# Combinatorial Methods in Bioinformatics
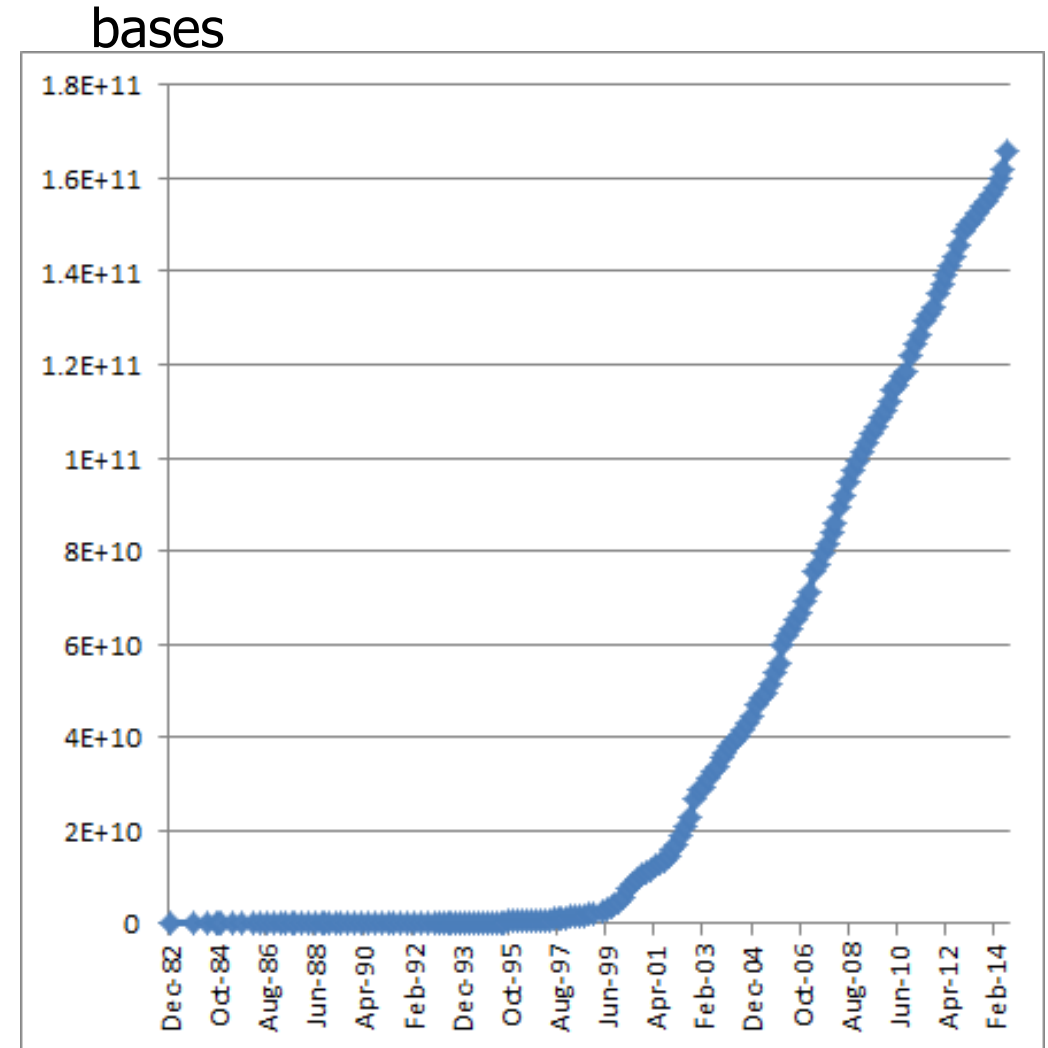
**Searching biological database**

Wing-Kin Sung, Ken 宋永健

ksung@comp.nus.edu.sg

# Biological database

- Biological database stores known sequences of proteins, RNAs and DNAs.

- There are many biological database. The most common database is genbank.

- The size of biological database increases exponentially.

- In Dec 2019, genbank stores 215 millions annotated genes and 388 billions base pairs.

bases

# Why searching biological database?

- Example:
  - In November 2005, an unknown virus from a child in an Amish community was suspected to be an intestinal virus.
  - Through database searching, the unknown sequence was found to be a polio virus.
  - This is the first polio case in the U.S. since 1999 and it demonstrates the power of database search.
- In molecular biology, it is common to search query sequences on various databases since sequence similarity often leads to functional similarity.
- For instance, in 2002, GenBank alone served 100,000 search queries daily which had risen to over 140,000 queries daily by 2004.

# Problem definition

- Consider a database D of genomic sequences (or protein sequences)

- Given a query string Q,
  - we look for a string S in D which is the closest match to the query string Q using some score function.
  - There are two choices of the score function:
    - **Semi-global alignment score**: the best possible alignment score between a substring A of S and Q
    - **Local alignment score**: the best possible alignment score between a substring A of S and a substring B of Q.

# Example

- Database D:
  - $S_1$=ACTCGGATGGT
  - $S_2$=CACCAGTACGCA
  - $S_3$=TCAGGTGGTA
- Query Q=GGTTGG
- Assume match scores +1, mismatch/indel scores -1

- Solution: There are two hits with the highest score 4 (same solutions for both local or semi-global alignment)
  - $S_1$=ACTCGGATGGT          GGTTGG
  - $S_2$=CACCAGTACGCA          GGATGG
  - $S_3$=TCAGGTGGTA

          GGTTGG
          GGT-GG

# Measurement of the goodness of a search algorithm

- Efficiency
  - Measure the running time of an algorithm.
- Sensitivity
  - Sensitivity is the ratio of the number of true positives found by the algorithm over the actual number of true positives.
  - An algorithm with 100% sensitivity finds all true positives.
- A good search algorithm should be both sensitive and specific

# Different approaches

- Exhaustive search algorithms
  - E.g. Smith-Waterman Algorithm and BWT-SW
- Heuristic algorithms
  - E.g. FastA, BLAST, BLAT and PatternHunter
- Filter-based algorithms
  - E.g. QUASAR and LSH

# Exhaustive search algorithm

- Input:
  - the database D (total length: n) and
  - the query Q (length: m)
- Output: all closest matches (based on local alignment or semi-global alignment)

**Algorithm**

- For every sequences S in the database,
  - Compute the best alignment between S and Q by dynamic programming
- Return all alignments with the best score

- Time: $\sum_{S \in D} m|S| = O(nm)$.
- This is a brute force algorithm. So, it is the most sensitive algorithm.

# What is FastA?

- Given a database D and a query Q,
  - FastA does local alignment with all database sequences and returns the database sequences with the highest local alignment score
  - Its assumption is that good local alignment should have some exact match subsequences.
- History of FastA
  - 1983: Wilbur-Lipman algorithm
  - 1985: FastP
  - 1988: FastA

# Example



The input is a substring of human Sox2 protein sequence

# Example



It can find the correct alignment!

# FastP

- Lipman and Pearson observed that the occurrences of indels are far less than that of replacements.

- Hence, FastP focused to identify high scoring alignments without indels.

- For each sequence S in the database D, the highest scoring ungapped alignment between S and Q is identified in three steps:

    1. Identification of hotspots
    2. Locating diagonal runs
    3. Rescoring the best diagonal runs

# Step 1 of FastP: Identification of hotspots

- Every k-tuple (length-k substring) of a database sequence is called a hotspot if the k-tuple appears in the query sequence.

- The larger the value of k, the algorithm is faster but less sensitivity
  - Usually, k= 4-6 for DNA sequence and
  - k= 1-2 for protein sequence.

```
Query: CAACTTGCCT


Seq1:  TCGGTTGCGTAGGTCCG

Seq2:  TTGCGTAGGTACAACATGCCTCGT

Seq3:  TCGAAGTAGCCGTCGTC

Seq4:  TAGTAACTGGCCTAGTCGTAGTC
```

# Step 1 of FastP: Identification of hotspots

- Hotspots can be found in O(n+m) time using table lookup as follows.

1. Build a lookup table that contains all k-tuples.
2. For each k-tuple in query, flag it in the lookup table
3. For each k-tuple in database, if it is flagged in the lookup table, a hotspot is identified.

# Step 2 of FastP: Find the 10 best diagonal runs

- If two nearby hotspots have the same distance in both the query sequence and the database sequence, they can be merged to form a diagonal run.

- Each hot spot is assigned a positive score. Interspot space is given a negative score that decrease with length.

- The score of a diagonal run is the sum of scores for hot spots and interspot spaces

- This steps identifies the 10 highest scoring diagonal runs for each database sequence



Query: CAACTTGCCT

Seq1: TCGGTTGCGTAGGTCCG

Seq2: TTGCGTAGGTACAACATGCCTCGT

Seq3: TCGAAGTAGCCGTCGTC

Seq4: TACTAACTGGCCTAGTCGTAGTC

# Step 3 of FastP: Rescore the 10 best diagonal runs

- Step 3: Rescore the 10 best diagonal runs for each database sequence
  - Using the substitution matrix for amino acid (or nucleotide), the diagonal runs are rescored.
  - Sub-region of each diagonal run whose similarity score is maximum is determined.
  - For each database sequence, the score of the best of the 10 sub-regions is called the init1 score.
- In other words, the init1 score is expected to be the best ungapped local alignment score between the database sequence and the query.

- After these 3 steps, we rank and report the database sequences according to their init1 scores.

# FastA (I)

- FastA using the same first 3 steps of FastP.
- Then, it executes two additional steps, namely, Steps 4 and 5.

# Step 4 of FastA: Attempts to join the sub-regions by allowing indels

- For the 10 sub-regions in Step 3, discard those with scores smaller than a given threshold

- For the remaining sub-regions, attempts to join them by allowing indels

- The score of the combined region is the sum of the scores of the sub-regions minus the penalty for gaps

- The best score can be computed using dynamic programming and it is called initn score.

# Step 5 of FastA: Banded Smith-Waterman DP

- Sequences with initn smaller than a threshold are discarded

- For the remaining sequences, apply banded Smith-Waterman dynamic programming to complete the score opt.

- Rank the sequences based on their opt scores.

# Conclusion for FlastA

- FastA is much faster than Smith-Waterman.

- It is less sensitive than
  Smith-Waterman Algorithm.

# What is BLAST?

- BLAST = Basic Local Alignment Search Tool
- Input:
  - A database D of sequences
  - A sequence s
- Aim of BLAST:
  - Compare s against all sequences in D as fast as possible using heuristics.
- Disadvantage of BLAST:
  - To be fast, it scarifies the accuracy. Thus, less sensitive.

## Standard Nucleotide BLAST

blastn    blastp    blastx    tblastn    tblastx

### Enter Query Sequence

BLASTN programs search nucleotide databases using a nucleotide query. more...

Reset page    Bookmark

Enter accession number(s), gi(s), or FASTA sequence(s) ⓘ    Clear    Query subrange ⓘ

```
GCGGCGCAAGATGGCCCAGGAGAACCCCAAGATGCACAACTCGGAGATCAGCAAGCGCCTGGGCGCCGAG
TGGAAACTTTTGTCGGAGACGGAGAAGCGGCCGTTCATCGACGAGGCTAA
```

From

To

**BLAST results will be displayed in a new format by default**
You can always switch back to the Traditional Results page.

Or, upload file    Choose File   No file chosen   ⓘ

Job Title    Nucleotide Sequence

Enter a descriptive title for your BLAST search ⓘ

☐ Align two or more sequences ⓘ

### Choose Search Set

Database    ● Standard databases (nr etc.): ○ rRNA/ITS databases ○ Genomic + transcript databases

Nucleotide collection (nr/nt)    ▼ ⓘ

Organism
Optional    Enter organism name or id--completions will be suggested    ☐ exclude  [+]
Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown ⓘ

Exclude
Optional    ☐ Models (XM/XP) ☐ Uncultured/environmental sample sequences

Limit to
Optional    ☐ Sequences from type material

Entrez Query
Optional    You Tube Create custom database
Enter an Entrez query to limit search ⓘ

### Program Selection

Optimize for    ● Highly similar sequences (megablast)
○ More dissimilar sequences (discontiguous megablast)
○ Somewhat similar sequences (blastn)
Choose a BLAST algorithm ⓘ

**BLAST**    Search database Nucleotide collection (nr/nt) using Megablast (Optimize for highly similar sequences)
☐ Show results in a new window

[+] Algorithm parameters

BLAST ® » blastn suite » results for RID-ZH1KAPW6016

Home   Recent Results   Saved Strategies   Help

‹ Edit Search      Save Search      Search Summary ⌄

? How to read this report?   ▶ BLAST Help Videos   ⟲ Back to Traditional Results Page

| Job Title | Nucleotide Sequence |
|---|---|
| RID | ZH1KAPW6016   Search expires on 12-18 14:40 pm   Download All ⌄ |
| Program | BLASTN ?   Citation ⌄ |
| Database | nt   See details ⌄ |
| Query ID | lcl|Query_63617 |
| Description | None |
| Molecule type | dna |
| Query Length | 120 |
| Other reports | Distance tree of results   MSA viewer ? |

