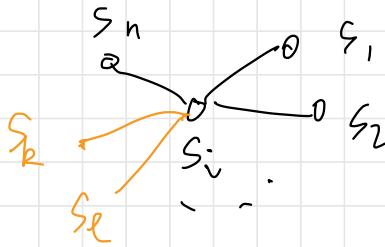

Seeding Methods in Homology Search

Review:

- ① multiple alignment (exact) algorithm.
- ② Relative entropy score for each column.

$$\sum_{a \in \Sigma} n(a) \cdot \log \frac{n(a)/n}{q(a)}$$

- ③ Approximation algorithm for SP-score
with Triangular Inequality.



$$\begin{aligned} d_i(s_k, s_l) &\leq d(s_i, s_k) + d(s_i, s_l) \\ &\leq d^*(s_i, s_k) + d^*(s_i, s_l) \end{aligned}$$

A similarity between mouse and human genomes

HSP.

```
GCNTACACGTCACCATCTGTGCCACCACNCATGTCTCTAGTGATCCCTCATAAGTTCCAACAAAGTTTGC
|| |||| | ||| ||| | | ||||| ||||| | ||||| | | ||||
GCCTACACACCGCCAGTTGTG-TTCCTGCTATGTCTCTAGTGATCCCTGAAAAGTTCCAGCGTATTTTGC

GAGTACTCAACACCAACATTGATGGGCAATGGAAAATAGCCTTCGCCATCACACCATTAAGGGTGA----
|| ||||| ||||| | |||| | ||||| || ||||| | | |||
GAATACTCAACAGCAACATCAACGGGCAGCAGAAAATAGGCTTTGCCATCACTGCCATTAAGGATGTGGG

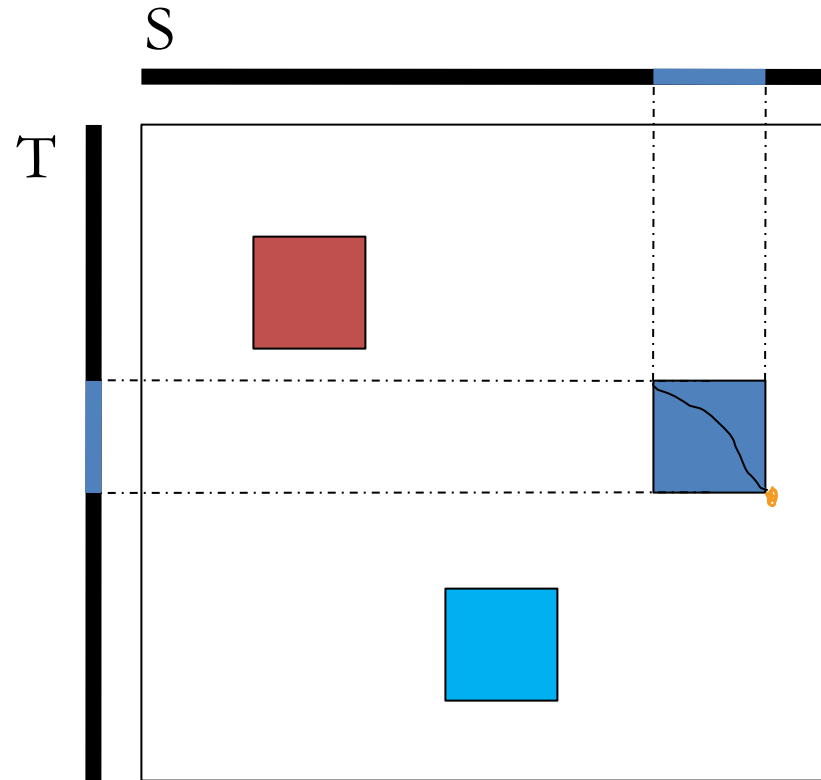
-----TGTTGAGGAAAGCAGACATTGACCTCACCGAGAGGGCAGGCGAGCTCAGGTA
|| ||||| ||||| || ||||| || ||||| || |||| |
TTGACAGTACACTCATAGTGTGAGGAAAGCTGACGTTGACCTACCAAGTGGGCAGGAGAACTCACTGA

GGATGAGGTGGAGCATATGATCACCATCATAAGAACTCAC-----CAAGATTCCAGACTGGTTCTTG
|| |||| | | |||| |||| || |||| | |||| | ||||| |||||
GGATGAGATGGAACGTGTGATGACCATTATGCAGAATCCATGCCAGTACAAGATCCAGACTGGTTCTTG
```

Smith-Waterman is the most accurate method.

