

Algorithms in Bioinformatics

Divide & Conquer



Recap:

- HW due
(feel free to submit Thurs.)
- Next HW - up by Thursday
- Today: Ch 7

Today: Divide + Conquer

Another variant of recursion:

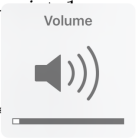
- Divide: usual recursion, but into halves, or n/c size subproblems
- Conquer:
↳ How to recombine

Classic example:

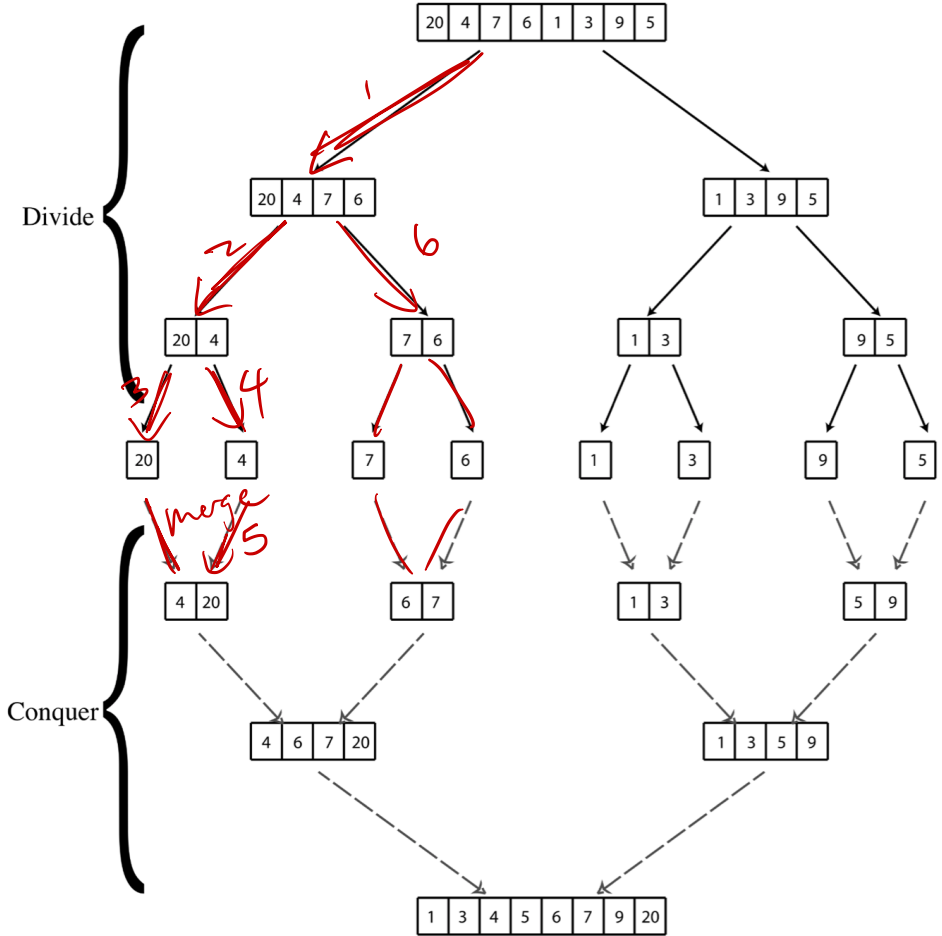
Merge Sort (video)

```
MERGESORT(c)
1  n ← size of c
2  if n = 1
3    return c
4  left ← list of first n/2 elements of c
5  right ← list of last n - n/2 elements of c
6  sortedLeft ← MERGESORT(left)
7  sortedRight ← MERGESORT(right)
8  sortedList ← MERGE(sortedLeft, sortedRight)
9  return sortedList
```

```
MERGE(a, b)
1  n1 ← size of a
2  n2 ← size of b
3  an1+1 ← ∞
4  bn2+1 ← ∞
5  i ← 1
6  j ← 1
7  for k ← 1 to n1 + n2
8    if ai < bj
9      ck ← ai
10     i ← i + 1
11   else
12     ck ← bj
13     j ← j + 1
14  return c
```