**Filter Results**

Organism   *only top 20 will appear*                    ☐ exclude

⬇ Download ⌄        GenBank  Graphics

**Homo sapiens SRY-box transcription factor 2 (SOX2), mRNA**

Sequence ID: NM_003106.4   Length: 2512   Number of Matches: 1

Range 1: 601 to 720 GenBank   Graphics                ▼ Next Match  ▲ Previous Match

| Score | Expect | Identities | Gaps | Strand |
|---|---|---|---|---|
| 222 bits(120) | 2e-54 | 120/120(100%) | 0/120(0%) | Plus/Plus |

```
Query  1     GCGGCGCAAGATGGCCCAGGAGAACCCCAAGATGCACAACTCGGAGATCAGCAAGCGCCT  60
             ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  601   GCGGCGCAAGATGGCCCAGGAGAACCCCAAGATGCACAACTCGGAGATCAGCAAGCGCCT  660

Query  61    GGGCGCCGAGTGGAAACTTTTGTCGGAGACGGAGAAGCGGCCGTTCATCGACGAGGCTAA  120
             ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  661   GGGCGCCGAGTGGAAACTTTTGTCGGAGACGGAGAAGCGGCCGTTCATCGACGAGGCTAA  720
```

| Descriptions | Graphic Summary | Alignments | Taxon |
|---|---|---|---|

**Sequences producing significant alignments**

☑ select all   *100 sequences selected*

|  | Description |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|
| ☑ | PREDICTED: Gorilla gorilla gorilla SRY-box transcription factor 2 (SOX2), mRNA | 222 | 222 | 100% | 2e-54 | 100.00% | XM_031009512.1 |
| ☑ | PREDICTED: Nomascus leucogenys SRY-box transcription factor 2 (SOX2), mRNA | 222 | 222 | 100% | 2e-54 | 100.00% | XM_030822867.1 |
| ☑ | Homo sapiens SRY-box transcription factor 2 (SOX2), mRNA | 222 | 222 | 100% | 2e-54 | 100.00% | NM_003106.4 |
| ☑ | PREDICTED: Pan troglodytes SRY-box 2 (SOX2), mRNA | 222 | 222 | 100% | 2e-54 | 100.00% | XM_516895.7 |
| ☑ | Human ORFeome Gateway entry vector pENTR223-SOX2, complete sequence | 222 | 222 | 100% | 2e-54 | 100.00% | LT738208.1 |
| ☑ | Homo sapiens isolate MA5 transcription factor SOX-2 (SOX2) gene, complete cds | 222 | 222 | 100% | 2e-54 | 100.00% | KU342033.1 |

Nucleotide

| Nucleotide ▾ | | **Search** |

Advanced

Help

GenBank ▾

Send to: ▾

**Change region shown** ▼

# Homo sapiens SRY-box transcription factor 2 (SOX2), mRNA

**Customize view** ▼

NCBI Reference Sequence: NM_003106.4

FASTA    Graphics

```
LOCUS       NM_003106              2512 bp    mRNA    linear   PRI 08-OCT-2019
DEFINITION  Homo sapiens SRY-box transcription factor 2 (SOX2), mRNA.
ACCESSION   NM_003106
VERSION     NM_003106.4
KEYWORDS    RefSeq; MANE Select.
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
            Catarrhini; Hominidae; Homo.
REFERENCE   1  (bases 1 to 2512)
  AUTHORS   Omori H, Sato K, Nakano T, Wakasaki T, Toh S, Taguchi K, Nakagawa T
            and Masuda M.
  TITLE     Stress-triggered YAP1/SOX2 activation transcriptionally reprograms
            head and neck squamous cell carcinoma for the acquisition of
            stemness
  JOURNAL   J. Cancer Res. Clin. Oncol. 145 (10), 2433-2444 (2019)
   PUBMED   31485767
   REMARK   GeneRIF: The stress-triggered activation of YAP1/SOX2
            transcriptionally reprograms HNSCC for the acquisition of stemness.
REFERENCE   2  (bases 1 to 2512)
```

# History of BLAST

- 1990: Birth of BLAST1
  - It is very fast and dedicate to the search of local similarities **without** gaps
  - Altschul et al, Basic local alignment search tool. J. Mol. Biol., 215(3):403-410, 1990.
  - The most highly cited paper in 1990 and the third most highly cited paper in 1983-2002.
- 1996-1997: Birth of BLAST2
  - BLAST2 allows insertion of gaps
  - BLAST2 have two versions. Developed by two groups of authors independently
    - 1997: NCBI-BLAST2 (National Center for Biotechnology Information)
    - 1996: WU-BLAST2 (Washington University)

# BLAST1

- A heuristic method which searches for local similarity without gap

- Basic idea:
  - Find hits between query and database sequences
  - Extend the hits to get the ungapped local alignments

- BLAST1 can be divided into three steps:
  - Step 1: Query preprocessing
  - Step 2: Generation of hits by database scanning
  - Step 3: Extension of hits

# Step 1: Query preprocessing (DNA)

- Step 1 finds w-tuples that are similar to some w-tuples in the query Q.

- These w-tuples are called neighbor

- For DNA, w=11 by default. A w-tuple X is a neighbor of a w-tuple Y if X=Y.

- Query preprocessing builds a table of neighbors of w-tuples in the query sequence.

Q=GTCATCATG

| w-tuple | positions |
|---------|-----------|
| ATCA | 4 |
| CATC | 3 |
| CATG | 6 |
| GTCA | 1 |
| TCAT | 2, 5 |

# Step 1: Query preprocessing (protein)

- For protein, w-tuple X is a neightbor of another w-tuple Y if $\sum_{i=1}^{w} \delta(X[i], Y[i]) \geq T$. By default, w=3 and T=13.
- For example, the (3,13)-neighbor of NII is { NII, NIV, NVI } while the (3,13)-neighbor of IIN is { IIN, IVN, VIN }.

Q=NIIN

| w-tuple | positions |
|---------|-----------|
| IIN     | 2         |
| IVN     | 2         |
| NII     | 1         |
| NIV     | 1         |
| NVI     | 1         |
| VIN     | 2         |

# Step 2: Generation of hits by database scanning

- For every database sequence S in the database D,
  - Scan every position q in the sequence S, if there is an exact match between the w-tuple at position q in S and a w-tuple in the query, a hit is found.
- A hit is characterized by the positions in both query and DB sequences.

Q=GTCATCATG

>seq1
CCGCTCATGATGATCA

The list of hits:
- (5 of seq1, 2 of Q)
- (5 of seq1, 5 of Q)
- (13 of seq1, 4 of Q)

| w-tuple | positions |
|---------|-----------|
| ATCA    | 4         |
| CATC    | 3         |
| CATG    | 6         |
| GTCA    | 1         |
| TCAT    | 2, 5      |

# Step 3: Extension of hits (I)

- For every hit, extend it in both directions, without gap.
- Equivalently, in the dotplot, the hit is extended along the diagonal.



seq1

q of DB

p of query

Q

```
Q=      GTCATCATG
Seq1=CCGCTCATGATGATCA
```

# Step 3: Extension of hits (II)

- The extension is extended greedy. The extension is stopped as soon as the score decreases by more than X(parameter of the program) from the highest value reached so far.

- Example: Assume the match score is +5, the mismatch score is -4, and the extension truncation threshold is X=13. Suppose the hit between S and T is GTCG.

```
S = GTCTCACCGTCGCACGACATCTC
T = GCGGAACTGTCGAACAATCCTCT
```

From 6 to -10, drop more than X. Truncate.

-10 -6  -2 +2 +6 +1 -4
-4  -4  -4 -4 +5 +5 -4

-4  +1 +6 +2 +7 +3 -1 -5  -9
-4  +5 +5 -4 +5 -4  -4 -4  -4

From 7 to -9, drop more than X. Truncate.

# Step 3: Extension of hits (III)

From 6 to -10, drop more than X. Truncate.

From 7 to -9, drop more than X. Truncate.

```
     -10 -6  -2 +2 +6 +1 -4          -4  +1 +6 +2 +7 +3 -1 -5  -9
       -4 -4  -4 -4 +5 +5 -4         -4 +5 +5 -4 +5 -4  -4 -4  -4
S  =  GTCTCACCGTCGCACGACATCTC
T  =  GCGGAACTGTCGAACAATCCTCT
```

- To obtain the highest ungapped alignment score , we trim the first 4 aligned positions and the last 4 aligned positions. Then, the resultant ungapped alignment is

```
ACCGTCGCACGA
ACTGTCGAACAA
```
+5 +5 -4 +5 +5 +5 +5 -4  +5 +5 -4 +5

- This is known as the high scoring segment pair (HSP), its score is 28.

# Step 3: Extension of hits (IV)

- If the HSP score is better than or equal to $\alpha$ (parameter of the program), it will be reported.

- For every database sequence, it may have multiple HSPs. The best scoring HSP is called the maximal segment pair (MSP).

# NCBI-Blast2

- Report gapped local alignments.

- The first 2 steps are the same as BLAST1.

- Two major differences:
  - Two-hits requirement (implemented for protein)
  - Gapped extension

# Step 3: Two-hits requirement

- To extend a hit, we require that there is another hit on the same diagonal within a distance smaller than A

- By default, A=40

- Note: Two-hits requirement is implemented for protein sequences (not DNA).

# Comparing One hit (3,13)-neighbor and Two hit (3, 11)-neighbor



+'s are 15 (3, 13) hits.
.'s are 22 (3, 11) hits.
The two lines are pairs of (3,11) hits.

The two-hit method (A=40) is more sensitive for HSPs with score > 33 bits.

# Step 4: Gapped extension (I)

- Obtain a set of HSPs by performing ungapped extension similar to Step 3 of BLAST1

- Among the generated HSP, we perform gapped extension for those with score > some threshold

- Note: The ungapped extension is a filter. It reduces the number of gapped extensions.

# Step 4: Gapped extension (II)

- For each HSP, the length-11 segment with the highest alignment score is identified. (If the HSP is shorter than 11, the whole segment is taken.)
- Gapped extension is started from the middle of the segment.

# Illustration

S = CCCTAACAGTTGACAATCCC
T = TGGTACGAGTTGAGCATGGT

S = TAAC-AGTTGA-CAAT
T = T-ACGAGTTGAGCA-T

Assume match scores +2
Mismatch/indel score -1.

|   |    | _ | G | T | A | G | C | A | T | G | G | T |
|---|----|----|----|----|----|----|----|----|----|----|----|----|
|   |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 | -1 | -2 | -3 | -4 | -5 | -6 | -7 | -8 | -9 | -10 |
| G | 1  | -1 | 2 | 1 | 0 | -1 | -2 | -3 | -4 | -5 | -6 | -7 |
| T | 2  | -2 | 1 | 4 | 3 | 2 | 1 | 0 | -1 | -2 | -3 | -4 |
| A | 3  | -3 | 0 | 3 | 6 | 5 | 4 | 3 | 2 | 1 | 0 | -1 |
| C | 4  | -4 | -1 | 2 | 5 | 5 | 7 | 6 | 5 | 4 | 3 | -1 |
| A | 5  | -5 | -2 | 1 | 4 | 4 | 6 | 9 | 8 | 7 | 6 | 5 |
| A | 6  | -6 | -3 | 0 | 3 | 3 | 6 | 8 | 7 | 6 | 5 | 5 |
| T | 7  | -7 | -4 | -1 | 2 | 2 | 4 | 7 | 10 | 9 | 8 | 7 |
| C | 8  | -8 | -5 | -2 | 1 | 1 | 4 | 6 | 9 | 9 | 8 | 7 |
| C | 9  | -9 | -6 | -3 | 0 | 0 | 3 | 5 | 8 | 8 | 8 | 7 |
| C | 10 | -10 | -7 | -4 | -1 | -1 | 2 | 4 | 7 | 7 | 7 | 7 |

# Step 4: Gapped extension (II)

- Running DP starting from the middle of the segment is slow.

- To reduce time, banded dynamic programming can be used.

- To further speedup, BLAST proposed to use X-drop algorithm for gapped extension

- X-drop is a modified semi-global alignment algorithm
  - It fill in only entries that are no more than X below the best alignment score yet found.



Banded DP                    X-drop algorithm

# X-drop algorithm



- Let antidiagonal-d be the entries V[i,j] in the table where i+j=d.
- The X-drop algorithm fills in entries in antidiagonal-d in increasing order of d and maintain the best alignment score seen so far in the variable $\tau$.

- For each antidiagonal, X-drop algorithm keeps entries that are consistent.
- An entry V[i,j] is consistent if
  - Either V[i-1,j-1], V[i-1,j] or V[i,j-1] is consistent; and
  - V[i,j] is no more than X below the best alignment score seen so far (i.e. V[i,d-i] $\geq \tau$-X).
- The X-drop algorithm will stop when the algorithm cannot generate consistent entries.

# Consistent entries

- X-drop algorithm computes V[i,j] on antidiagonal-d if:
  - V[i-1,j-1] on antidiagonal-(d-2) is consistent; or
  - V[i, j-1] or V[i-1,j] on antidiagonal-(d-1) is consistent.

- $V[i,j] = \max \begin{cases} V[i-1,j-1] \; + \delta(S[i], T[d-i]) \\ \quad\quad V[i-1,j] + \delta(S[i], \_) \\ \quad V[i,j-1] + \delta(\_, T[d-i]) \end{cases}$

- After V[i, j] is computed, V[i,j] is consistent if V[i,j] $\geq \tau$-X.

# X-drop algorithm

- 1. Set V[0,0]=0 as a consistent entry, $\tau$=0.
- 2. For d = 1, 2, …
  - 2.1. If both antidiagonal-(d-1) and antidiagonal-(d-2) do not have consistent entries, then stop;
  - 2.2. For every V[i,j] on antidiagonal-d such that V[i-1,j-1], V[i-1,j] or V[i,j-1] is a consistent entry,
    - Compute V[i,j]
    - If V[i,j] $\geq \tau$-X, set V[i,j] as a consistent entry
  - 2.3 Set $\tau$ = max {$\tau$, max{V[i,j] | V[i,j] is a consistent entry on antidiagonal-d } }

# Illustration (X-drop approach)

Perform semi-global alignment starting from the middle of the hit.
Since this example is symmetric, we only show the DP on the right segment.