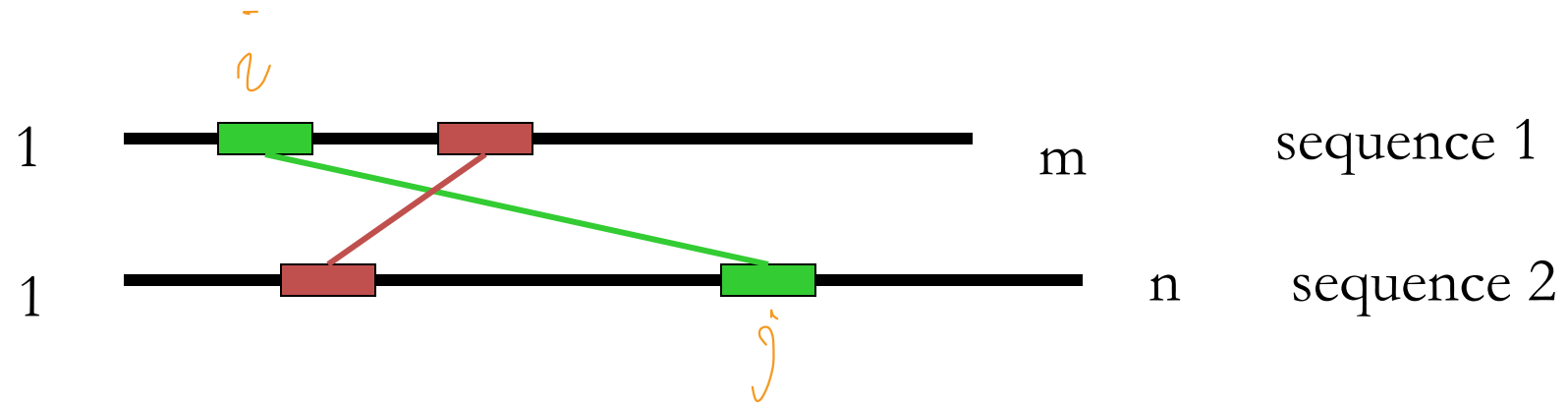
Time complexity: $O(mn)$.

Smith-Waterman Algorithm



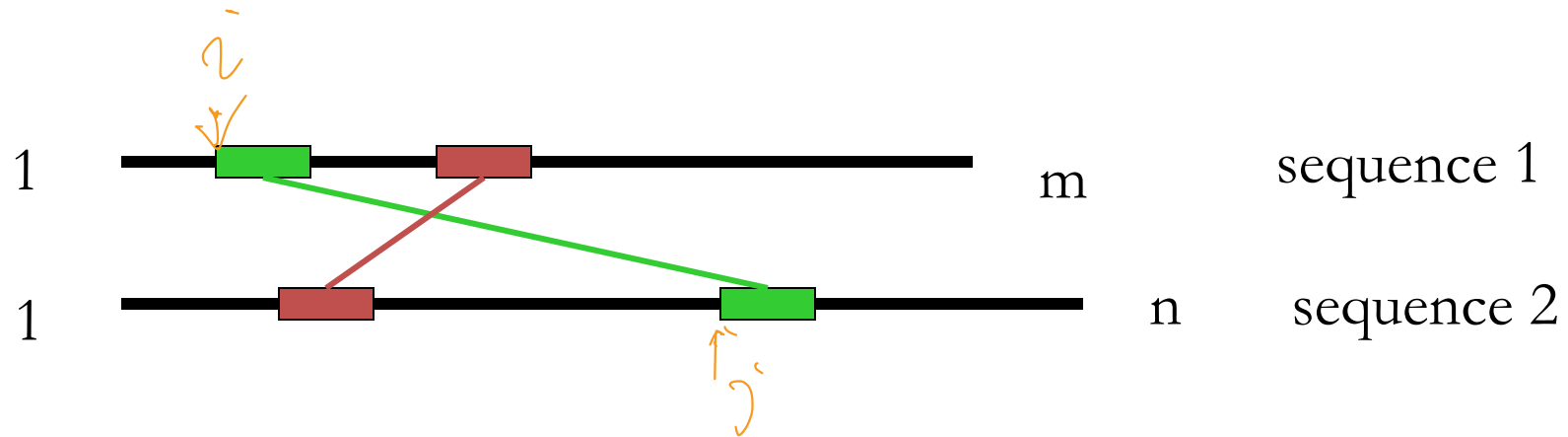
- The old algorithm requires $O(mn)$ and is too slow.
- Human v.s. mouse: $3 \times 10^9 \times 3 \times 10^9 = 9 \times 10^{18}$

Similarity Search



- Most similarities (local alignments) are very short relative to the genomes.

