A|=|B|=n$

Remember:  $a \perp b \Leftrightarrow \forall i \forall j \forall 1 \leq i \leq d (a[i] \cdot b[j] = 0)$

Coordinate Gadgets: Simulate behavior of  $a[i] \cdot b[j]$



number  $a[i]$   $\rightarrow$  string  $a[i]^A$   
 $b[i]$   $\rightarrow$  string  $b[i]^B$

Now  $LCS(a[i]^A, b[i]^B)$  can be written as  $f(a_i, b_i)$  with  $f(0) \neq f(1)$

Vector Gadgets: Simulate orthogonality

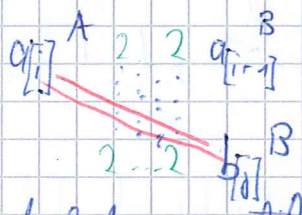
Concatenate coordinate gadgets, pad with new symbol 2

$VG(a) = a_1^A \ 2 \dots 2 \ a_2^A \ 2 \dots 2 \ a_3^A \ 2 \dots 2 \ \dots \ a_d^A$

$VG(b) = b_1^B \ 2 \dots 2 \ b_2^B \ 2 \dots 2 \ b_3^B \ 2 \dots 2 \ \dots \ b_d^B$

Claim: no LCS matches symbols in  $a_i^A$  with symbols in  $b_j^B$  when  $i \neq j$

Proof: assume otherwise (w.l.o.g.  $i < j$ )



We miss at least one possibility for matching two 2-sequences

Thus we will match at most  $(d-2) \cdot 4d$  symbols 2 and  $\leq 3d$  symbols 0/1

But:  $LCS(VG(a), VG(b)) \geq (d-1) \cdot 4d > (d-2) \cdot 4d + 3d$   $\checkmark$

Claim: Some LCS matches all 2's

Proof Sketch: Either 2 is skipped in both  $VG(a)$  and  $VG(b)$

$2 \ 1 \ 2 \dots 2$   
 $\downarrow \times \downarrow$   
 $2 \ 2 \dots$   
 Then, can increase subsequence by matching

Or 2 is skipped in only one of  $VG(a)$  and  $VG(b)$

$2 \ 2 \dots$   
 $\downarrow \ 2 \ \downarrow$   
 $2 \ 2 \dots$   
 Then we can also increase length of subsequ. by matching







Thus:  $\exists C'$  s.t.  $LCS(VG'(a), VG'(b)) = \begin{cases} C'+2 & \text{if } a \perp b \\ C' & \text{o.w.} \end{cases}$  ( $C' = C + 10d^2$ )

OR-Gadget: are there  $a \in A, b \in B$  s.t.  $a \perp b$ ?

Write  $A = \{a_1, \dots, a_n\}$   $B = \{b_1, \dots, b_n\}$

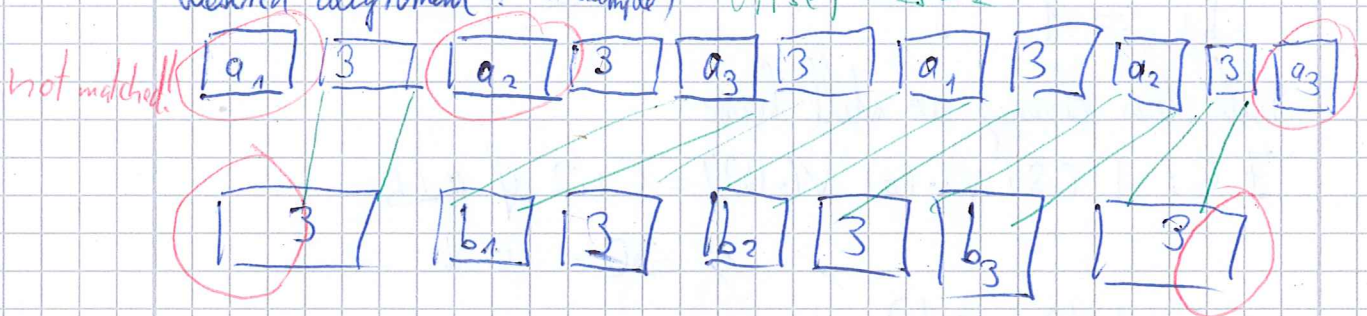
Want to construct two strings as follows:

every vector gadget appears twice

x:  $VG'(a_1) \underbrace{3 \dots 3}_{\text{fresh symbol } 100d^2} VG'(a_2) \underbrace{3 \dots 3}_{\text{length } 100d^2} \dots VG'(a_n) \underbrace{3 \dots 3}_{\text{length } 100d^2} VG'(a_1) \underbrace{3 \dots 3}_{\text{length } 100d^2} \dots VG'(a_n)$

y:  $\underbrace{33 \dots 3}_{\text{length } 100d^2} VG'(b_1) \underbrace{3 \dots 3}_{\text{every vector gadget appears once (additional padding at front and end)}} VG'(b_2) \underbrace{3 \dots 3}_{\text{length } 100d^2} \dots VG'(b_n) \underbrace{3 \dots 3}_{\text{length } 100d^2}$