S = CCCTAACAG̲T̲T̲GACAATCCC
T = TGGTACGAG̲T̲T̲GAGCATGGT

$\tau$=0
d=0
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|   |    | _ | G | T | A | G | C | A | T | G | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
|   |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 |   |   |   |   |   |   |   |   |   |   |
| G | 1  |   |   |   |   |   |   |   |   |   |   |   |
| T | 2  |   |   |   |   |   |   |   |   |   |   |   |
| A | 3  |   |   |   |   |   |   |   |   |   |   |   |
| C | 4  |   |   |   |   |   |   |   |   |   |   |   |
| A | 5  |   |   |   |   |   |   |   |   |   |   |   |
| A | 6  |   |   |   |   |   |   |   |   |   |   |   |
| T | 7  |   |   |   |   |   |   |   |   |   |   |   |
| C | 8  |   |   |   |   |   |   |   |   |   |   |   |
| C | 9  |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Illustration (X-drop approach)

Perform semi-global alignment starting from the middle of the hit.
Since this example is symmetric, we only show the DP on the right segment.

S = CCCTAACAG<u>TT</u>GACAATCCC
T = TGGTACGAG<u>TT</u>GAGCATGGT

$\tau=0$
d=1
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|   |    | _ | G | T | A | G | C | A | T | G | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|----|
|   |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 | -1 |   |   |   |   |   |   |   |   |    |
| G | 1  | -1 |   |   |   |   |   |   |   |   |   |    |
| T | 2  |   |   |   |   |   |   |   |   |   |   |    |
| A | 3  |   |   |   |   |   |   |   |   |   |   |    |
| C | 4  |   |   |   |   |   |   |   |   |   |   |    |
| A | 5  |   |   |   |   |   |   |   |   |   |   |    |
| A | 6  |   |   |   |   |   |   |   |   |   |   |    |
| T | 7  |   |   |   |   |   |   |   |   |   |   |    |
| C | 8  |   |   |   |   |   |   |   |   |   |   |    |
| C | 9  |   |   |   |   |   |   |   |   |   |   |    |
| C | 10 |   |   |   |   |   |   |   |   |   |   |    |

# Illustration (X-drop approach)

S = CCCTAACAG<u>TT</u>GACAATCCC
T = TGGTACGAG<u>TT</u>GAGCATGGT

τ=2
d=2
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|   |    | _ | G | T | A | G | C | A | T | G | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
|   |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 | -1 | -2 |   |   |   |   |   |   |   |   |
| G | 1  | -1 | 2 |   |   |   |   |   |   |   |   |   |
| T | 2  | -2 |   |   |   |   |   |   |   |   |   |   |
| A | 3  |   |   |   |   |   |   |   |   |   |   |   |
| C | 4  |   |   |   |   |   |   |   |   |   |   |   |
| A | 5  |   |   |   |   |   |   |   |   |   |   |   |
| A | 6  |   |   |   |   |   |   |   |   |   |   |   |
| T | 7  |   |   |   |   |   |   |   |   |   |   |   |
| C | 8  |   |   |   |   |   |   |   |   |   |   |   |
| C | 9  |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Illustration (X-drop approach)

Perform semi-global alignment starting from the middle of the hit.
Since this example is symmetric, we only show the DP on the right segment.

S = CCCTAACAG<u>TT</u>GACAATCCC
T = TGGTACGAG<u>TT</u>GAGCATGGT

$\tau=2$
d=3
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

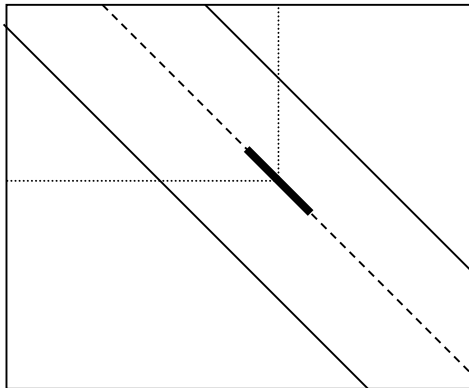|   |    | _ | G | T | A | G | C | A | T | G | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
|   |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 | -1 | -2 |   |   |   |   |   |   |   |   |
| G | 1  | -1 | 2 | 1 |   |   |   |   |   |   |   |   |
| T | 2  | -2 | 1 |   |   |   |   |   |   |   |   |   |
| A | 3  |   |   |   |   |   |   |   |   |   |   |   |
| C | 4  |   |   |   |   |   |   |   |   |   |   |   |
| A | 5  |   |   |   |   |   |   |   |   |   |   |   |
| A | 6  |   |   |   |   |   |   |   |   |   |   |   |
| T | 7  |   |   |   |   |   |   |   |   |   |   |   |
| C | 8  |   |   |   |   |   |   |   |   |   |   |   |
| C | 9  |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Illustration (X-drop approach)

Perform semi-global alignment starting from the middle of the hit.
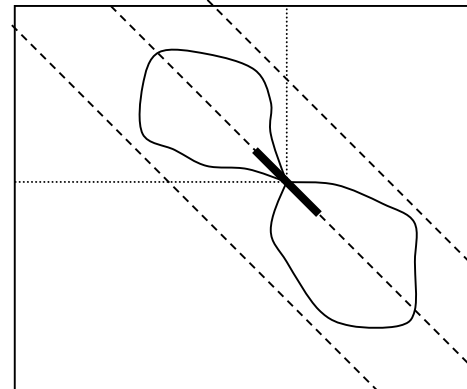Since this example is symmetric, we only show the DP on the right segment.

S = CCCTAACA<u>GTT</u>GACAATCCC
T = TGGTACGA<u>GTT</u>GAGCATGGT

$\tau=4$
d=4
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|   |    | _  | G  | T  | A  | G | C | A | T | G | G | T  |
|---|----|----|----|----|----|---|---|---|---|---|---|----|
|   |    | 0  | 1  | 2  | 3  | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0  | -1 | -2 |    |   |   |   |   |   |   |    |
| G | 1  | -1 | 2  | 1  | 0  |   |   |   |   |   |   |    |
| T | 2  | -2 | 1  | 4  |    |   |   |   |   |   |   |    |
| A | 3  |    | 0  |    |    |   |   |   |   |   |   |    |
| C | 4  |    |    |    |    |   |   |   |   |   |   |    |
| A | 5  |    |    |    |    |   |   |   |   |   |   |    |
| A | 6  |    |    |    |    |   |   |   |   |   |   |    |
| T | 7  |    |    |    |    |   |   |   |   |   |   |    |
| C | 8  |    |    |    |    |   |   |   |   |   |   |    |
| C | 9  |    |    |    |    |   |   |   |   |   |   |    |
| C | 10 |    |    |    |    |   |   |   |   |   |   |    |

# Illustration (X-drop approach)

S = CCCTAACAG‾T‾T‾GACAATCCC
T = TGGTACGAG‾T‾T‾GAGCATGGT

$\tau=9$
d=12
Black: Consistent
Green: Not consistent

|   |    | _  | G  | T  | A | G | C | A | T | G | G | T  |
|---|----|----|----|----|---|---|---|---|---|---|---|----|
|   |    | 0  | 1  | 2  | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0  | -1 | -2 |   |   |   |   |   |   |   |    |
| G | 1  | -1 | 2  | 1  | 0 |   |   |   |   |   |   |    |
| T | 2  | -2 | 1  | 4  | 3 | 2 |   |   |   |   |   |    |
| A | 3  |    | 0  | 3  | 6 | 5 | 4 |   |   |   |   |    |
| C | 4  |    |    | 2  | 5 | 5 | 7 | 6 | 5 |   |   |    |
| A | 5  |    |    |    | 4 | 4 | 6 | 9 | 8 |   |   |    |
| A | 6  |    |    |    |   |   | 5 | 8 |   |   |   |    |
| T | 7  |    |    |    |   |   |   |   |   |   |   |    |
| C | 8  |    |    |    |   |   |   |   |   |   |   |    |
| C | 9  |    |    |    |   |   |   |   |   |   |   |    |
| C | 10 |    |    |    |   |   |   |   |   |   |   |    |

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

# Illustration (X-drop approach)

S = CCCTAACA<u>GTT</u>GACAATCCC
T = TGGTACGA<u>GTT</u>GAGCATGGT

$\tau=9$
d=13
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|   |    | _ | G | T | A | G | C | A | T | G | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
|   |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 | -1 | -2 |   |   |   |   |   |   |   |   |
| G | 1  | -1 | 2 | 1 | 0 |   |   |   |   |   |   |   |
| T | 2  | -2 | 1 | 4 | 3 | 2 |   |   |   |   |   |   |
| A | 3  |   | 0 | 3 | 6 | 5 | 4 |   |   |   |   |   |
| C | 4  |   |   | 2 | 5 | 5 | 7 | 6 | 5 |   |   |   |
| A | 5  |   |   |   | 4 | 4 | 6 | 9 | 8 | 7 |   |   |
| A | 6  |   |   |   |   |   | 5 | 8 | 7 |   |   |   |
| T | 7  |   |   |   |   |   |   | 7 |   |   |   |   |
| C | 8  |   |   |   |   |   |   |   |   |   |   |   |
| C | 9  |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Illustration (X-drop approach)

Perform semi-global alignment starting from the middle of the hit.
Since this example is symmetric, we only show the DP on the right segment.

S = CCCTAACA<u>GTT</u>GACAATCCC
T = TGGTACGA<u>GTT</u>GAGCATGGT

$\tau=10$
d=14
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|   |    | _ | G | T | A | G | C | A | T | G | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
|   |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 | -1 | -2 |   |   |   |   |   |   |   |   |
| G | 1  | -1 | 2 | 1 | 0 |   |   |   |   |   |   |   |
| T | 2  | -2 | 1 | 4 | 3 | 2 |   |   |   |   |   |   |
| A | 3  |   | 0 | 3 | 6 | 5 | 4 |   |   |   |   |   |
| C | 4  |   |   | 2 | 5 | 5 | 7 | 6 | 5 |   |   |   |
| A | 5  |   |   |   | 4 | 4 | 6 | 9 | 8 | 7 |   |   |
| A | 6  |   |   |   |   |   | 5 | 8 | 7 | 7 |   |   |
| T | 7  |   |   |   |   |   |   |   | 7 | 10 |   |   |
| C | 8  |   |   |   |   |   |   |   |   |   |   |   |
| C | 9  |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Illustration (X-drop approach)

Perform semi-global alignment starting from the middle of the hit.
Since this example is symmetric, we only show the DP on the right segment.

S = CCCTAACA<u>GTT</u>GACAATCCC
T = TGGTACGA<u>GTT</u>GAGCATGGT

$\tau=10$
d=15
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|   |    | _ | G | T | A | G | C | A | T | G | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
|   |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 | -1 | -2 |   |   |   |   |   |   |   |   |
| G | 1  | -1 | 2 | 1 | 0 |   |   |   |   |   |   |   |
| T | 2  | -2 | 1 | 4 | 3 | 2 |   |   |   |   |   |   |
| A | 3  |   | 0 | 3 | 6 | 5 | 4 |   |   |   |   |   |
| C | 4  |   |   | 2 | 5 | 5 | 7 | 6 | 5 |   |   |   |
| A | 5  |   |   |   | 4 | 4 | 6 | 9 | 8 | 7 |   |   |
| A | 6  |   |   |   |   |   | 5 | 8 | 7 | 7 |   |   |
| T | 7  |   |   |   |   |   |   | 7 | 10 | 9 |   |   |
| C | 8  |   |   |   |   |   |   |   | 9 |   |   |   |
| C | 9  |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Illustration (X-drop approach)

S = CCCTAACAG̲T̲T̲GACAATCCC
T = TGGTACGAG̲T̲T̲GAGCATGGT

$\tau=10$
d=16
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|   |    | _ 0 | G 1 | T 2 | A 3 | G 4 | C 5 | A 6 | T 7 | G 8 | G 9 | T 10 |
|---|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|
| _ | 0  | 0   | -1  | -2  |     |     |     |     |     |     |     |      |
| G | 1  | -1  | 2   | 1   | 0   |     |     |     |     |     |     |      |
| T | 2  | -2  | 1   | 4   | 3   | 2   |     |     |     |     |     |      |
| A | 3  |     | 0   | 3   | 6   | 5   | 4   |     |     |     |     |      |
| C | 4  |     |     | 2   | 5   | 5   | 7   | 6   | 5   |     |     |      |
| A | 5  |     |     |     | 4   | 4   | 6   | 9   | 8   | 7   |     |      |
| A | 6  |     |     |     |     |     | 5   | 8   | 7   | 7   |     |      |
| T | 7  |     |     |     |     |     |     | 7   | 10  | 9   | 8   |      |
| C | 8  |     |     |     |     |     |     |     | 9   | 9   |     |      |
| C | 9  |     |     |     |     |     |     |     |     | 8   |     |      |
| C | 10 |     |     |     |     |     |     |     |     |     |     |      |

# Illustration (X-drop approach)

S = CCCTAACAG<u>TT</u>GACAATCCC
T = TGGTACGAG<u>TT</u>GAGCATGGT

$\tau=10$
d=17
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

| | | _ | G | T | A | G | C | A | T | G | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0 | 0 | -1 | -2 | | | | | | | | |
| G | 1 | -1 | 2 | 1 | 0 | | | | | | | |
| T | 2 | -2 | 1 | 4 | 3 | 2 | | | | | | |
| A | 3 | | 0 | 3 | 6 | 5 | 4 | | | | | |
| C | 4 | | | 2 | 5 | 5 | 7 | 6 | 5 | | | |
| A | 5 | | | | 4 | 4 | 6 | 9 | 8 | 7 | | |
| A | 6 | | | | | | 5 | 8 | 7 | 7 | | |
| T | 7 | | | | | | | 7 | 10 | 9 | 8 | |
| C | 8 | | | | | | | | 9 | 9 | 8 | |
| C | 9 | | | | | | | | 8 | 8 | | |
| C | 10 | | | | | | | | | | | |

# Illustration (X-drop approach)

Perform semi-global alignment starting from the middle of the hit.
Since this example is symmetric, we only show the DP on the right segment.