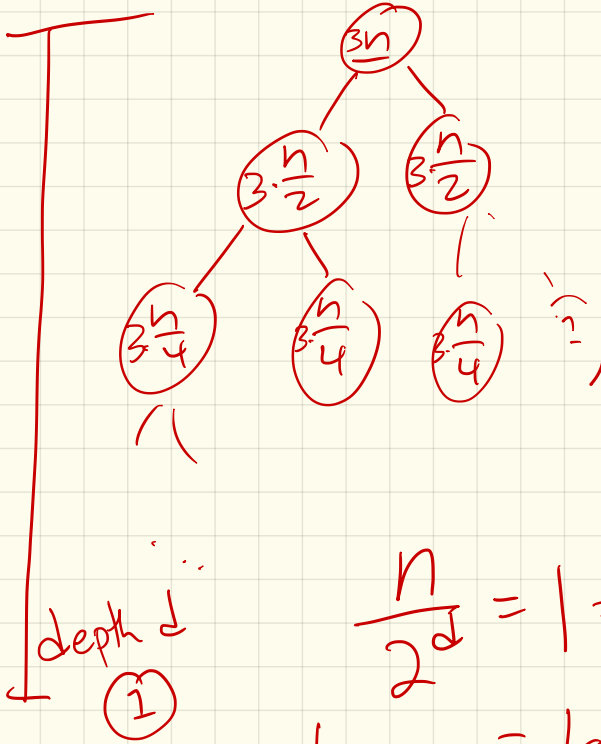


Picture :



Analysis:

$$T(n) = 2T\left(\frac{n}{2}\right) + \underline{3n}$$



amount of work per node on level i

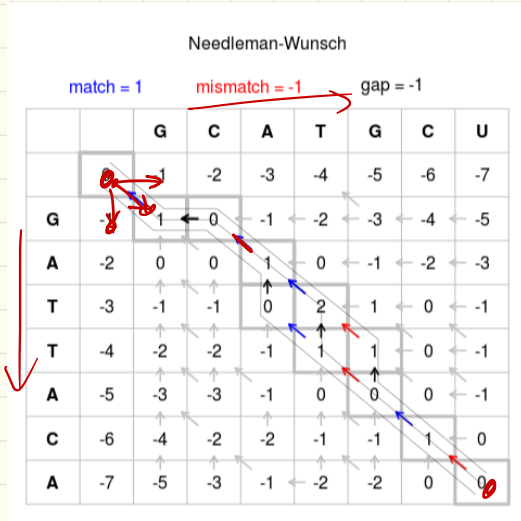
$$\sum_{i=0}^d \left(3 \cdot \frac{n}{2^i}\right) \cdot 2^i$$

\parallel

$$O(n \log n)$$

$$\frac{n}{2^d} = 1 \Rightarrow n = 2^d$$
$$\log_2 n = \log_2 2^d$$
$$= d \log_2 2$$

More biological: back to sequence alignment
 Recall our alignment structure



(Can use arbitrary scores)

Recursion:

$$S_{i,j} = \max \begin{cases} S_{i-1,j} + \delta(i^{\text{th}} \text{ item}, -) \\ S_{i,j-1} + \delta(-, j^{\text{th}}) \\ S_{i-1,j-1} + \text{penalty} \end{cases}$$

Runtime & space:

