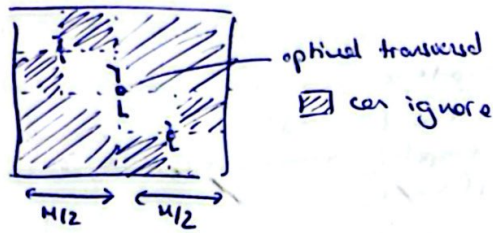


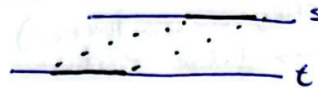
Lecture 3 Database search

Continuation of traceback from L2: Finding optimal path using linear space



Local Alignment

A local alignment of string s and t is an alignment of a substring of s to a substring of t .



Applications:

small domains may be the only conserved search rearrangements

Global Algn. (Needleman-Wunsch)

Init: $F(0,0) = 0$

It:

$$F(i,j) = \max \begin{cases} F(i-1,j) - d \\ F(i,j-1) - d \\ F(i-1,j-1) + s(x_i, y_j) \end{cases}$$

Ter: \searrow

Local Algn. (Smith-Waterman)

Init: $F(0,j) = F(i,0) = 0$

It:

$$F(i,j) = \max \begin{cases} 0 \\ F(i-1,j) - d \\ F(i,j-1) - d \\ F(i-1,j-1) + s(x_i, y_j) \end{cases}$$

Ter: anywhere

Semi global

• no penalization for terminations

Init: $F(0,j) = c$
 $F(i,0) = c$

It: $\max \begin{cases} F(i-1,j) - d \\ F(i,j-1) - d \\ F(i-1,j-1) + s(x_i, y_j) \end{cases}$

Ter: $\downarrow \rightarrow$

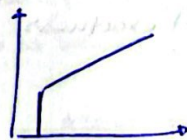
Gap penalty functions

Generalized: Linear gap penalty: pay same for each gap.

• exponential: encodes the length of the gap $O(w)$ and pays less as it grows!



• Affine gap: start gap (price = p) no (price = q) $w(k) = p + q \cdot k$



Need another matrix to see whether you risk in a gap or not.

• Lengths (mod 3) gap penalty for protein-coding regions

↳ prefer leaving 1aa than disrupting the reading frame

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

or

A A A - - - - T C A is more likely

Linear-time string matching

- Looking for exact matches (no gaps)
- Karp-Rabin algorithm (probabilistic linear time)
→ Interpret string numerically

T: 2 3 5 9 0 2 3 1 4 1 5 2 6 7 3 9 9
 \downarrow $u_2 = 35,902$ \downarrow y_7
 $y_1 = 23,590$ $x = y_7$

P: 3 1 4 1 5 → $x = 31,415$

- The u 's are too big → take (modulo) to make computation easier.
(hashing: one way function)
- Taking the mod → collisions → false positives.
- So when I have a hit → go and check!

Computing t_{s+1} from t_s in constant time

$\boxed{13 \mid 1 \mid 4 \mid 1 \mid 5 \mid 2 \mid}$

$y_1 = 31,415$ → $y_2 = 14,152$

old digit

$$14,152 = (31,415 - 3 \times 10,000) \times 10 + 2$$

$$14,152 \bmod 13 = \boxed{8}$$

new digit
1
left shift

The BLAST algorithm and inexact matching

Search algorithm!

Assumes common ancestry

Finds alignment

Evo interpretation → $\min(\text{Cost}(\pi))$

Query must be very fast! → most seq. will be unrelated → then analyzed.

1. reject idpair ≤ 90
why bother?