S = CCCTAACA<u>GTT</u>GACAATCCC
T = TGGTACGA<u>GTT</u>GAGCATGGT

$\tau$=10
d=18
Black: Consistent
Green: Not consistent

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

|     |     | _ | G | T | A | G | C | A | T | G | G | T |
|-----|-----|---|---|---|---|---|---|---|---|---|---|---|
|     |     | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _   | 0   | 0 | -1 | -2 |   |   |   |   |   |   |   |   |
| G   | 1   | -1 | 2 | 1 | 0 |   |   |   |   |   |   |   |
| T   | 2   | -2 | 1 | 4 | 3 | 2 |   |   |   |   |   |   |
| A   | 3   |   | 0 | 3 | 6 | 5 | 4 |   |   |   |   |   |
| C   | 4   |   |   | 2 | 5 | 5 | 7 | 6 | 5 |   |   |   |
| A   | 5   |   |   |   | 4 | 4 | 6 | 9 | 8 | 7 |   |   |
| A   | 6   |   |   |   |   |   | 5 | 8 | 7 | 7 |   |   |
| T   | 7   |   |   |   |   |   |   | 7 | 10 | 9 | 8 |   |
| C   | 8   |   |   |   |   |   |   |   | 9 | 9 | 8 |   |
| C   | 9   |   |   |   |   |   |   |   | 8 | 8 | 8 |   |
| C   | 10  |   |   |   |   |   |   |   |   |   |   |   |

# Illustration (X-drop approach)

S = CCCTAACA<u>GTT</u>GACAATCCC
T = TGGTACGA<u>GTT</u>GAGCATGGT

S = TAAC-A<u>GTT</u>GA-CAAT
T = T-ACGA<u>GTT</u>GAGCA-T

Assume match scores +2
Mismatch/indel score -1.
Assume X=1

Perform semi-global alignment starting from the middle of the hit.
Since this example is symmetric, we only show the DP on the right segment.

|   |    | _  | G  | T  | A  | G  | C  | A  | T  | G  | G  | T  |
|---|----|----|----|----|----|----|----|----|----|----|----|----|
|   |    | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 |
| _ | 0  | 0  | -1 | -2 |    |    |    |    |    |    |    |    |
| G | 1  | -1 | 2  | 1  | 0  |    |    |    |    |    |    |    |
| T | 2  | -2 | 1  | 4  | 3  | 2  |    |    |    |    |    |    |
| A | 3  |    | 0  | 3  | 6  | 5  | 4  |    |    |    |    |    |
| C | 4  |    |    | 2  | 5  | 5  | 7  | 6  | 5  |    |    |    |
| A | 5  |    |    |    | 4  | 4  | 6  | 9  | 8  | 7  |    |    |
| A | 6  |    |    |    |    |    | 5  | 8  | 7  | 7  |    |    |
| T | 7  |    |    |    |    |    |    | 7  | 10 | 9  | 8  |    |
| C | 8  |    |    |    |    |    |    |    | 9  | 9  | 8  |    |
| C | 9  |    |    |    |    |    |    | 8  | 8  | 8  |    |    |
| C | 10 |    |    |    |    |    |    |    |    |    |    |    |

# X-drop algorithm allowing affine gap penalty

- The previous algorithm is for linear gap.

- We can generalize the X-drop algorithm to handle affine gap penalty.

- Exercise!

# BLAST2 algorithm

1. Query preprocessing

2. Generation of hits by database scanning

3. Filter hits by 2-hit requirement (for protein only)

4. Generate HSP by ungapped extension. For each HSP with score > some threshold, perform gapped extension using X-drop algorithm.

# BLAST1 vs. NCBI-BLAST2

- BLAST1 spends 90% of its time on extension
- For NCBI-BLAST2, due to the two-hits requirement, the number of extensions is reduced.
  - NCBI-BLAST2 is about 3 times faster than BLAST1.

# BLAST program options

| Program | Query | Database | Alignment type |
| --- | --- | --- | --- |
| blastn | DNA | DNA | Search DNA query sequence in DNA database |
| blastp | Protein | Protein | Search protein query sequence in protein database |
| blastx | DNA | Protein | Convert DNA query sequence into protein sequences in all 6 reading frames. Search these translated proteins in protein database |
| tblastn | Protein | DNA | Search protein query sequence again protein sequences generated from the 6 reading frames of the DNA sequences in the DNA database |
| tblastx | DNA | DNA | Convert DNA query sequence into protein sequences in all 6 reading frames. Search these translated protein query sequence again protein sequences generated from the 6 reading frames of the DNA sequences in the DNA database |

# Statistics for local alignment

- An ungapped local alignment consists of a pair of equal length segments.

- BLAST finds the local alignments whose score cannot be improved by extension. Such local alignments are called high-scoring segment pairs or HSPs.

- To determine the significant of each HSP, BLAST gives E-value and bit score. Below, we give a brief discussion on them.

# Assumption

- We required the expected score for aligning a random pair of residues/bases to be negative.

- Otherwise, the longer the alignment, the higher is the score independent of whether the segments aligned are related or not.

|   | A | C | G | T |
|---|---|---|---|---|
| A | 2 | -1 | -1 | -1 |
| C | -1 | 2 | -1 | -1 |
| G | -1 | -1 | 2 | -1 |
| T | -1 | -1 | -1 | 2 |

Good!

Expected score =
$$\frac{1}{16}\left(4 * 2 + 12 * (-1)\right) = -\frac{1}{4}$$

|   | A | C | G | T |
|---|---|---|---|---|
| A | 5 | -1 | -1 | -1 |
| C | -1 | 5 | -1 | -1 |
| G | -1 | -1 | 5 | -1 |
| T | -1 | -1 | -1 | 5 |

Not Good!

Expected score =
$$\frac{1}{16}\left(4 * 5 + 12 * (-1)\right) = \frac{1}{2}$$

# Raw Score for an ungapped alignment in BLAST

- There are 12 matches and 3 mismatches
- Raw score = 12*2 − 3*3 = 15.

GAACGTACTCCTACG
GCACTTACTGCTACG

|   | A | C | G | T |
|---|---|---|---|---|
| **A** | 2 | -3 | -3 | -3 |
| **C** | -3 | 2 | -3 | -3 |
| **G** | -3 | -3 | 2 | -3 |
| **T** | -3 | -3 | -3 | 2 |

BLAST score matrix

Expected score =
$$\frac{1}{16}\left(4 * 2 + 12 * (-3)\right) = -1.75$$

# E-value

- Let S be the raw score of a HSP between a length-m query sequence and a length-n database sequence.
- E-value of the HSP is defined as the expected number of random alignments whose score is bigger than the raw score S.
- Denote E be such E-value.
- When E is small, the HSP is unlikely to be a random alignment.
- By default, BLAST reports the HSP if $E \leq 10$.

# E-value (II)

- When both m (length of query) and n (length of database sequence) are sufficiently long,
  - E follows the extreme distribution (Gumbel distribution). We have:
    - $E = Kmne^{-\lambda S}$

      for some parameters K and $\lambda$ which depends on the similarity matrix $\delta$ and the expected frequencies of the residues/bases.
    - For more information on estimating K and $\lambda$, please read
      - http://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-3.html
      - http://oreilly.com/catalog/blast/chapter/ch04.pdf

# E-value (III)

- We can justify the equation $E = Kmne^{-\lambda S}$ intuitively:
    - Double the length of either sequence will double the expected number of HSPs. (i.e. $E \propto nm$)
    - Double the score S will exponentially reduce the expected number of HSPs. (i.e. $\ln E \propto S$ or $E \propto e^{-\lambda S}$)
- Hence, $E \propto mne^{-\lambda S}$.

Query sequence    m

Database sequence    n

# Example

|   | A | C | G | T |
|---|---|---|---|---|
| A | 2 | -3 | -3 | -3 |
| C | -3 | 2 | -3 | -3 |
| G | -3 | -3 | 2 | -3 |
| T | -3 | -3 | -2 | 2 |

BLAST matrix $\delta$

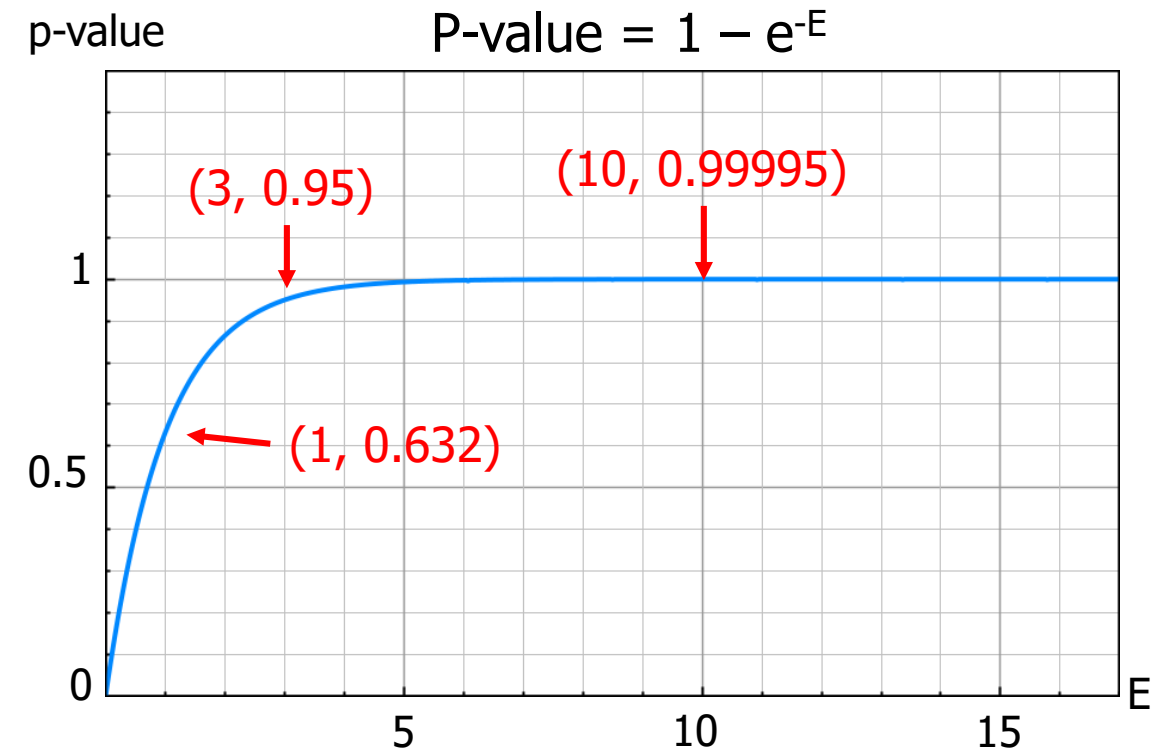- The raw score for the following HSP is S=15.

GAACGTACTCCTACG
GCACTTACTGCTACG

- For matrix $\delta$, ln K $\approx$ 0.42 and $\lambda \approx$ 0.626

- Suppose m (query length) = 20, n (database sequence length) = 70.

- The E-value is $E = Kmne^{-\lambda S} = 0.178035469$

# Bit score

- The raw score S of an alignment depends on the scoring system.
- Without knowing the scoring system, the raw score is meaningless.

- The bit score is a normalized raw score, which is independent of the scoring system. It is defined as:

$$S' = \frac{\lambda S - \ln K}{\ln 2}$$

- By definition, $2^{-S'} = Ke^{-\lambda S}$. Hence, $E = Kmne^{-\lambda S} = mn2^{-S'}$.
- E-value can be computed without parameters depending on the scoring system.

# Example

| | A | C | G | T |
|---|---|---|---|---|
| **A** | 2 | -3 | -3 | -3 |
| **C** | -3 | 2 | -3 | -3 |
| **G** | -3 | -3 | 2 | -3 |
| **T** | -3 | -3 | -2 | 2 |

Matrix $\delta$

- The raw score for the following HSP is S=15.

    GAACGTACTCCTACG
    GCACTTACTGCTACG

- For matrix $\delta$, ln K $\approx$ 0.42 and $\lambda \approx$ 0.626

- Suppose m (query length) = 20, n (database sequence length) = 70.

- The bit score is $S' = \frac{\lambda S - \ln K}{\ln 2} = 12.94097452$

- Note that $E = Kmne^{-\lambda S} = 0.178035469$ and $mn2^{-S'} = 0.178035469$

# P-value

- Let X be the number of random alignments with score $\geq$ S.
- X ~ Poisson(E) where $E = Kmne^{-\lambda S}$ is the E-score.
- Pr(exactly x random alignments with score $\geq$ S) $= \dfrac{e^{-E}E^x}{x!}$
  - where $E = Kmne^{-\lambda S}$ is the E-score

- Definition: P-value is the probability of getting at least 1 random alignment with score $\geq$ S.
- Hence, p-value = Pr(X $\geq$ 1) $= 1 - \Pr(X = 0) = 1 - e^{-E}$.

# Why there is no p-value in BLAST?

- In BLAST, p-value is not shown since we expect p-value and E-value are approximately the same when E<0.01 while p-value is almost 1 when E>10.

p-value $\quad\quad\quad$ P-value $= 1 - e^{-E}$



Note:
 When E=1, p-value is $1-e^{-1} = 0.632$.
 When E=3, p-value is $1-e^{-3} = 0.95$.
 When E=10, p-value is $1-e^{-10} = 0.99995$
 When E increases, p-value is approaching 1.
 When E>10, p-value$\approx$1.