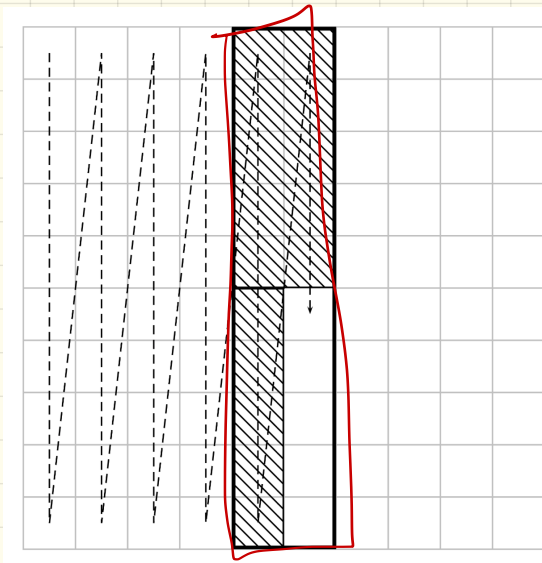
$$O(m \cdot n)$$

However: If just want score,
can get space down
to linear.

How?

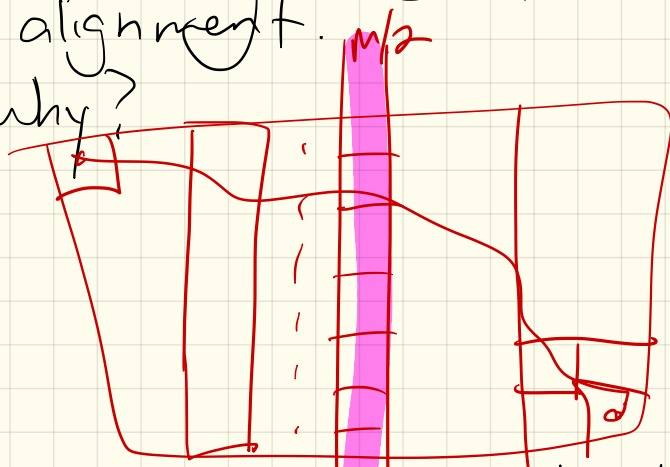
only need up, left,
& corner to fill
in a cell

⇒ column by column:
store 2 of them



Unfortunately - can't recover alignment.

Why?

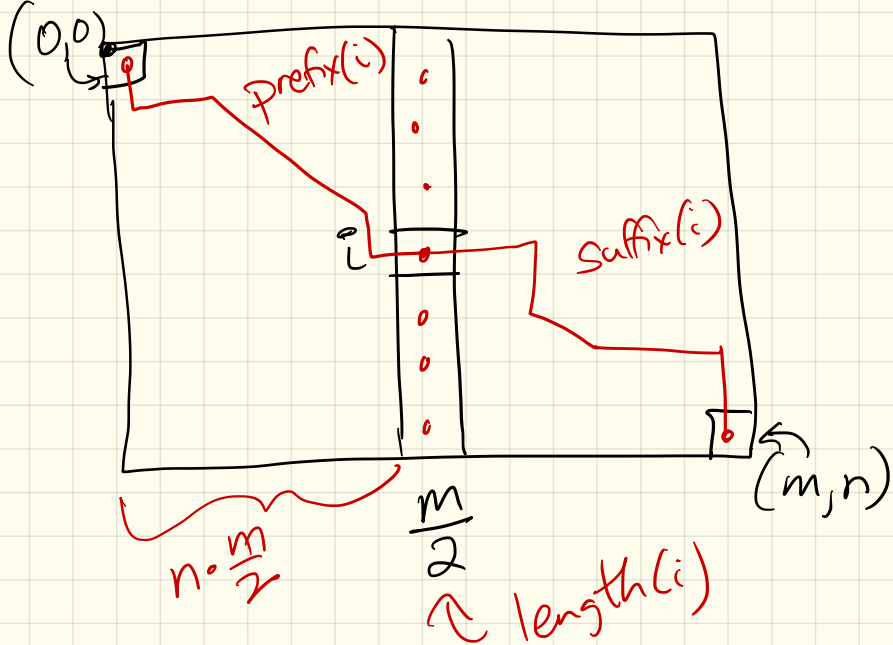


Think in terms of divide & conquer now:

define $\text{length}(i) :=$ cost of best alignment through cell $(i, \frac{m}{2})$

One of these must be the best path:

it goes from $(0, 0)$
to (m, n)



$$\text{length}(i) = \text{prefix}(i) + \text{suffix}(i)$$

How to get $\text{prefix}(i)$?

