

How to implement that?

- With BLASTP:
 - Build an automaton that reflects all string close to short strings in T (the short sequence)
 - Scan S (the longer sequence), looking for matches.
- We do not study the classic ways to match multiple patterns efficiently.
If interested, you can read at
https://en.wikipedia.org/wiki/Aho%E2%80%93Corasick_algorithm

A Simpler Way

- There is another way:
 - 1) For every 3-mer, find all “neighboring” 3-mers that, score at least +11 (or whatever). Build these into a data structure NeighborList.
 - 2) Build a hash table H for S of its 3-mers, just like for the nucleotide case
 - 3) For every 3-mer x in T, retrieve all neighbors from NeighborList. For each neighbor, query H to find hits in S.

NeighborList is a small structure: there are only 8000 3-mers.

Which sequence to index?

- That's actually a tough question.
- Here's a typical scenario:
 - S is the human genome (length n)
 - P_1 is a short protein sequence (length m_1)
 - P_2 is another short sequence (length m_2)
- If we're smart, build an index for S, *once*, and then look up the short sequences in it.
- Added time for P_2 is more like $O(m_2)$, not $O(n+m_2)$.

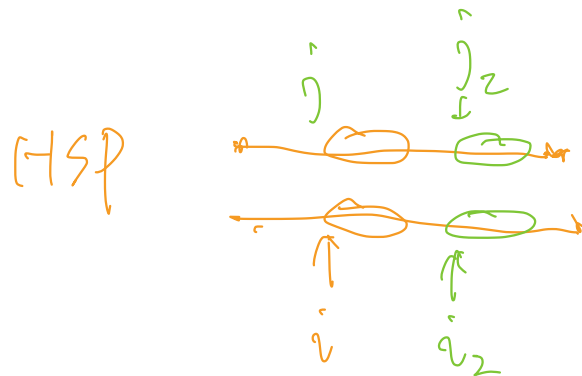
More on indexing

- But memory is a concern:
- Indexing the human genome is expensive!
- Oh, wait. No, it isn't, not anymore... you probably should index the longer sequence.
- BLASTN (1990) indexes the query, not the database.
- BLAT (2000) indexes the database, not the query.

- BLASTP also indexes the query.

Extensions to this idea

- Two-hit BLAST:
- Require two seeds (probably shorter) that are nearer than k from each other, and base the alignment on their enclosing box.
- Potentially even fewer false positives, but one has to use shorter seeds. There's quite a tradeoff here.



$$\hat{j} - i = \hat{j}_2 - i_2$$

Q&Q $|i_2 - i| \leq W$

Wrap-up

- Local alignment slow when sequences are large
- Use 11 consecutive matches as hits
 - How these hits are found efficiently
 - What to do after hits are found
- Spaced seeds better
 - How sensitivity is computed and how optimal seed is found
 - How hits are found for spaced seed
- Multiple spaced seed.
- Protein seeds.
- Two hits.