 When E=0.1, p-value is $1 - e^{-0.1} = 0.095$.
 When E=0.01, p-value is $1 - e^{-0.01} = 0.00995$.
 When E<0.01, $1-e^{-E}\approx$E.

# Local alignment with gaps

- There is no solid theoretical foundation for local alignment with gaps.
- Moreover, experimental results suggested that the theory for ungapped local alignment can be applied to the gapped local alignment as well.

# Completeness of BLAST (I)

- BLAST is the most popular solution for finding local alignments. It is well-known that BLAST is heuristics and it will miss solution.

- We would like to check how many good alignments are missed by BLAST.

- We extracted 2000 mRNA sequences from each of the 4 different species. We aligned them on human genome. Then, we checked how many significant alignments are missed by BLAST.

# Completeness of BLAST (II)

| E-value ($\leq$) | Chimpanzee Missing % | Mouse Missing % | Chicken Missing % | Zebrafish Missing % | All 4 species Missing % |
|---|---|---|---|---|---|
| $1.0 \times 10^{-16}$ | 0.00 | 0.03 | 0.05 | 0.06 | 0.01 |
| $1.0 \times 10^{-15}$ | 0.00 | 0.03 | 0.05 | 0.06 | 0.02 |
| $1.0 \times 10^{-14}$ | 0.00 | 0.04 | 0.06 | 0.06 | 0.02 |
| $1.0 \times 10^{-13}$ | 0.00 | 0.03 | 0.07 | 0.14 | 0.02 |
| $1.0 \times 10^{-12}$ | 0.01 | 0.04 | 0.10 | 0.17 | 0.03 |
| $1.0 \times 10^{-11}$ | 0.02 | 0.05 | 0.11 | 0.28 | 0.05 |
| $1.0 \times 10^{-10}$ | 0.02 | 0.07 | 0.13 | 0.39 | 0.06 |
| $1.0 \times 10^{-9}$ | 0.03 | 0.09 | 0.16 | 0.60 | 0.08 |
| $1.0 \times 10^{-8}$ | 0.05 | 0.11 | 0.25 | 0.77 | 0.12 |
| $1.0 \times 10^{-7}$ | 0.10 | 0.19 | 0.31 | 0.81 | 0.18 |
| $1.0 \times 10^{-6}$ | 0.17 | 0.31 | 0.45 | 1.08 | 0.28 |
| $1.0 \times 10^{-5}$ | 0.32 | 0.47 | 0.70 | 1.45 | 0.45 |
| $1.0 \times 10^{-4}$ | 0.57 | 0.88 | 0.99 | 1.81 | 0.75 |
| $1.0 \times 10^{-3}$ | 0.99 | 1.36 | 1.25 | 2.25 | 1.17 |
| $1.0 \times 10^{-2}$ | 1.69 | 2.11 | 1.68 | 2.61 | 1.84 |
| $1.0 \times 10^{-1}$ | 2.70 | 2.97 | 2.33 | 2.86 | 2.76 |

- 2000 queries for each species.
- BLAST only missed 0.06% of those 8000 queries (with E-value smaller than $1.0 \times 10^{-10}$).
- In conclusion, BLAST is accurate enough in most cases, yet the few alignments missed could be critical for biological research.

# Variation of BLAST

- MegaBLAST
- BLAT
- PatternHunter
- PSI-BLAST

# MegaBLAST

- BLAST targets to find alignments with similarity ≈ 89%.

- For many database queries, the query sequences are highly similar with the database sequences.

- MegaBLAST targets to find alignments with similarity > 95%. (For example, compare homologous sequences between human and chimpanzee.)

- Up to 10 times faster than the traditional BLAST.

- Now, it is the default database search method in NCBI.

# MegaBLAST uses longer w-tuples

- MegaBLAST targets to find high similarity alignment.

- Long exact match hits are expected to occur between query and database sequences.

- Hence, MegaBLAST uses longer w-tuple to find hits.

- Unlike BLAST, which sets w=11. MegaBLAST sets w=28.

- This modification reduces the number of hits. However, it also reduces the sensitivity.

# Penalty for gaps

- Since MegaBLAST targets high homology alignment, we expect it has less long gaps.

- MegaBLAST uses linear gap penalty instead of affine gap penalty.

# Why $r_{ind} = r_{mis} - r_{mat}/2$ ?

- For efficiency purpose, MegaBLAST computes differences between two sequences instead of computing their similarity.

- Below lemma shows that the similarity score can be computed by the difference score when $r_{ind} = r_{mis} - r_{mat}/2$.


- **Lemma**: Suppose $r_{ind} = r_{mis} - r_{mat}/2$. An alignment of S[1..n] and T[1..m] with edit distance d has score $(n+m)r_{mat}/2 - d(r_{mat} - r_{mis})$.

- Proof: Exercise!

# Correctness of the lemma

- $r_{ind} = r_{mis} - r_{mat}/2 \Leftrightarrow 2\,r_{ind} + r_{mat} = 2\,r_{mis}$.

- For any two mismatches in an alignment, we can replace it by two indels and one match. Then, both the alignment score and the edit distance are the same.

```
......A......G......              ......A-......G......
......C......T......              ......-C......G......
```

Score = $2\,r_{mis}$          Score = $2\,r_{ind} + r_{mat}$
Edit dist = 2              Edit dist = 2

# MegaBLAST score matrix

| | A | C | G | T |
|---|---|---|---|---|
| A | 1 | -2 | -2 | -2 |
| C | -2 | 1 | -2 | -2 |
| G | -2 | -2 | 1 | -2 |
| T | -2 | -2 | -2 | 1 |

MegaBLAST score matrix

- MegaBLAST targets 95% similarity (i.e. $\tau$=0.95).

- Match/Mismatch ratio = $\dfrac{-\log(4*\tau)}{\log\left(4*\frac{(1-\tau)}{3}\right)} = \dfrac{-\lg(4*0.95)}{\lg(4*0.05/3)} \approx 0.5$.

- MegaBLAST sets the default match score is $r_{mat}$=1 and mismatch score is $r_{mis}$=-2.

- MegaBLAST sets $r_{ind} = r_{mis} - r_{mat}/2 = -2 - 1/2 = -2.5$.

# Example

$$V[i,j] = \max \begin{cases} V[i-1,j-1] + \delta(S[i],T[j]) \\ V[i-1,j] - 2.5 \\ V[i,j-1] - 2.5 \end{cases}$$

Match score $r_{mat} = +1$
Mismatch score $r_{mis} = -2.$
Indel score $r_{ind} = -2.5$

$\delta(x,y) = 1$ if x=y; otherwise, $\delta(x,y) = -2.$

S  =  GTACAATCCC
T  =  GTAGCAATGT

S  =  GTA-CAAT
T  =  GTAGCAAT

|   |    | _ | G | T | A | G | C | A | A | T | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
| **V** |  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0 | 0 | -2.5 | -5 | -7.5 | -10 | -12.5 | -15 | -17.5 | -20 | -22.5 | -25 |
| G | 1 | -2.5 | 1 | -1.5 | -4 | -6.5 | -9 | -11.5 | -14 | -16.5 | -19 | -21.5 |
| T | 2 | -5 | -1.5 | 2 | -0.5 | -3 | -5.5 | -8 | -10.5 | -13 | -15.5 | -18 |
| A | 3 | -7.5 | -4 | -0.5 | 3 | 0.5 | -2 | -4.5 | -7 | -9.5 | -12 | -14.5 |
| C | 4 | -10 | -6.5 | -3 | 0.5 | 1 | 1.5 | -1 | -3.5 | -6 | -8.5 | -11 |
| A | 5 | -12.5 | -9 | -5.5 | -2 | -1.5 | -1 | 2.5 | 0 | -2.5 | -5 | -7.5 |
| A | 6 | -15 | -11.5 | -8 | -4.5 | -4 | -3.5 | 0 | 3.5 | 1 | -1.5 | -4 |
| T | 7 | -17.5 | -14 | -10.5 | -7 | -6.5 | -6 | -2.5 | 1 | 4.5 | 2 | -0.5 |
| C | 8 | -20 | -16.5 | -13 | -9.5 | -9 | -8.5 | -5 | -1.5 | 2 | 2.5 | 0 |
| C | 9 | -22.5 | -19 | -15.5 | -12 | -11.5 | -11 | -7.5 | -4 | -0.5 | 0 | 0.5 |
| C | 10 | -25 | -21.5 | -18 | -14.5 | -14 | -13.5 | -10 | -6.5 | -3 | -2.5 | -2 |

# Example

$$D[i,j] = \min \begin{cases} D[i-1,j-1] + \delta(S[i], T[j]) \\ D[i-1,j] + 1 \\ D[i,j-1] + 1 \end{cases}$$

$\delta(x,y) = 0$ if x=y; otherwise, $\delta(x,y) = 1$.

Let D[i,j] be the minimum number of differences between the alignment of S[1..i] and T[1..j].

S = GTACAATCCC
T = GTAGCAATGT

S = GTA-CAAT
T = GTAGCAAT

| D | | _ | G | T | A | G | C | A | A | T | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| G | 1 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| T | 2 | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 3 | 3 | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| C | 4 | 4 | 3 | 2 | 1 | 1 | 1 | 2 | 3 | 4 | 5 | 6 |
| A | 5 | 5 | 4 | 3 | 2 | 2 | 2 | 1 | 2 | 3 | 4 | 5 |
| A | 6 | 6 | 5 | 4 | 3 | 3 | 3 | 2 | 1 | 2 | 3 | 4 |
| T | 7 | 7 | 6 | 5 | 4 | 4 | 4 | 3 | 2 | 1 | 2 | 3 |
| C | 8 | 8 | 7 | 6 | 5 | 5 | 5 | 4 | 3 | 2 | 2 | 3 |
| C | 9 | 9 | 8 | 7 | 6 | 6 | 6 | 5 | 4 | 3 | 3 | 3 |
| C | 10 | 10 | 9 | 8 | 7 | 7 | 7 | 6 | 5 | 4 | 4 | 4 |

# Example

$$V[i,j] = \frac{(i+j)}{2} - 3D[i,j]$$

Match score $r_{mat}$ = +1
Mismatch score $r_{mis}$ = -2.
Indel score $r_{ind}$ = -2.5

| | _ | G | T | A | G | C | A | A | T | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|
| V | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ 0 | 0 | -2.5 | -5 | -7.5 | -10 | -12.5 | -15 | -17.5 | -20 | -22.5 | -25 |
| G 1 | -2.5 | 1 | -1.5 | -4 | -6.5 | -9 | -11.5 | -14 | -16.5 | -19 | -21.5 |
| T 2 | -5 | -1.5 | 2 | -0.5 | -3 | -5.5 | -8 | -10.5 | -13 | -15.5 | -18 |
| A 3 | -7.5 | -4 | -0.5 | 3 | 0.5 | -2 | -4.5 | -7 | -9.5 | -12 | -14.5 |
| C 4 | -10 | -6.5 | -3 | 0.5 | 1 | 1.5 | -1 | -3.5 | -6 | -8.5 | -11 |
| A 5 | -12.5 | -9 | -5.5 | -2 | -1.5 | -1 | 2.5 | 0 | -2.5 | -5 | -7.5 |
| A 6 | -15 | -11.5 | -8 | -4.5 | -4 | -3.5 | 0 | 3.5 | 1 | -1.5 | -4 |
| T 7 | -17.5 | -14 | -10.5 | -7 | -6.5 | -6 | -2.5 | 1 | 4.5 | 2 | -0.5 |
| C 8 | -20 | -16.5 | -13 | -9.5 | -9 | -8.5 | -5 | -1.5 | 2 | 2.5 | 0 |
| C 9 | -22.5 | -19 | -15.5 | -12 | -11.5 | -11 | -7.5 | -4 | -0.5 | 0 | 0.5 |
| C 10 | -25 | -21.5 | -18 | -14.5 | -14 | -13.5 | -10 | -6.5 | -3 | -2.5 | -2 |

| | _ | G | T | A | G | C | A | A | T | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|
| D | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ 0 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| G 1 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| T 2 | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A 3 | 3 | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| C 4 | 4 | 3 | 2 | 1 | 1 | 1 | 2 | 3 | 4 | 5 | 6 |
| A 5 | 5 | 4 | 3 | 2 | 2 | 2 | 1 | 2 | 3 | 4 | 5 |
| A 6 | 6 | 5 | 4 | 3 | 3 | 3 | 2 | 1 | 2 | 3 | 4 |
| T 7 | 7 | 6 | 5 | 4 | 4 | 4 | 3 | 2 | 1 | 2 | 3 |
| C 8 | 8 | 7 | 6 | 5 | 5 | 5 | 4 | 3 | 2 | 2 | 3 |
| C 9 | 9 | 8 | 7 | 6 | 6 | 6 | 5 | 4 | 3 | 3 | 3 |
| C 10 | 10 | 9 | 8 | 7 | 7 | 7 | 6 | 5 | 4 | 4 | 4 |

# Property of D table

- $D[i-1,j-1] \leq D[i,j] \leq D[i-1,j-1]+1$

- Diagonal k consists of entries D[i,j] such that i − j = k.

- Property: The values are monotonic increasing for every diagonal.