Similarity Search



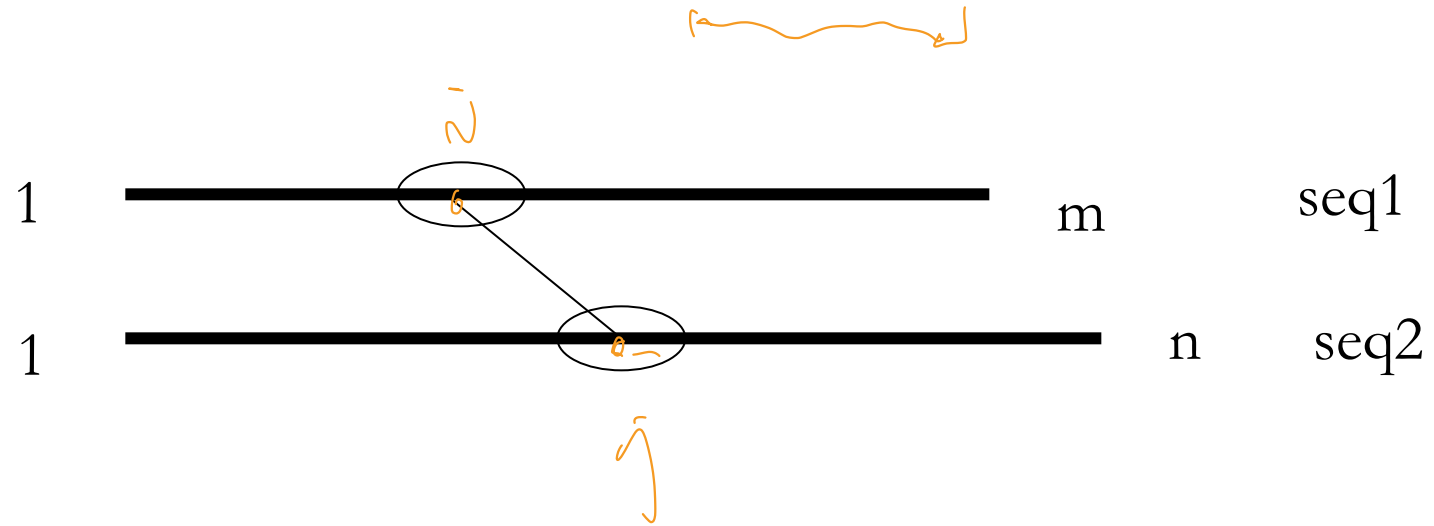
- For every pairs of (i, j) , build a local alignment around it.
 - $O(mnT)$
 - Not better than Smith-Waterman.
- But this leads to an important idea...

Main Idea

- Most pairs of (i, j) are useless. We only want to try local alignments on the “promising” pairs of (i, j) .
- In the context of sequence similarity search in bioinformatics, these “promising” pairs are called “seeds” or “hits”.
- We need
 - a proper definition of hits.
 - some efficient way to enumerate the hits faster than trying every pair of (i, j) .