Desired alignment: (example) offset  $\Delta = 2$



In general, if we align  $VG'(b_j)$  with  $VG'(a_{(\Delta+j) \bmod n})$  for some offset  $\Delta$

(will all be matched)

$$LCS(x, y) \geq \underbrace{\#3\text{'s in upper string}}_{(2n-1) \cdot 100d^2} + \sum_{j=1}^{\max_n} \underbrace{LCS(VG'(a_{(\Delta+j) \bmod n}), VG'(b_j))}_{\geq (n-1) \cdot C' + C' + 2 = nC' + 2}$$

$$= (2n-1) \cdot 100d^2 + nC' + 2$$

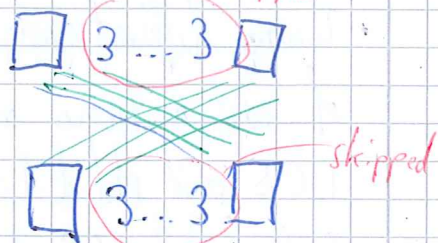
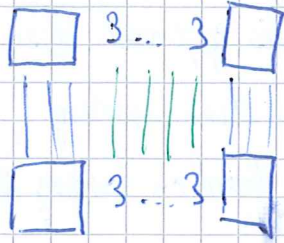
If there is an orthogonal pair, we can pick  $a^* \in A, b^* \in B$ , there is some offset  $\Delta$  s.t.  $VG'(a^*)$  and  $VG'(b^*)$  will be aligned (and then  $LCS(VG'(a^*), VG'(b^*)) = C' + 2$ )

Remains to show:

If there is no orthogonal pair, then  $LCS(x, y) \leq (2n-1)100d^2 + nC$



Observation: LCS of  $x$  and  $y$  will have no "crossings" skipped

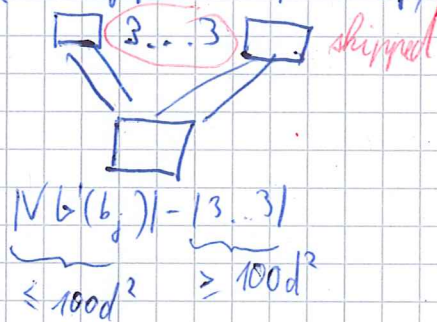


subsequence length  $\geq 100d^2$

length  $\leq 2 \cdot (C+2) < 100d^2$

We can now bound LCS in case of no orthogonal pair as follows:

$$LCS(x, y) \leq (2n-1)100d^2 + \sum_{j=1}^n \begin{cases} 0 & \text{if } VG'(b_j) \text{ is not matched} \\ C & \text{if } VG'(b_j) \text{ is matched to one } VG'(a_j) \\ |VG'(b_j) - \{3...3\}| & \text{if } VG'(b_j) \text{ is matched to several } VG'(a_j)'s \end{cases}$$



$\leq (2n-1)100d^2 + n \cdot C$

We have thus proved the following:

Lemma:  $A, B$  contains orthogonal pair iff  $LCS(x, y) = (2n-1)100d^2 + n(C+2)$

~~If  $A, B$  contains orth~~

More precisely: If  $A, B$  contains orthogonal pair, then  $LCS(x, y) = \dots + 2$

If  $A, B$  does not contain OP, then  $LCS(x, y) = \dots$

Time Complexity: ~~not~~  $|x| \leq |y| \leq O(d^2 \cdot n)$

(Time to construct  $|x|, |y|$ :  $O(d^2 \cdot n)$ )

Thus, if there were an  $O(n^{2-\delta})$ -algorithm for LCS for

( $n = \max\{|x|, |y|\}$ ), then there would be an

$O((d^2 \cdot n)^{2-\delta}) = O(n^{2-\delta} \text{ poly}(d))$  ( $n = |A| = |B|$ ) algorithm

for the OP Problem

Theorem: Assuming OVH, there is no  $O(n^{2-\delta})$  algorithm (for constant  $\delta > 0$ ) for the LCS problem.



## Remarks:

- Using different gadgets, ~~approx~~ this approach also shows hardness of edit distance and dynamic time wrapping
- LCS and edit distance are even hard on binary strings (alphabet:  $\{0, 1\}$ )