Diagonal 0

Diagonal 5

|   |    | _ | G | T | A | G | C | A | A | T | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
| **D** |    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| G | 1  | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| T | 2  | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 3  | 3 | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| C | 4  | 4 | 3 | 2 | 1 | 1 | 1 | 2 | 3 | 4 | 5 | 6 |
| A | 5  | 5 | 4 | 3 | 2 | 2 | 2 | 1 | 2 | 3 | 4 | 5 |
| A | 6  | 6 | 5 | 4 | 3 | 3 | 3 | 2 | 1 | 2 | 3 | 4 |
| T | 7  | 7 | 6 | 5 | 4 | 4 | 4 | 3 | 2 | 1 | 2 | 3 |
| C | 8  | 8 | 7 | 6 | 5 | 5 | 5 | 4 | 3 | 2 | 2 | 3 |
| C | 9  | 9 | 8 | 7 | 6 | 6 | 6 | 5 | 4 | 3 | 3 | 3 |
| C | 10 | 10 | 9 | 8 | 7 | 7 | 7 | 6 | 5 | 4 | 4 | 4 |

# Proof: D[i-1,j-1] $\leq$ D[i,j] $\leq$ D[i-1,j-1]+1

- **Lemma**: D[i,j] $\leq$ D[i-1,j-1]+1, D[i-1,j]+1, D[i,j-1]+1.
- Proof: D[i,j]= min{D[i-1,j-1]+$\delta$(S[i],T[j]), D[i-1,j]+1, D[i,j-1]+1}

- **Lemma**: D[i,j]$\leq$D[i+1,j+1].
- Proof: By induction. Suppose D[x,y] $\leq$D[x+1,y+1] for x+y<i+j.
- By contrary, suppose D[i,j]>D[i+1,j+1].
- Since D[i+1,j+1]=min{D[i,j]+$\delta$(S[i+1],T[j+1]), D[i+1,j]+1, D[i,j+1]+1}, we have either (1) D[i+1,j+1]=D[i+1,j]+1 or (2) D[i+1,j+1]=D[i,j+1]+1.
  - Case 1: By induction, D[i+1,j]$\geq$D[i,j-1].
  - Since D[i,j-1]+1 $\geq$D[i,j], we have D[i+1,j+1]=D[i+1,j]+1$\geq$D[i,j-1]+1 $\geq$D[i,j]. Contradiction.
  - Case 2: By induction, D[i,j+1]$\geq$D[i-1,j].
  - Since D[i-1,j]+1 $\geq$D[i,j], we have D[i+1,j+1]=D[i,j+1]+1$\geq$D[i-1,j]+1 $\geq$D[i,j]. Contradiction.
- The lemma is true.

# Diagonal and R(d,k)

- Diagonal k consists of entries D[i,j] such that $i - j = k$.
- Define $R(d,k) = \max\{ i \mid D[i,i-k]=d \}$.

- E.g.
  - Diagonal 0: $R(0,0)=3$, $R(1,0)=4$, $R(2,0)=8$, $R(3,0)=9$ and $R(4,0)=10$.
  - Diagonal 5: $R(0,5)=\ldots=R(4,5)=$invalid, $R(5,5)=8$, $R(6,5)=9$ and $R(7,5)=10$.

Diagonal 0

|   |    | _  | G | T | A | G | C | A | A | T | G | T  |
|---|----|----|---|---|---|---|---|---|---|---|---|----|
|   | **D** | 0  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| G | 1  | 1  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9  |
| T | 2  | 2  | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8  |
| A | 3  | 3  | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7  |
| C | 4  | 4  | 3 | 2 | 1 | 1 | 1 | 2 | 3 | 4 | 5 | 6  |
| A | 5  | 5  | 4 | 3 | 2 | 2 | 2 | 1 | 2 | 3 | 4 | 5  |
| A | 6  | 6  | 5 | 4 | 3 | 3 | 3 | 2 | 1 | 2 | 3 | 4  |
| T | 7  | 7  | 6 | 5 | 4 | 4 | 4 | 3 | 2 | 1 | 2 | 3  |
| C | 8  | 8  | 7 | 6 | 5 | 5 | 5 | 4 | 3 | 2 | 2 | 3  |
| C | 9  | 9  | 8 | 7 | 6 | 6 | 6 | 5 | 4 | 3 | 3 | 3  |
| C | 10 | 10 | 9 | 8 | 7 | 7 | 7 | 6 | 5 | 4 | 4 | 4  |

Diagonal 5

# Computing R(d,k)

- The computation of R(d,k) has two steps.

- 1. Find the biggest index i such that the alignment between S[1..i] and T[1..i-k] has d differences and the last aligned pair is not match.

- 2. Find R(d,k).


- Step 1 is computed based on the following lemma.

**Lemma**: Let i be the biggest index such that the alignment between S[1..i] and T[1..i-k] has d differences and the last aligned pair is not match. We have:

$$i = \max \left\{ \begin{array}{c} R(d - 1, k) + 1 \\ R(d - 1, k - 1) + 1 \\ R(d - 1, k + 1) \end{array} \right\}.$$

# Proof of the lemma

**Lemma**: Let i be the biggest index such that the alignment between S[1..i] and T[1..i-k] has d differences and the last aligned pair is not match. We have:

$$i = \max \left\{ \begin{array}{c} R(d-1, k) + 1 \\ R(d-1, k-1) + 1 \\ R(d-1, k+1) \end{array} \right\}.$$

- Proof: Three cases: Last aligned pair is (1) deletion, (2) mismatch and (3) insertion.
- The largest index i such that D[i,i-k]=d and the last aligned pair is (S[i],-), that is, deletion.
  - This means that i-1 is the biggest index such that D[i-1, i-k]=d-1. This means that i-1 = R(d-1, k-1).
- The largest index i such that D[i,i-k]=d and the last aligned pair is (S[i],T[i-k]) where S[i]≠T[i-k].
  - This means that i-1 is the biggest index such that D[i-1,i-1-k]=d-1. This means that i-1=R(d-1,k).
- The largest index i such that D[i,i-k]=d and the last aligned pair is (-,T[i-k]).
  - This means that i-1 is the biggest index such that D[i-1, i-k]=d-1. This means that i-1 = R(d-1, k-1).
- We choose the maximum index i among the three cases.

# Computing R(d,k)

- After step 1, i is the biggest index such that the alignment between S[1..i] and T[1..i-k] has d differences and the last aligned pair is not match.

- Suppose R(d,k)=j. This means that the alignment between S[i..j] and T[1..j-k] has d differences.

- If j>i, the alignment between S[i+1..j] and T[i-k+1..j-k] must be matches. Hence, S[i+1..j]=T[i-k+1..j-k].

- j can be found by Extend(i,k)

- Hence, R(d,k)= Extend(i, k).

```
Extend(i,k)
   while (S[i+1]=T[i-k+1]) {
      i++;
   }
   return i;
```

# Algorithm

- Set R(0,0)=extend(0, 0);
- d=0; $L_0=U_0=0$.
- Repeat
  - d=d+1
  - For k=$L_{d-1}$-1 to $U_{d-1}$+1
    - $i = \max \begin{Bmatrix} R(d-1,k)+1 \\ R(d-1,k-1)+1 \\ R(d-1,k+1) \end{Bmatrix}$
    - R(d,k)=Extend(i, k);
  - Set $(L_d,U_d)$ = be a subrange of $(L_{d-1}$-1,$U_{d-1}$+1) such that R(d,$L_d$) and R(d,$U_d$) are valid;
- Until $(L_d$-1,$U_d$+1) is an invalid range

# Example

- First, we handle d=0 (no mismatch).
- We set R(0,0) = Extend(0,0).
- We set $L_0=0, U_0=0$.

| D | _ 0 | G 1 | T 2 | A 3 | G 4 | C 5 | A 6 | A 7 | T 8 | G 9 | T 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| _ 0 | 0 | | | | | | | | | | |
| G 1 | | 0 | | | | | | | | | |
| T 2 | | | 0 | | | | | | | | |
| A 3 | | | | 0 | | | | | | | |
| C 4 | | | | | | | | | | | |
| A 5 | | | | | | | | | | | |
| A 6 | | | | | | | | | | | |
| T 7 | | | | | | | | | | | |
| C 8 | | | | | | | | | | | |
| C 9 | | | | | | | | | | | |
| C 10 | | | | | | | | | | | |

# Example

- Next, we handle d=1 (i.e. 1 difference)
- We need to find R(1,k) for $-1 = L_0 - 1 \leq k \leq U_0 + 1 = 1$.
- By $i = \max \begin{cases} R(d-1, k-1) + 1 \\ R(d-1, k) + 1 \\ R(d-1, k+1) \end{cases}$ and R(d,k)=Extend(i,k),

  we have:
- k=-1: i=R(0,0)=3, R(1,-1) = Extend(3, -1) = 7
- k=0: i=R(0,0)+1=4, R(1,0)=Extend(4, 0)=4
- k=1: i=R(0,0)+1=4, R(1,1)=Extend(4, 1)=4
- So, $L_1 = -1$ and $U_1 = 1$

|   |   | _ | G | T | A | G | C | A | A | T | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **D** |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0 | 0 | 1 |   |   |   |   |   |   |   |   |   |
| G | 1 | 1 | 0 | 1 |   |   |   |   |   |   |   |   |
| T | 2 |   | 1 | 0 | 1 |   |   |   |   |   |   |   |
| A | 3 |   |   | 1 | 0 | 1 |   |   |   |   |   |   |
| C | 4 |   |   |   | 1 | 1 | 1 |   |   |   |   |   |
| A | 5 |   |   |   |   |   |   | 1 |   |   |   |   |
| A | 6 |   |   |   |   |   |   |   | 1 |   |   |   |
| T | 7 |   |   |   |   |   |   |   |   | 1 |   |   |
| C | 8 |   |   |   |   |   |   |   |   |   |   |   |
| C | 9 |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Example

- Next, we handle d=2 (i.e. 2 difference)
- We need to find R(2,k) for -2=$L_1$-1 $\leq$ k $\leq$ $U_1$+1=2.
- By $i = \max \begin{cases} R(d-1, k-1) + 1 \\ R(d-1, k) + 1 \\ R(d-1, k+1) \end{cases}$ and R(d,k)=Extend(i,k),
  we have:
- k=-2: i=R(1,-1)=7, R(2,-2) = Extend(7, -2) = 7
- k=-1: i=max{R(1,-1)+1, R(1,0)}=8, R(2,-1)=Extend(8, -1)=8
- k=0: i=max{R(1,-1)+1, R(1,0)+1, R(1,1)}=8, R(2,0)=Extend(8, 0)=0
- K=1: i=max{R(1,1)+1,R(1,0)+1}=5, R(2,1)=Extend(5, 1)=5
- k=2: i=R(1,1)+1=5, R(2,2)=Extend(5, 2)=5
- So, $L_2$=-2 and $U_2$=2

| | | _ | G | T | A | G | C | A | A | T | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **D** | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0 | 0 | 1 | 2 | | | | | | | | |
| G | 1 | 1 | 0 | 1 | 2 | | | | | | | |
| T | 2 | 2 | 1 | 0 | 1 | 2 | | | | | | |
| A | 3 | | 2 | 1 | 0 | 1 | 2 | | | | | |
| C | 4 | | | 2 | 1 | 1 | 1 | 2 | | | | |
| A | 5 | | | | 2 | 2 | 2 | 1 | 2 | | | |
| A | 6 | | | | | | | 2 | 1 | 2 | | |
| T | 7 | | | | | | | | 2 | 1 | 2 | |
| C | 8 | | | | | | | | | 2 | 2 | |
| C | 9 | | | | | | | | | | | |
| C | 10 | | | | | | | | | | | |

# Example

- Next, we handle d=3 (i.e. 3 difference)
- We need to find R(3,k) for $-3=L_2-1 \leq k \leq U_2+1=3$.
- By $i = \max \begin{cases} R(d-1,k-1)+1 \\ R(d-1,k)+1 \\ R(d-1,k+1) \end{cases}$ and R(d,k)=Extend(i,k), we have:
- k=-3: i=R(2,-2)=7, R(3,-3) = Extend(7, -3) = 10
- k=-2: i=max{R(2,-2)+1, R(2,-1)}=8, R(3,-2)=Extend(8,-2)=8
- k=-1: i=max{R(2,-2)+1, R(2,-1)+1, R(2,0)}=9, R(3,-1)=Extend(9,-1)=9
- k=0: i=max{R(2,-1)+1, R(2,0)+1, R(2,1)}=9, R(3,0)=Extend(9,0)=9
- k=1: i=max{R(2,0)+1, R(2,1)+1, R(2,2)}=9, R(3,1)=Extend(9, 1)=9
- k=2: i=max{R(2,1)+1, R(2,2)+1}=6, R(3,2)=Extend(6,2)=6
- k=3: i=R(2,2)+1=6, R(3,3)=Extend(6,3)=6
- So, $L_3$=-3 and $U_3$=3



| | | _ | G | T | A | G | C | A | A | T | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **D** | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0 | 0 | 1 | 2 | 3 | | | | | | | |
| G | 1 | 1 | 0 | 1 | 2 | 3 | | | | | | |
| T | 2 | 2 | 1 | 0 | 1 | 2 | 3 | | | | | |
| A | 3 | 3 | 2 | 1 | 0 | 1 | 2 | 3 | | | | |
| C | 4 | | 3 | 2 | 1 | 1 | 1 | 2 | 3 | | | |
| A | 5 | | | 3 | 2 | 2 | 2 | 1 | 2 | 3 | | |
| A | 6 | | | | 3 | 3 | 3 | 2 | 1 | 2 | 3 | |
| T | 7 | | | | | | | 3 | 2 | 1 | 2 | 3 |
| C | 8 | | | | | | | | 3 | 2 | 2 | 3 |
| C | 9 | | | | | | | | | 3 | 3 | 3 |
| C | 10 | | | | | | | | | | | |

# Example

- When we handle d=10 (10 differences), we obtain the final table.

|   |   | _ | G | T | A | G | C | A | A | T | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **D** |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| G | 1 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| T | 2 | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 3 | 3 | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| C | 4 | 4 | 3 | 2 | 1 | 1 | 1 | 2 | 3 | 4 | 5 | 6 |
| A | 5 | 5 | 4 | 3 | 2 | 2 | 2 | 1 | 2 | 3 | 4 | 5 |
| A | 6 | 6 | 5 | 4 | 3 | 3 | 3 | 2 | 1 | 2 | 3 | 4 |
| T | 7 | 7 | 6 | 5 | 4 | 4 | 4 | 3 | 2 | 1 | 2 | 3 |
| C | 8 | 8 | 7 | 6 | 5 | 5 | 5 | 4 | 3 | 2 | 2 | 3 |
| C | 9 | 9 | 8 | 7 | 6 | 6 | 6 | 5 | 4 | 3 | 3 | 3 |
| C | 10 | 10 | 9 | 8 | 7 | 7 | 7 | 6 | 5 | 4 | 4 | 4 |

# X-drop criteria

- X-drop criteria sets V[i,j]=invalid if V[i,j]<$\tau$-X, where
  $\tau$ = max{ max{V[p,q] | p+q<i+j }, max{ V[p,q] + $\delta$(S[p+1],T[q+1])/2 | p+q=i+j-1} }.