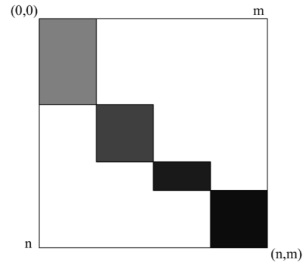
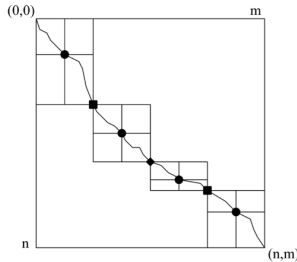
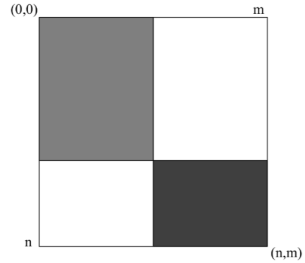
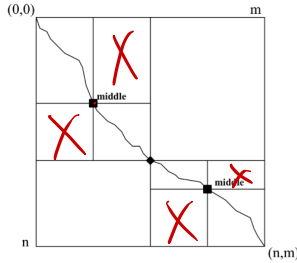
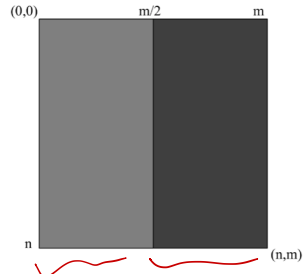
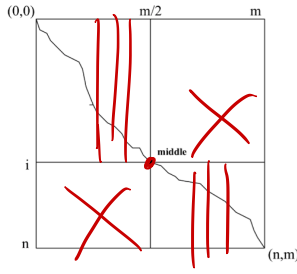
Chop off right half
+ run dynamic pro.

↳ $\frac{nm}{2}$ time, 2^n space

$\text{Suffix}(i)$: same, but "flip" the rectangle

Picture:

Linear-Space Sequence Alignment



Time: $O(mn)$
Space: $O(n)$

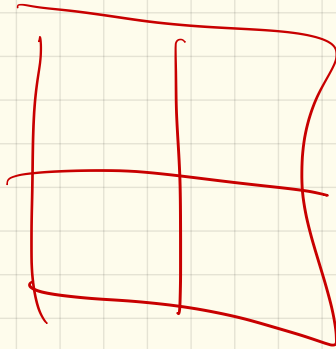
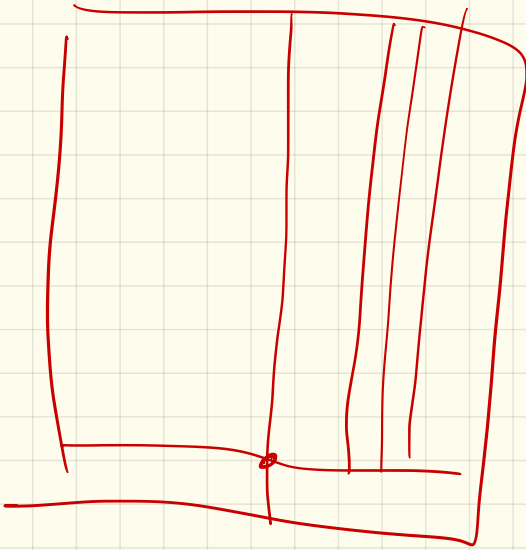
Runtime: mn time, $O(n)$ space

$$+ \frac{mn}{2}$$

$$+ \frac{mn}{4}$$

$$\leq 2mn$$

$$\sum_{i=0}^{\infty} \frac{1}{2^i} < 2$$
 keep using π



Block alignment (7.3)

Consider \vec{u} & \vec{v} , 2 DNA sequences, divided into blocks of length t

$$u = \underbrace{u_1 u_2 \dots u_t}_{\text{block 1}} \underbrace{u_{t+1} \dots u_{2t}}_{\text{block 2}} u_{2t+1} \dots u_{k \cdot t}$$

$$v = \underbrace{v_1 \dots v_t}_{\text{block 1}} \underbrace{v_{t+1} \dots v_{2t}}_{\text{block 2}} \dots v_{k \cdot t}$$

(ie if $t=3$, divide into codons)

We will seek an alignment where every block j in \vec{u} is either aligned with another entire block, or is inserted/deleted altogether

Goal: Find best such alignment.