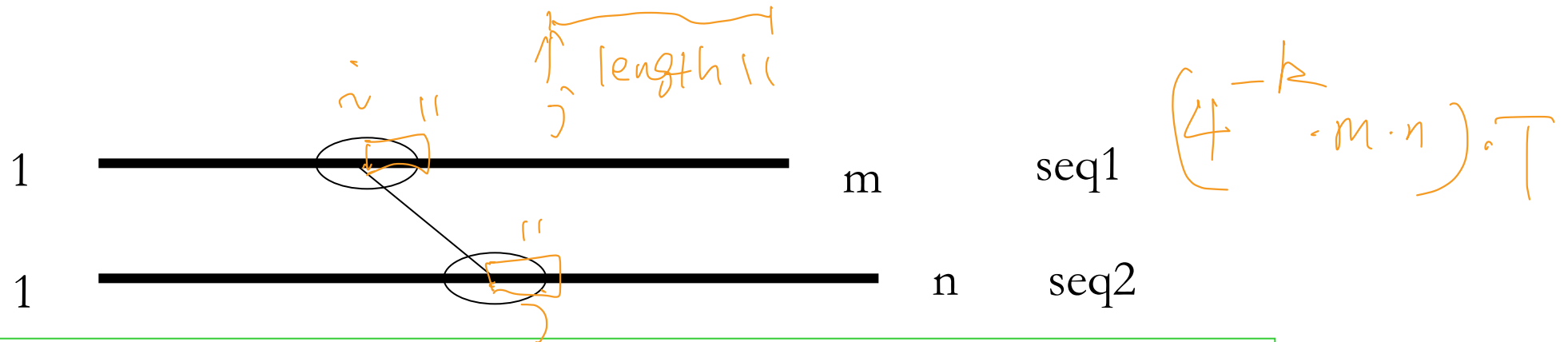
BLAST Uses Short Consecutive Match as Hits

GCNTACACGTCACCATCTGTGCCACCACNC**ATGTCTCTAGT**GATCCCTCATAAGTTCCAACAAAGTTTGC
|| |||| | ||| |||| | | ||||| ||||| | ||||| | | |||||
GCCTACACACCGCCAGTTGTG-TTCCTGCT**ATGTCTCTAGT**GATCCCTGAAAAGTTCCAGCGTATTTTGC



BLAST Uses Short Consecutive Match as Hits

GCNTACACGTCACCATCTGTGCCACCACNC**ATGTCTCTAGT**GATCCCTCATAAGTTCCAACAAAGTTTGC
 || |||| | ||| |||| | | ||||| ||||| | ||||| | |||||
 GCCTACACACCGCCAGTTGTG-TTCCTGCT**ATGTCTCTAGT**GATCCCTGAAAAGTTCCAGCGTATTTTGC



- Majority of (i,j) are random and probability of generating a random hit is small.
- For length-k seed, time complexity becomes $O(4^{-k}mnT)$
- By default, BLAST used $k=11$.
- What's the speed up factor for $k=11$?

$4^{-11} \approx \frac{1}{4 \text{ million}}$

The Idea behind Seeding

- A true similarity has a high chance of being hit. (Hitting at any position in the similarity will do.)
- A random pair (i, j) has low chance of being hit.
- Thus, if we use hit to filter (i, j) , we will
 - Detect most true similarities.
 - Not wasting time on random pairs of (i, j) .

The Data Structure for Finding Hit?

- for each k -mer, index table to remember all its occurrences in S.
- for each k -mer of T, find its hits in the index table.
- The index table can be a trie or a hash table.