2. Speed in preprocessing (offline)
are query arrives → fast!

3. Context-based hashing
index 10-mers
→ only $1/4^{10}$ will match
→ 1 in a million

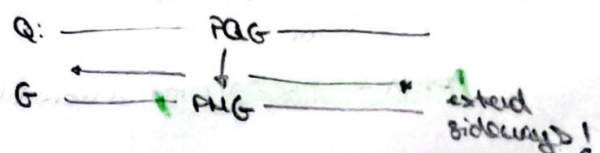
• Hashing

• Neighborhood search → find matches regardless of exactness

Overview

1. Query into overlapping words of length w
2. Neighborhood words for each word until threshold T
3. Look in table when neighbor words occur: seeds S
4. Extend seeds S until score $< X$

Query: ... P A G ...
 $w=3$



Given $a, b \in \mathbb{R}$:

$a \bmod b$ is the remainder of $\frac{a}{b}$

compute x

for i in $[1 \dots n]$:

compute y_i ← not constant, linear!

if $x == y_i$:

Print "match at $S[i]$ "

we need to compute y_i based on y_{i-1}

at least are hole with no pigeon.

BLAST works because \rightarrow Pigeonhole principle $\left[\begin{smallmatrix} R \\ \vdots \\ Q \end{smallmatrix} \right] \left[\begin{smallmatrix} \vdots \\ \vdots \\ \vdots \end{smallmatrix} \right] \left[\begin{smallmatrix} Q \\ \vdots \\ \vdots \end{smallmatrix} \right]$

- 1) Filter out low complexity regions (NNN...) or very common subsequences (uninformative)
- 2) Two-hit BLAST \rightarrow two smaller W-mers more likely than one larger
- 3) Non conservative k-mers (cubs)

RG-1KW \rightarrow R * 1 k *, R * k * W

\hookrightarrow Random projections \rightarrow n dimensions \rightarrow n-k dimensions

Probabilistic foundations of seq. alignment

Assigning scores/penalties for matches/mismatches

Nucleotide space:

	A	G	T	C
A	1	-1/2	-1	-1
G	-1/2	1	-1	-1
T	-1	-1	1	-1/2
C	-1	-1	-1/2	1

Protein space: BLOSUM matrix of sim scores

The score should represent:

Likelihood ratio $\left\{ \begin{array}{l} H_1: \text{random alignment} \rightarrow \text{chance} \\ H_2: \text{common ancestry} \end{array} \right.$

$$\frac{P(x, y | H_1)}{P(x, y | H_2)} \quad \left\{ \begin{array}{l} P(x, y | H_1) \\ P(x, y | H_2) \end{array} \right. \quad P(\text{homology}) = \frac{P(x, y | H_2)}{P(x, y | H_1)}$$

$$\text{Score} = -\log(P(\text{homology}))$$

★ additive metrics are exactly that:
• take logs for each aa and add.

$$P_r(x, y | U) = \prod_{i=1}^n q_{x_i} \prod_{i=1}^n q_{y_i}$$

• unrelated

$$P_r(x, y | R) = \prod_{i=1}^n p_{x_i y_i}$$

• related

• how can we decide U or R?

\rightarrow odds ratio

$$\frac{P_r(x, y | R)}{P_r(x, y | U)} = \frac{\prod_i p_{x_i y_i}}{\prod_i q_{x_i} \prod_i q_{y_i}} = \frac{\prod_i p_{x_i y_i}}{\prod_i q_{x_i} q_{y_i}}$$

$$\log\left(\frac{P_r(x, y | R)}{P_r(x, y | U)}\right) = \sum_i \log\left(\frac{p_{x_i y_i}}{q_{x_i} q_{y_i}}\right) = S$$

Bayes Theorem

$$P(+1 | x, y) = \frac{P(x, y | R) P(R)}{P(x, y | R) P(R) + P(x, y | U) P(U)}$$

$$L = \sigma(S')$$

$$S' = \frac{e^S}{1 + e^S} ; S' = S + \log \frac{P(R)}{P(U)}$$

$$O \sim O(n) \rightarrow O(n^2)$$

Deterministic linear-time exact string matching

Input: P: pattern
T: text

\rightarrow analyze int. redundancy

\rightarrow as soon as mismatch \rightarrow shift

e.g.

P = a b a

T = b a a b a c a b a b a d

* Make bigger shift!