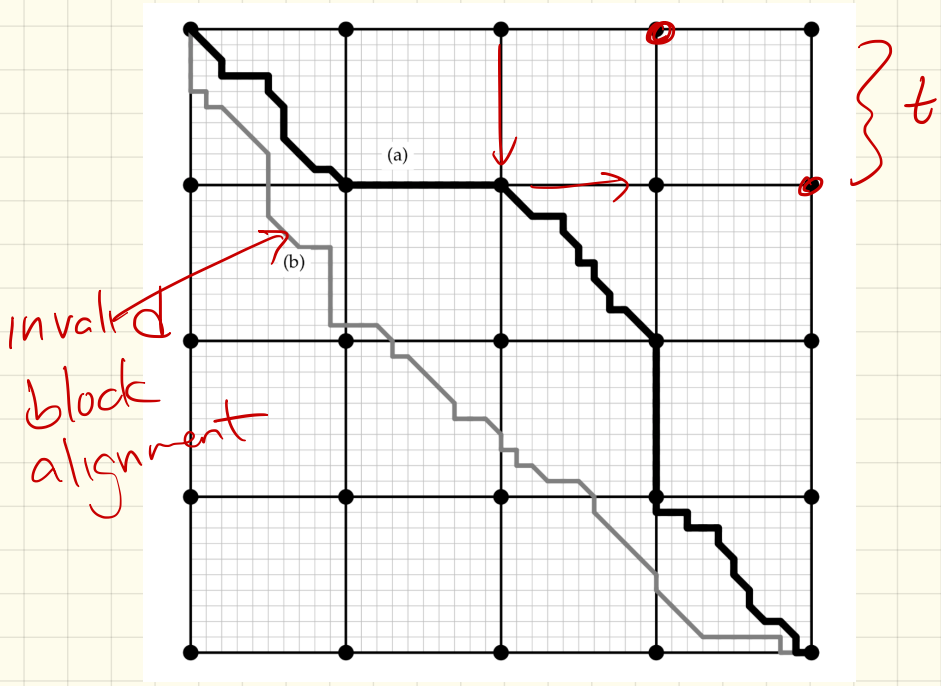
Block Alignment Problem:

Find the longest block path through an edit graph.

Input: Two sequences, u and v partitioned into blocks of size t .

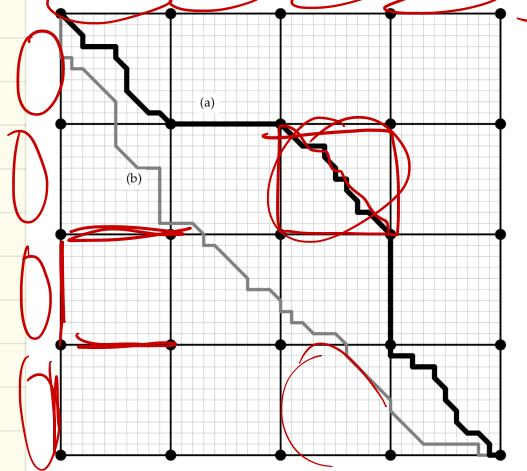
Output: The block alignment of u and v with the maximum score (i.e., the longest block path through the edit graph).