$k=2$

AA	→	0, 6
AC		
AG		
AT	→	1
CA		
CC		
CG		
CT	→	3
GA		
GC		
GG		
GT		
TA	→	5
TC	→	2
TG		
TT	→	4

S: AATCTTAA
 01234567 ← i

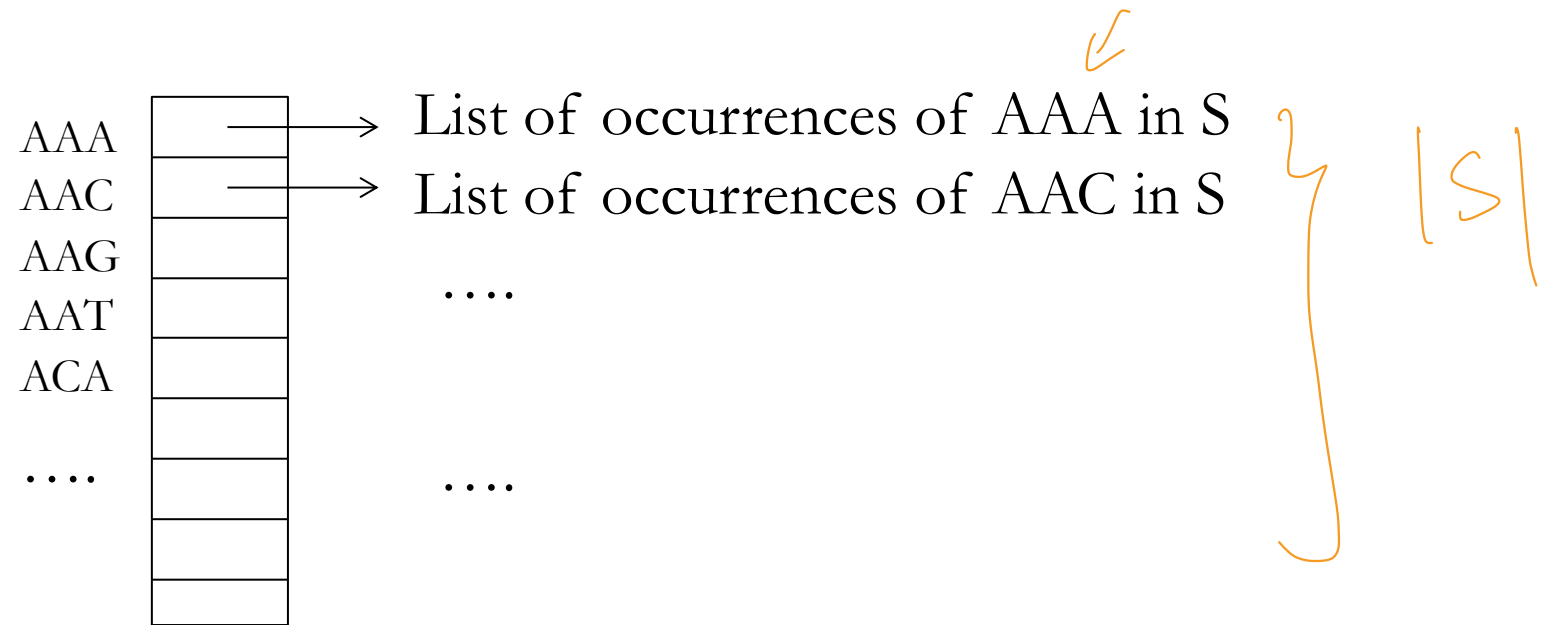
T: GAACTTA
 ↑ ↑
 i j

$$O(m+n+k)$$

↑
of hits

i, j
 $(0, 2)$
 $(6, 2)$

The Data Structure for Finding Hit?



Space complexity?

$$4^k$$

$$4^k \approx 4M, \quad k=11$$

Overall runtime

- Build the index using S : $O(n)$ time.
- Find matches between the index and sequence T : $O(m)$ time to scan T , plus we need to examine all of the N hits found. Let t be the examination time. Then $O(m + Nt)$.
- Overall runtime: $O(n + m + Nt)$.
*Another round of filtration:
 $N \cdot t' + N' \cdot t$*
- The term Nt is the most expensive part. Indexing overhead is small.
- In practice, most of the hits encountered are random hits.

$$(4^{-k} \cdot m \cdot n)$$

HSP extension

```

GCNTACACGTCACCATCTGTGCCACCAGCCATGTCCTAGTGATCCCTCATGGTGGCCAACAAAGTTTGC
      |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
TGCCTACACACCGCCAGTTGTGTTCCCTGCTATGTCCTAGTTATCCCTGAAAAGTTCCAGCGTATTTTGC
    
```

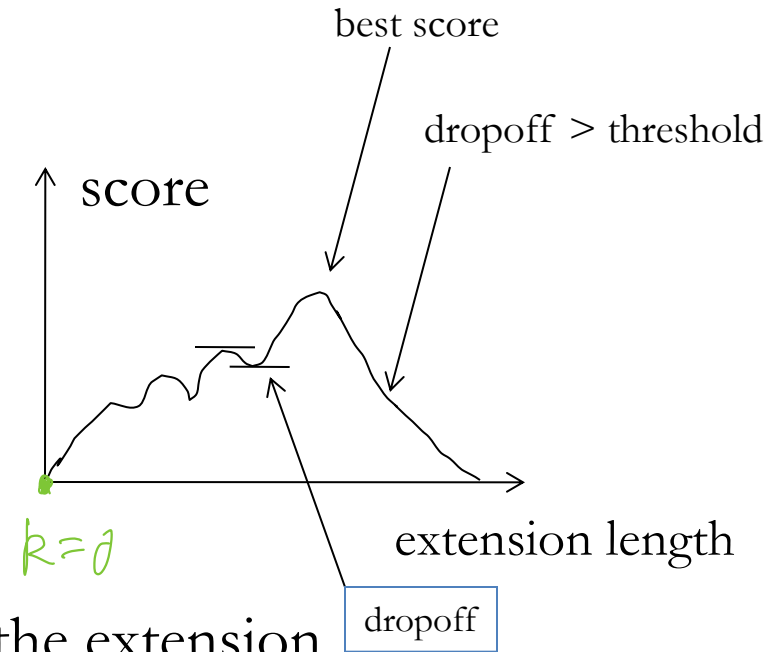
for k from 0 to ...

score += sc(S[i+k],T[j+k])

for k from 1 to ...

score += sc(S[i-k],T[j-k])

- But when to stop?
- Score will increase and decrease during the extension.
- Extension stops when drop off greater than threshold.