- However, $\tau$ is time consuming to compute.

- Below lemma gives an alternative.

**Lemma**: Suppose $D[i,j] = d$. Let $d' = d - \left\lceil \dfrac{X + \frac{r_{mat}}{2}}{r_{mat} - r_{mis}} \right\rceil - 1$. Let $\tau[d'] = $ max$\{V[p,q]|D[p,q] \leq d'\}$. V[i,j]<$\tau$-X if and only if V[i,j]<$\tau$[d']-X.

- ($\rightarrow$) Suppose V[i,j]<$\tau$-X. Then, there exist p,q such that p+q<i+j and V[i,j] < V[p,q] − X.

- We claim that such (p,q) pair satisfies $D[p,q] \leq d'$.

- D[p,q]=((p+q)$r_{mat}$/2 − V[p,q])/($r_{mat}$ − $r_{mis}$)

-       < ((p+q)$r_{mat}$/2 − V[i,j] − X)/($r_{mat}$ − $r_{mis}$)

-       < ((i+j)$r_{mat}$/2 − V[i,j] − X)/($r_{mat}$ − $r_{mis}$)

-       = D[i,j] − (X + $r_{mat}$/2)/($r_{mat}$ − $r_{mis}$)       $\boxed{\text{V[i,j] = (n+m)}r_{mat}\text{/2 − D[i,j] (}r_{mat}\text{ − }r_{mis}\text{)}}$

-       = d − (X + $r_{mat}$/2)/($r_{mat}$ − $r_{mis}$) = d'.

- ($\leftarrow$) Suppose V[i,j]<$\tau$[d']-X.
- Note that $\tau$[d']=V[p,q] for some D[p,q]$\leq$d'.
- If p+q<i+j, then V[i,j]<$V[p,q] - X \leq \tau - X$.

$$d - d' = \left\lfloor \frac{X + \frac{r_{mat}}{2}}{r_{mat} - r_{mis}} \right\rfloor + 1 > \frac{X + \frac{r_{mat}}{2}}{r_{mat} - r_{mis}}$$

- If (p+q)$\geq$i+j, Select (p',q') on antidiagonal-(i+j-1) such that p-p'=q-q'. Note that D[p',q']$\leq$d'.
- V[p',q']-V[i,j] = $(p'+q')r_{mat}/2 - D[p',q'] (r_{mat} - r_{mis}) - (i+j)r_{mat}/2 - d (r_{mat} - r_{mis})$
-  $\geq$ - $r_{mat}/2 + (d - d') (r_{mat} - r_{mis})$
-  $\geq -r_{mat}/2 + \left( \frac{X + \frac{r_{mat}}{2}}{r_{mat} - r_{mis}} \right) (r_{mat} - r_{mis}) > X$

- For antidiagonal-k, let i=B(d,k) and i'=R(d,k).
- This means that D[i,i-k]=D[i+1,i+1-k]=…=D[i',i'-k]=d.
- Let $d' = d - \left\lfloor \dfrac{X + \frac{r_{mat}}{2}}{r_{mat} - r_{mis}} \right\rfloor - 1$.
- Lemma: If V[i,i-k]$\geq\tau$[d']-X, then V[i+s,i+s-k] $\geq\tau$[d']-X for 0$\leq$s$\leq$i'-i.
- Proof: It follows from the fact that V[i+s,i+s-k] = V[i,i-k]+s * r$_{mat}$.

# Algorithm

- Set i=extend(0, 0); R(0,0)=i; $\tau'$=V[i, i]; $\tau[0]$=V[i, i];

- d=0; $L_0=U_0=0$.

- Repeat

  - d=d+1; $d' = d - \left\lfloor \dfrac{X+\frac{r_{mat}}{2}}{r_{mat}-r_{mis}} \right\rfloor - 1$;

  - For k=$L_{d-1}$-1 to $U_{d-1}$+1

    - If V[i,i-k]<$\tau[d']$

      - R(d,k)=invalid

    - else

      - $i = \max \begin{Bmatrix} R(d-1, k) + 1 \\ R(d-1, k-1) + 1 \\ R(d-1, k+1) \end{Bmatrix}$

      - i'=Extend(i, k); R(d, k)=i'; $\tau'$ = max{$\tau'$, V[i', i'-k] };

  - $\tau[d]=\tau'$

  - Set $(L_d,U_d)$ = be a subrange of $(L_{d-1}$-1,$U_{d-1}$+1) such that R(d,$L_d$) and R(d,$U_d$) are valid;

- Until $(L_d$-1,$U_d$+1) is an invalid range

# Illustration of greedy algorithm with X-drop

$$V[i,j] = \frac{(i+j)}{2} - 3D[i,j]$$

- Suppose X=1.

- Then, $d' = d - \left\lfloor \dfrac{X+\frac{r_{mat}}{2}}{r_{mat}-r_{mis}} \right\rfloor - 1 = d - \left\lfloor \dfrac{1+\frac{1}{2}}{1-(-2)} \right\rfloor - 1 = d - 1$

- First, we handle d=0 (no mismatch).
- We set R(0,0) = Extend(0,0)=3.
- We set $L_0$=0, $U_0$=0. $\tau[0]$=0. $\tau'$=V[3,3]=3.

|   |    | _ | G | T | A | G | C | A | A | T | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
|   | D  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 |   |   |   |   |   |   |   |   |   |   |
| G | 1  |   | 0 |   |   |   |   |   |   |   |   |   |
| T | 2  |   |   | 0 |   |   |   |   |   |   |   |   |
| A | 3  |   |   |   | 0 |   |   |   |   |   |   |   |
| C | 4  |   |   |   |   |   |   |   |   |   |   |   |
| A | 5  |   |   |   |   |   |   |   |   |   |   |   |
| A | 6  |   |   |   |   |   |   |   |   |   |   |   |
| T | 7  |   |   |   |   |   |   |   |   |   |   |   |
| C | 8  |   |   |   |   |   |   |   |   |   |   |   |
| C | 9  |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Diagonal and R(d,k)

- Next, we handle d=1 (i.e. 1 difference)

- $d' = d - 1 = 0$

- We need to find R(1,k) for -1=$L_0$-1 $\leq$ k $\leq$ $U_0$+1=1.

- By $i = \max \begin{cases} R(d-1, k-1) + 1 \\ R(d-1, k) + 1 \\ R(d-1, k+1) \end{cases}$ and R(d,k)=extend(i,i-k),

  we have:

- k=-1: i=R(0,0)=3, V[3,3+1]=0.5, R(1,-1) = Extend(3, -1) = 7, V[7,7+1]=4.5

- k=0: i=R(0,0)+1=4, V[4,4]=1, R(1,0)=Extend(4, 0)=4

- k=1: i=R(0,0)+1=4, V[4,3]=0.5, R(1,1)=Extend(4, 1)=4

- So, $L_1$=-1 and $U_1$=1

- $\tau'$=max{0, 4.5, 1, 0.5}=4.5

- $\tau[1]$=4.5

|  |  | _ | G | T | A | G | C | A | A | T | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | D | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0 | 0 |  |  |  |  |  |  |  |  |  |  |
| G | 1 |  | 0 |  |  |  |  |  |  |  |  |  |
| T | 2 |  |  | 0 |  |  |  |  |  |  |  |  |
| A | 3 |  |  |  | 0 | 1 |  |  |  |  |  |  |
| C | 4 |  |  |  | 1 | 1 | 1 |  |  |  |  |  |
| A | 5 |  |  |  |  |  |  | 1 |  |  |  |  |
| A | 6 |  |  |  |  |  |  |  | 1 |  |  |  |
| T | 7 |  |  |  |  |  |  |  |  | 1 |  |  |
| C | 8 |  |  |  |  |  |  |  |  |  |  |  |
| C | 9 |  |  |  |  |  |  |  |  |  |  |  |
| C | 10 |  |  |  |  |  |  |  |  |  |  |  |

# Diagonal and R(d,k)

- Next, we handle k=2 (i.e. 2 difference)

- We need to find R(2,k) for $-2=L_1-1 \leq k \leq U_1+1=2$.

- By $i = \max \begin{cases} R(d-1, k-1)+1 \\ R(d-1, k)+1 \\ R(d-1, k+1) \end{cases}$ and R(d,k)=extend(i,i-k),

  we have:

- k=-2: i=R(1,-1)=7, V[7,7+2]=2, R(2,-2) = invalid

- k=-1: i=max{R(1,-1)+1, R(1,0)}=8, V[8,8+1]=2.5, R(2,-1)= invalid

- k=0: i=max{R(1,-1)+1, R(1,0)+1, R(1,1)}=8, V[8,8]=2, R(2,0)= invalid

- k=1: i=max{R(1,1)+1,R(1,0)+1}=5, V[5,4]=-1.5<$\tau[1]$-X, R(2,1)=invalid

- k=2: i=R(1,1)+1=5, V[5,3]=-2 <$\tau[1]$-X, R(2,2)=invalid

- So, $(L_2, U_2)$ is an invalid range.

|   |    | _ | G | T | A | G | C | A | A | T | G | T |
|---|----|---|---|---|---|---|---|---|---|---|---|---|
|   | D  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| _ | 0  | 0 |   |   |   |   |   |   |   |   |   |   |
| G | 1  |   | 0 |   |   |   |   |   |   |   |   |   |
| T | 2  |   |   | 0 |   |   |   |   |   |   |   |   |
| A | 3  |   |   |   | 0 | 1 |   |   |   |   |   |   |
| C | 4  |   |   |   | 1 | 1 | 1 |   |   |   |   |   |
| A | 5  |   |   |   | X | X | 2 | 1 |   |   |   |   |
| A | 6  |   |   |   |   |   |   | 2 | 1 |   |   |   |
| T | 7  |   |   |   |   |   |   |   | 2 | 1 | X |   |
| C | 8  |   |   |   |   |   |   |   |   | X | X |   |
| C | 9  |   |   |   |   |   |   |   |   |   |   |   |
| C | 10 |   |   |   |   |   |   |   |   |   |   |   |

# Back-tracking

- Find the entry with score $\tau[1]=4.5$, which is V[7,8].

- By back-tracking, we have

```
GTAGCAAT
GTA-CAAT
```

# BLAT

- BLAT aims to improve the efficiency of BLAST.
- Only for DNA.
- The main trick is to index the database and put the index in the main memory

- The algorithm is as follows.
1. For every w-tuple in the query, identify hits using the index of the database (w=11)
2. Accept hits satisfying the two-hit requirement.
3. Gapped extension.

- Note that BLAT is less sensitive than BLAST, but more sensitive than MegaBLAST.

# Main trick of BLAT

- BLAST cannot build index of human genome since it it big.
- BLAT's index stores the positions of non-overlapping w-tuples in memory.

Database = ACTTGTACTTGTACTTGTA

### Index of all w-mers

| w-mer | positions |
|-------|-----------|
| ACTT | 1, 7, 13 |
| CTTG | 2, 8, 14 |
| GTAC | 5, 11 |
| TACT | 6, 12 |
| TGTA | 4, 10 |
| TTGT | 3, 9, 15 |
| TGTA | 16 |

Store n integers

### Index of w-mers at positions iw+1

| w-mer | positions |
|-------|-----------|
| ACTT | 1, 13 |
| GTAC | 5 |
| TTGT | 9 |

Store n/w integers

# About the inventor: Jim Kent



- Education: University of California, Santa Cruz

- Awards: Overton Prize, Benjamin Franklin Award

# PatternHunter

- PatternHunter can only apply to DNA
- PatternHunter is similar to BLAST. Moreover, it uses gapped w-tuple.
  - For w=11, they use 111010010100110111
  - Example,

```
111010010100110111
ACTCCGATATGCGGTAAC
|||-|--|-|--||-|||
ACTTCACTGTGAGGCAAC
```

- They found that gapped w-tuple can increase the sensitivity while increase the efficiency.

# Advantage of gapped w-tuple (I)

- Increase sensitivity
  - Gapped w-tuples are more independent.
  - Examples:
    - Two adjacent ungapped 11-tuples share 10 symbols
      - ```
        11111111111
        11111111111
        ```
        1/4 chances to have 2nd hit next to the 1st hit
    - Two adjacent gapped 11-tuples share 5 symbols
      - ```
        111010010100110111
        111010010100110111
        ```
        $1/4^6$ chances to have 2nd hit next to the 1st hit
  - If the w-tuples are more independent, the probability of having at least one hit in a homologous region is higher.