Picture: $n = m = 40$
 $t = 10$



Solution: Try each pair of blocks + compute alignment score $\beta_{i,j}$

↳ calling global alignment



$$\frac{n}{t} \times \frac{n}{t}$$

time

$$= O\left(\frac{n^2}{t^2}\right)$$

Then $S_{i,j}^* = \max \left\{ \begin{array}{l} S_{i-1,j} - \sigma_{\text{block}} \\ S_{i,j-1} - \sigma_{\text{block}} \\ S_{i-1,j-1} - \beta_{i,j} \end{array} \right.$

↑
goes from 1 to $\frac{n}{t}$

Runtime:

each block $\beta_{i,j}$: $t \times t$ block
 $\Rightarrow O(t^2)$

$\frac{n}{t} \times \frac{n}{t}$ pairs
 $\Rightarrow O\left(\frac{n^2}{t^2}\right)$

$\Rightarrow O\left(\frac{n^2}{t^2}\right) \cdot O(t^2)$
 $\Rightarrow O(n^2)$

Interesting speedup: (4-Russians technique)

Let $t = \log n$

Don't compute all $\frac{n}{t} \times \frac{n}{t}$ subproblems

Instead, consider all possible length t nucleotides

↳ $4^t \times 4^t$ total pairs

Store the cost of aligning in a table

Runtime: If $t \approx \frac{\log n}{\sqrt{4}}$,

$$4^t \times 4^t = n^{1/2} \times n^{1/2} = O(n)$$

Resulting table, Score, has only $4^t \times 4^t$ entries:

$O(n)$ size

Runtime:

- $O(n)$ entries in table
- each table entry takes $O(\log^2 n)$ time

Then same dynamic program:

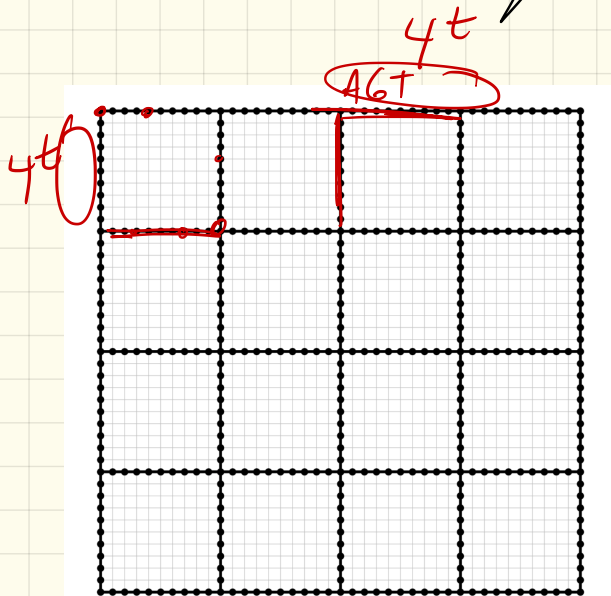
$$s_{i,j} = \max \begin{cases} s_{i-1,j} - \sigma_{block} \\ s_{i,j-1} - \sigma_{block} \\ s_{i-1,j-1} + \text{Score}(i\text{th block of } v, j\text{th block of } u) \end{cases}$$

Then $\frac{n}{t} \times \frac{n}{t}$ time for dynamic program, plus $O(\log n)$ lookup time in table (since need to find where it is stored)

Final runtime: $O\left(\frac{n^2}{\log n}\right)$

This can actually be used for
LCS problem, also.

Idea: break up into $t \times t$
squares again



n/t rows

n/t columns

total:

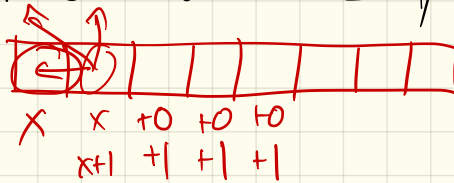
$$\frac{n^2}{t}$$

(can't just use top left &
bottom right corner
on each)

Brute force:

$4^t \times 4^t$ sequences,
and all possible scores
for 1st column & row

1st row: what is possible?



So: encode as binary
vector instead, 2^t

Size of table:

$$4^t \times 4^t \times 2^t \times 2^t = O(n^2)$$

$\hookrightarrow O\left(\frac{n^2}{\log n}\right)$ ✗