HSP Extension

- How long will the extension continue after reaching best score?
- Assumptions:
 - After reaching best score, sequence becomes random.
 - match=1 and mismatch=-1
- Expected score on each additional base is -0.5.
- If dropoff=k, then after 2k bases, the expected dropoff will reach k.
- Conclusion: Not too long.

prob	score
$\frac{1}{4}$	1
$\frac{3}{4}$	-1

Example of missing a target

R

- Fail:

```
GAGTACTCAACACCAACATTAGTGGGCAATGGAAAAT
|| (|||||) ||||| | |||||   |||||
GAATACTCAACAGCAACATCAATGGGCAGCAGAAAAT
```

- Dilemma
 - **Sensitivity** – needs shorter seeds
 - the success rate of finding a homology
 - **Speed** – needs longer seeds
 - Mega-BLAST uses seeds of length 28.

PatternHunter uses “spaced seeds”

- 111*1**1*1**11*111 (called a spaced seed)
 - Eleven required matches (**weight**=11, **length** = 18)
 - Seven “don’t care” positions



```
GAGTACTCAACACCAACATTAGTGGCAATGGAAAAT...
|| ||||| ||||| || ||||| |||||
GAATACTCAACAGCAACACTAATGGCAGCAGAAAAT...
111*1**1*1**11*111
```

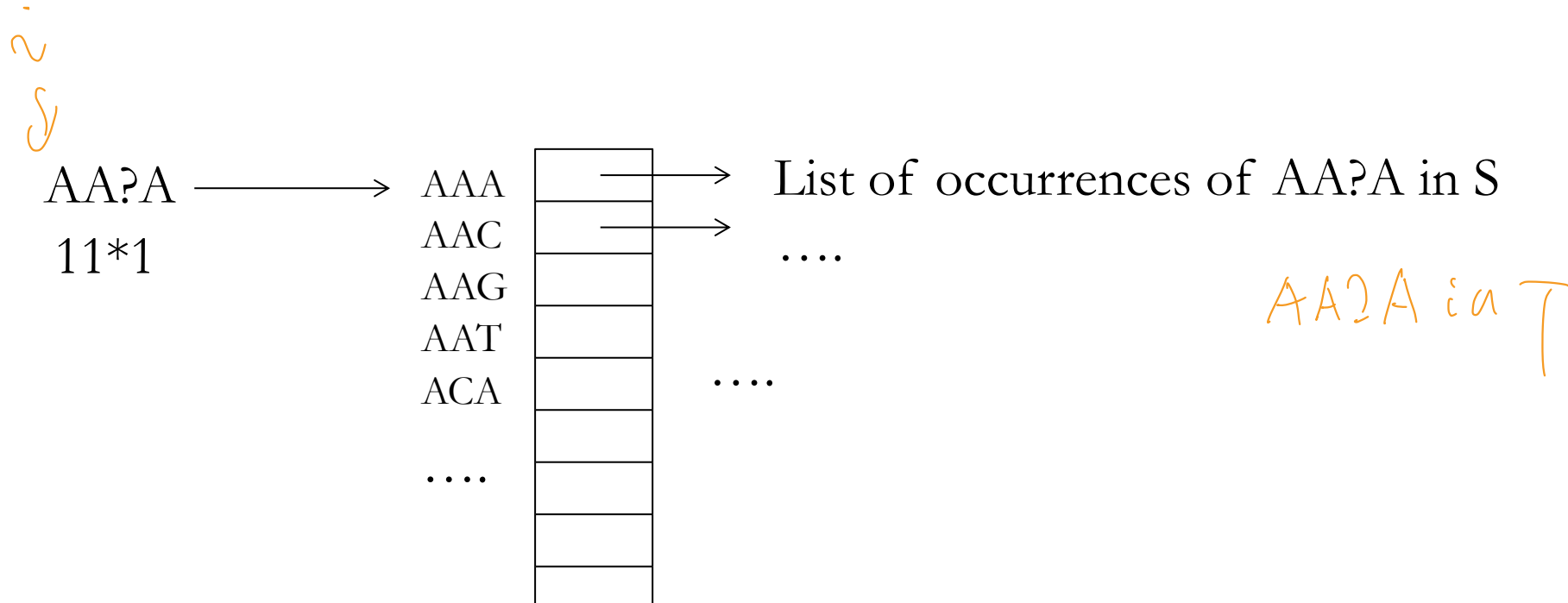
- Hit = all the required matches are satisfied.
- BLAST’s seed = 11111111111