# Advantage of gapped w-tuple (II)

- Reduce the number of hits.
  - For the same query length (says, 64),
    - It covers by 54 ungapped 11-tuples
    - It covers by 47 gapped 11-tuples
  - So, the number of hits is smaller.
- Thus, the efficiency is increased!

# PatternHunter I

Proposition. The expected number of hits of a weight-$W$ length-$M$ model within a length-$L$ region of similarity $p$ is $(L - M + 1) * p^W$

Proof.

For any fixed position, the prob of a hit is $p^W$.

There are $L - M + 1$ candidate positions.

The proposition follows.

# Implication

- For *L = 1017*
  - BLAST seed expects *(1017 − 11 + 1) \* p^{11} = 1007 \* p^{11}* hits
  - But ~*1/4* of these overlap each other. So likely to have only ~*750 \* p^{11}* distinct hits
  - Our example spaced seed expects *(1017 − 18 + 1) \* p^{11} = 1000 \* p^{11}* hits
  - But only *1/4^6* of these overlap each other. So likely to have ~*1000 \* p^{11}* distinct hits

Spaced seeds likely to be more sensitive & more efficient

# Sensitivity of PatternHunter I



Image credit: Li

# More for PatternHunter

- To further improve the efficiency,
  - PatternHunter uses a variety of advanced data structures including priority queues, a variation of red-black tree, queues, hash tables.
  - PatternHunter also uses a new method of sequence alignment.
- To further improve the accuracy,
  - PatternHunter II suggested to use multiple gapped seeds.
  - They show that the accuracy can approach smith-waterman algorithm while the speed 3000 times faster than smith-waterman.
- PatternHunter II is both faster and sensitive than BLAST, MegaBLAST.

# About the Inventor: Ming Li

- Ming Li
  - Canada Research Chair Professor of Bioinformatics, University Professor, Univ of Waterloo
  - Fellow, Royal Society of Canada. Fellow, ACM. Fellow, IEEE.

# PSI-BLAST (Position Specific Iterated BLAST)

- PSI-BLAST is an implementation of BLAST for finding protein families. It allows us to detect distant homology.

- Input: a protein sequence
  - Using BLAST, we get a set of sequences that align with the query protein with E-score below a threshold, 0.01 (by default).
  - Align the selected sequences
  - Generate a PSSM profile from the multiple alignment
  - Iterate until no significant alignment found,
    - Using a modified BLAST, search the database with the PSSM profile.
    - Align the selected sequences
    - Generate a PSSM from the multiple alignment

- This version automatically combines statistically significant alignments produced by BLAST into a position-specific score matrix (PSSM).

- It is much more sensitive to weak but biologically relevant sequence similarities

Input: Protein sequence

↓

BLAST

↓

Select high
similarity sequences

↓

Multiple alignment
And
Compute PSSM

↓

PSSM of the
aligned sequences

# Idea of PSI-BLAST (I)

- 1. Find a set of similar sequences of the query Q by BLAST
- 2. Use the similar sequences to generate a representative and set it to be Q.
- 3. Repeat 1 and 2.



```
Seq1: ATCGTTACTG
Seq2: CCCGTATCTG
Seq3: GACGTCCCTC
Seq4: TGCGTGGCTA
Seq5: CCCGTAACTT
Seq6: ATCGTCTCTA
PSSM: NNCGTNNCTN
```

Set Q to be the PSSM and rerun BLAST.

# Idea of PSSM (II)

- The PSSM will capture all the conserved amino acid residues (those residues form conserved regions like domains).

- When we run BLAST again, we can find the distant homologous sequences.



```
Seq1:  ATCGTTACTG
Seq2:  CCCGTATCTG
Seq3:  GACGTCCCTC
Seq4:  TGCGTGGCTA
Seq5:  CCCGTAACTT
Seq6:  ATCGTCTCTA
PSSM:  NNCGTNNCTN
```

Set Q to be the PSSM and rerun BLAST.

PSSM captures the conserved regions of all sequences.

# 1. Find a set of sequences similar to the query

- Using BLAST 2.0, we get a set of sequences that align with the query protein with E-score below a threshold, 0.01 (by default).

- In the set of selected sequences, some have >98% identical. We keep one copy for those selected sequences which are >98% identical.

# 2. Multiple sequence alignment of the selected sequences

- Using the query sequence as the template, we aligned the selected sequences.
- All gap characters inserted into the query sequence are ignored.
- Note:
  - the length of the alignment is the same as the query sequence.
  - Some columns of the multiple sequence alignment may include nothing except the query.

query

# 3. Generate a PSSM profile from the alignment

- Given the multiple alignment of length n,
  - We generate the position-specific score matrix (PSSM) profile, which is a 20xn matrix.
  - For each column and each residue a in the profile, we generate a log-odds score $\log(O_{ia}/P_a)$.
    - where $O_{ia}$ is the observed frequency of residue a at position i and $P_a$ is the expected frequency respectively of the residue a.
- Since number of sequences may be small, data-dependent pseudo frequency is added to $O_{ia}$.

# Example

- This is an example PSSM M.

- For position p, amino acid a,
$M(p, a) = \sum_{b=1}^{20} W(p, b)\delta(a, b)$, where

  - $\delta(a, b)$ is the substitution matrix,

  - $W(p, b) = \dfrac{n(p,b)}{n(p)}$,

  - n(p) is the number of probes in position p and

  - n(p,b) is the number of probes with amino acid b at position p.



FIG. 1. The concept of a profile. (a) A flow diagram of profile analysis. (b) A 49-residue sample profile for the immunoglobulin variable-region domain, generated from the four-probe sequences shown at the left (see Fig. 2b for details). The profile is shown in the box. The rightmost column of the profile gives the penalty for insertion/deletion (+/−). Positions 31–47 of the profile are omitted from the figure for clarity. Notice that where gaps appear in some of the probe sequences, the insertion/deletion penalty is lower than elsewhere.

Profile analysis: Detection of distantly related proteins. PNAS 1987.

# 4. Run BLAST again with the PSSM profile

- We apply a modified BLAST to the PSSM profile.
  - Basically, when we compare a position of the PSSM and a residue in the database, we use the corresponding log-odds score in that position.

- Repeat until we satisfy.

# QUASAR

- QUASAR stands for Q-gram Alignment based on Suffix ARrays
- It is a good searching tool for identifying strong similar strings.
- Problem:
  - Input: a database D, a query S, k, w
  - Output: a set of (x, y) where
    - x and y are length-w substring in D and S, respectively
    - edit_dist(x, y) $\leq$ k

# Observation

```
gcagactgctac  k=3
gccgacagccac  w=12
              q=3
              t=w+1-(k+1)q
               =12+1-(3+1)3
               =1
```

Lemma:

Given two length-w sequences X and Y, if they have edit distance $\leq$ k, then they share at least t common q-grams (length-q substrings) where t = w+1-(k+1)q.

Proof:

- Suppose X and Y has r differences.
- X has (w+1-q) q-grams
- Note that a q-gram in X overlaps with some difference iff X and Y does not share that q-gram
- For each difference, there are at most q q-grams overlap with the difference. In total, rq q-grams overlap with the r differences
- Thus, X and Y share (w+1-q – rq) q-grams, which is bigger than w+1-(k+1)q.

- We make use of this observation to do filtering!

# Algorithm for finding potential approximate matches of S in D

- For every X = S[i..i+w-1]  of the query where i=1, 2, …
  - For every length-w substring Y in D, associate a counter with it and initialize it to zero
  - For each q-gram Q in X,
    - Find the hitlist, that is, the list of positions in D so that Q occurs
    - Increment the counter for every length-w substring Y which contains Q
  - For every length-w substring Y in D with counter > t, X and Y are potential approximate match. We check it using an alignment algorithm!

# Illustration of the algorithm

# How to get the hitlist?

- Based on the data-structures
  - A suffix array SA of the database D is the lexicographically ordered sequence of all suffixes in D.
  - An auxiliary array idx where for each q-gram Q, idx[Q] is the start of the hitlist for Q!

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| Database D = | C | A | G | C | A | C | T |

| i | SA[i] | |
|---|---|---|
| 1 | 5 | ACT |
| 2 | 2 | AGCACT |
| 3 | 4 | CACT |
| 4 | 1 | CAGCACT |
| 5 | 6 | CT |
| 6 | 3 | GCACT |
| 7 | 7 | T |

idx(AC)

idx(AG)

idx(CA)

idx(CT)

idx(GC)

# Speedup Feature 1: Window shifting

- In the previuos algorithm, building the counters list for S[i..w+i-1] is time consuming!

- Suppose the counters list for S[1..w] is given, can we determine the counters list for S[2..w+1] easily?
  - Idea: For every length-w string Y in D,
    - Decrement counter for Y if it contains q-gram S[1..q]
    - Increment counter for Y if it contains q-gram S[w-q+2..w+1]

- The window shifting idea reduce the time complexity.

Q-gram S[1..q]

S[1..w]

Q-gram S[w-q+2..w+1]

S[2..w+1]

# Speedup Feature 2: Block addressing

- Another problem: too many counters

- Solution (Block addressing scheme):
  - Instead of associate a counter for every length-w substring Y in D
  - The database D is divided into blocks of size b (b $\geq$ 2w). Each block is associating a counter.
  - If a block contains more than t q-grams, this block has to be checked for approximate matches using an alignment algorithm

# Weakness of QUASAR

- Extensive memory requirement
  - Construction phase:
    - Memory space = 9n (where n is DB size)
  - Query phase:
    - Memory space = 5n
- Not suitable for distant homologous sequences
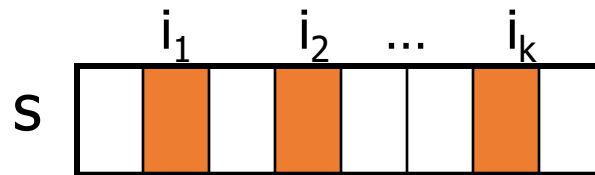
# Locality-Sensitive Hashing (LSH)

**LSH-ALL-PAIRS**

- Input: biosequence database D

- Aim: find pairs of w-mers that differ by at most d substitutions (ungapped local alignment) in a collection of biosequences D.
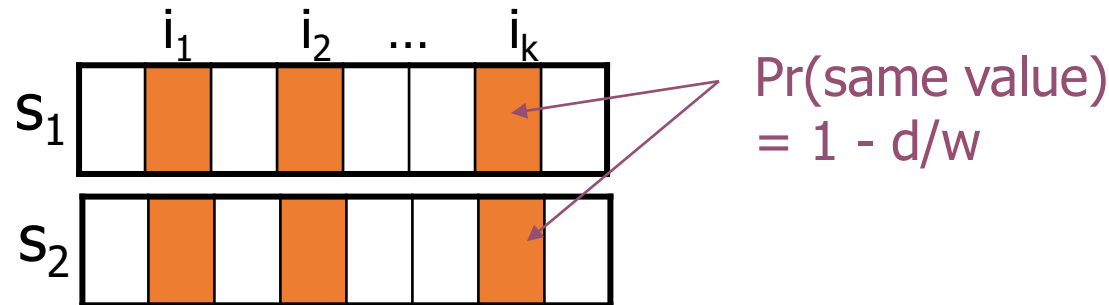
# Locality-sensitive hash function

- Consider an w-mers s,
  - choose k indices $i_1$, $i_2$, ..., $i_k$ uniformly from the set {1, 2, ..., w}
  - Define $\pi(s) = (s[i_1], s[i_2], ..., s[i_k])$. This function is called the locality-sensitive hash function

# Property of locality-sensitive hash function (I)

- Consider two w-mers $s_1$ and $s_2$,
  - the more similar are they, the higher probability that $\pi(s_1) = \pi(s_2)$.
- More precisely, if the hamming distance of $s_1$ and $s_2$ = d,
  - $\Pr[\pi(s_1) = \pi(s_2)] = \prod_{j=1,...,k} \Pr[s_1[i_j] = s_2[i_j]]$
    $= (1 - d/w)^k$

# Property of locality-sensitive hash function (II)

- Hence, $s_1$ and $s_2$ are similar if
  - $\pi(s_1) = \pi(s_2)$
- However, we may have false positive and false negative
  - False positive: $s_1$ and $s_2$ are dissimilar but $\pi(s_1) = \pi(s_2)$.
    - False positive can be distinguished from true positive by computing hamming distance between $s_1$ and $s_2$
  - False negative: $s_1$ and $s_2$ are similar but $\pi(s_1) \neq \pi(s_2)$.
    - We cannot detect false negative.
    - We can only reduce the number of false negative by repeating the test using different $\pi()$ functions

# LSH-ALL-PAIRS

**Algorithm:**

1. Generate m random locality-sensitive hash functions $\pi_1(\ )$, $\pi_2(\ )$, ..., $\pi_m(\ )$.
2. For every w-mer s in the database, compute $\pi_1(s)$, $\pi_2(s)$, ..., $\pi_m(s)$.
3. For every pair of w-mers s and t such that $\pi_j(s) = \pi_j(t)$ for some j,
   - If hamming distance(s, t) < d, report (s, t)-pair.

# Conclusion

- This lecture presents some database searching methods.
- In fact, there are many other methods. For examples:
  - CAFÉ, FLASH, RAMdb, FD, suffix tree, suffix array, compressed suffix array

# More information

- Altschul, Madden, Schaffer, Zhang, Zhang, Miller and Lipman. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. NAR, 1997.

- Zhang, Schwartz, Wagner and Miller. A Greedy Algorithm for Aligning DNA sequences. Journal of Computational Biology, 2000.

- The list of database used by blast
  - ftp://ftp.ncbi.nlm.nih.gov/blast/db/