Notes about the notation

- A homology/similarity region's actual sequences do not matter, the match/mismatch matters.
- Therefore, a region is often denoted by a binary 0-1 sequence,
1101111100111011101111
- A hit is then as follows:
- **1101111100111011101111**
- **111*1**1*1**11*111**

The Data Structure for Finding Hit

- The same as consecutive seed. Except that now we have a length l , weight w seed. E.g. 11^*1 .
 - Each l -mer, take the w letters out and put in index table.
- The index table can be a hash table.



Time Complexity Comparison

- Lemma: for random sequence S and T with lengths m and n , the expected number of hits for weight w , length l seed is

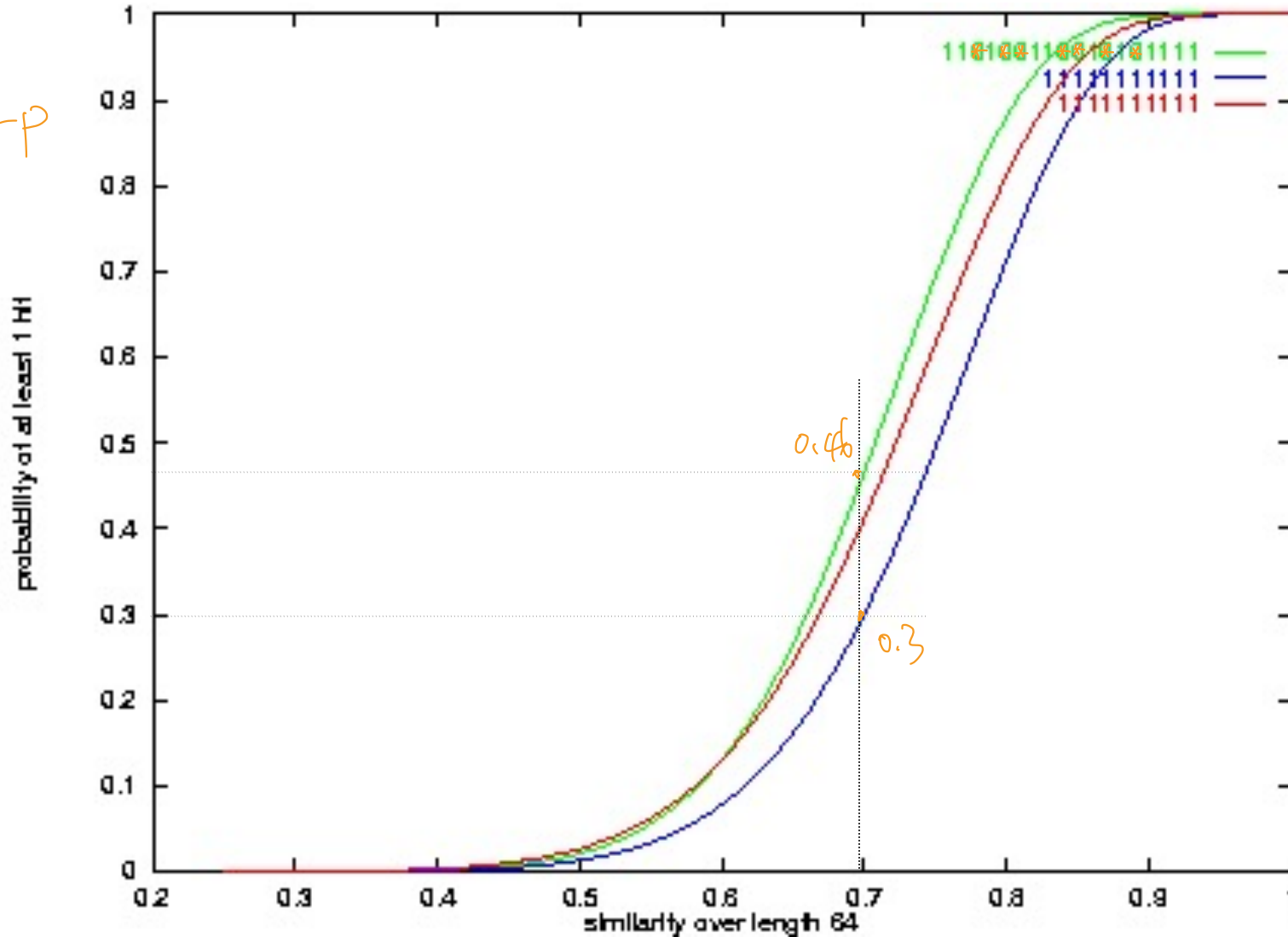
$$(m - l + 1)(n - l + 1)4^{-w}$$

- Because usually l is much shorter than S and T, this is approximately $4^{-w}mn$
- That is, the expected number of hits in random regions only depends on the weight, but not the shape of the seed. So does the running time.
- So, speed-wise, spaced seed is similar to consecutive seed.
- What about the sensitivity?

64
HSP

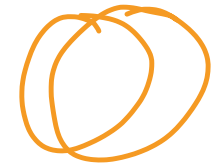
Simulated sensitivity curves

$\text{prob}(1) = p$
 $\text{prob}(0) = 1 - p$



Similarity level ϕ

Why spaced seeds are better?

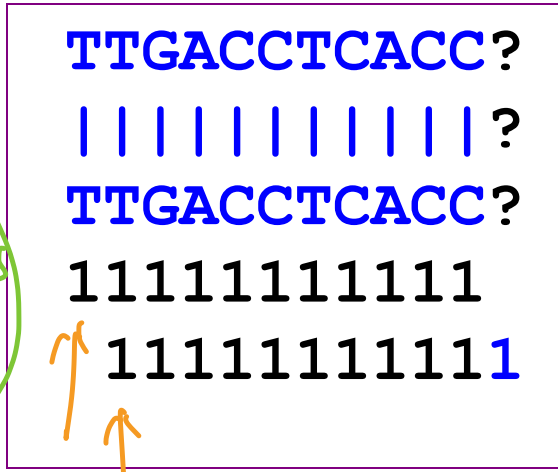


HSP.

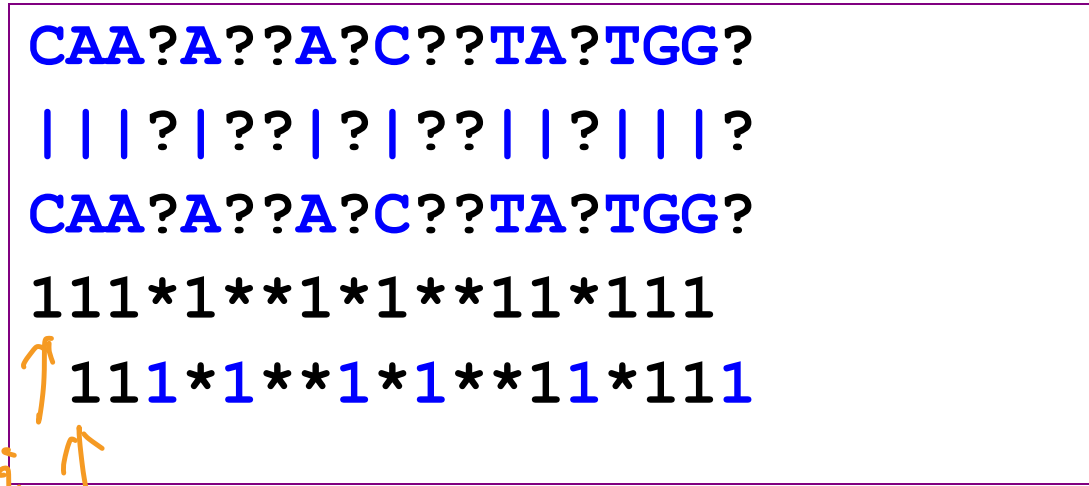


p''

P_r (HSP containing
at least 1 hit)



i $i+1$



i $i+1$

- BLAST's seed usually uses more than one hits to detect one homology (redundant)
- Spaced seeds uses fewer hits to detect one homology (efficient)

