# BLAST

Basic Local Alignment Search Tool

So useful – it is now a verb in the literature

# Goals for Today:

What is BLAST and why is it important?

Principles of the algorithm

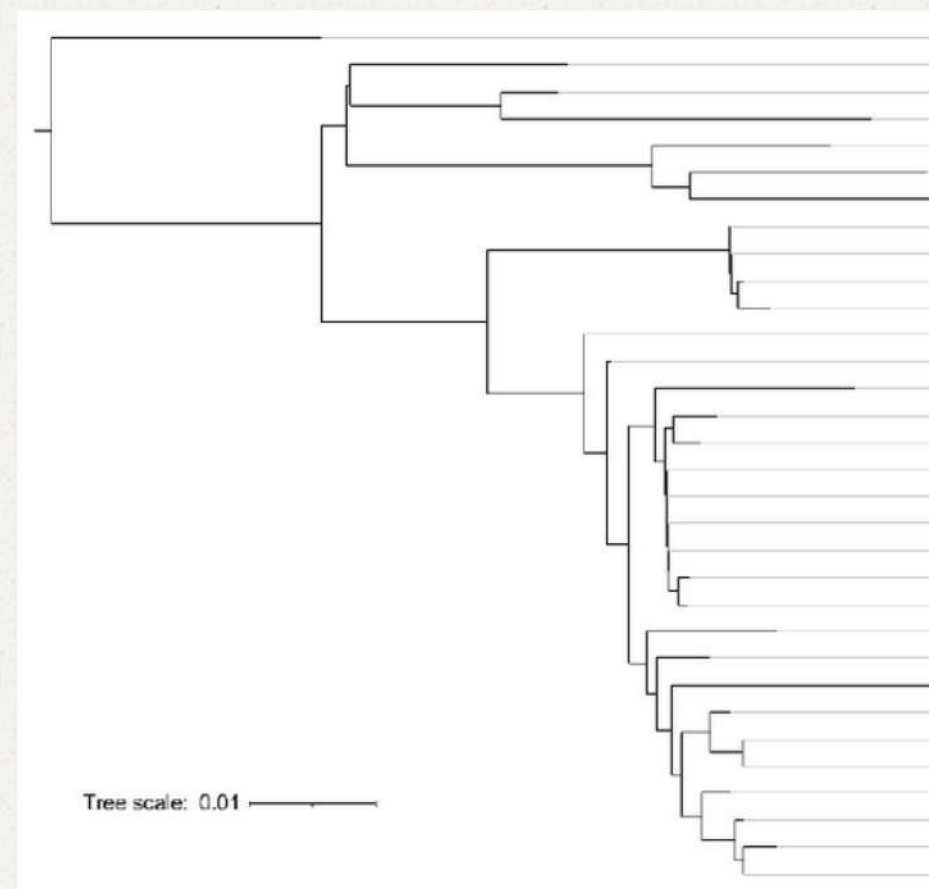Online Examples

Command Line Implementation

# A Lot of Blasting

Van Noorden et al. 2014, *Nature*

- • Where is BLAST on this list?
  - • Altschul et al. 1990
    - • #12 – 38,380 citations
      - • 53,672 (Web of Science 6/1/2019)
  - • Altschul et al. 1997
    - • #14 – 36,410 citations
      - • 48,001 (Web of Science 6/1/2019)
  - • Combined: 4th!

# Goals

- Search a query against a database
  - Identify species
  - Locate domains
  - Assess function
  - Establish phylogeny
  - Mapping
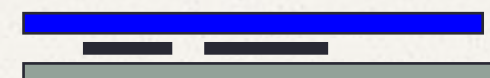
Tree scale: 0.01

# What is BLAST?

- Sequence searching algorithm

- Finds the best local alignments

- Calculates statistical significance

- Similarity suggests homology

- Less sensitive than Smith-Waterman, but FASTER!

### Global vs Local Alignment



- Global alignment: entire sequences

- Local alignment: segments of sequences

- Local alignment often the most relevant
  - Depends on biological assumptions

# Blast Types

| Name | Query | Database |
|------|-------|----------|
| blastn | nucleotide | nucleotide |
| blastp | protein | protein |
| blastx | nucleotide | protein |
| tblastx | nucleotide | nucleotide |
| tblastn | protein | nucleotide |
| PSI-blast | protein | protein |

# Blast Databases: Protein

| Name | Host | Description |
| --- | --- | --- |
| nr | NCBI | Non-redundant, general |
| Refseq_protein | NCBI | Annotated and curated protein collection |
| SwissProt | SIB | Manually curated and reviewed proteins form UniProt |
| Trembl | EBI | Automatically annotated, non-reviewed proteins |
| PDB | Rutgers/UCSD/UCSC | Proteins with 3D structural information |

# Blast Databases: Nucleotide

| Name | Host | Description |
|---|---|---|
| nt | NCBI | Non-redundant, general |
| Refseq_RNA | NCBI | Annotated and curated RNA sequence collection |
| Refseq_Genomics | NCBI | Sequenced and curated genomes |
| EST | NCBI | Expressed sequence tags |
| UNIVEC | NCBI | Vector contaminant database |
| WGS | NCBI | Draft, whole genome shotgun sequence assemblies |
| SRA | NCBI | Raw NGS datasets |
| Many more databases, e.g. barcoding, viral, tRNA, etc, custom-built databases | | |

# How it works: Making words

<div>

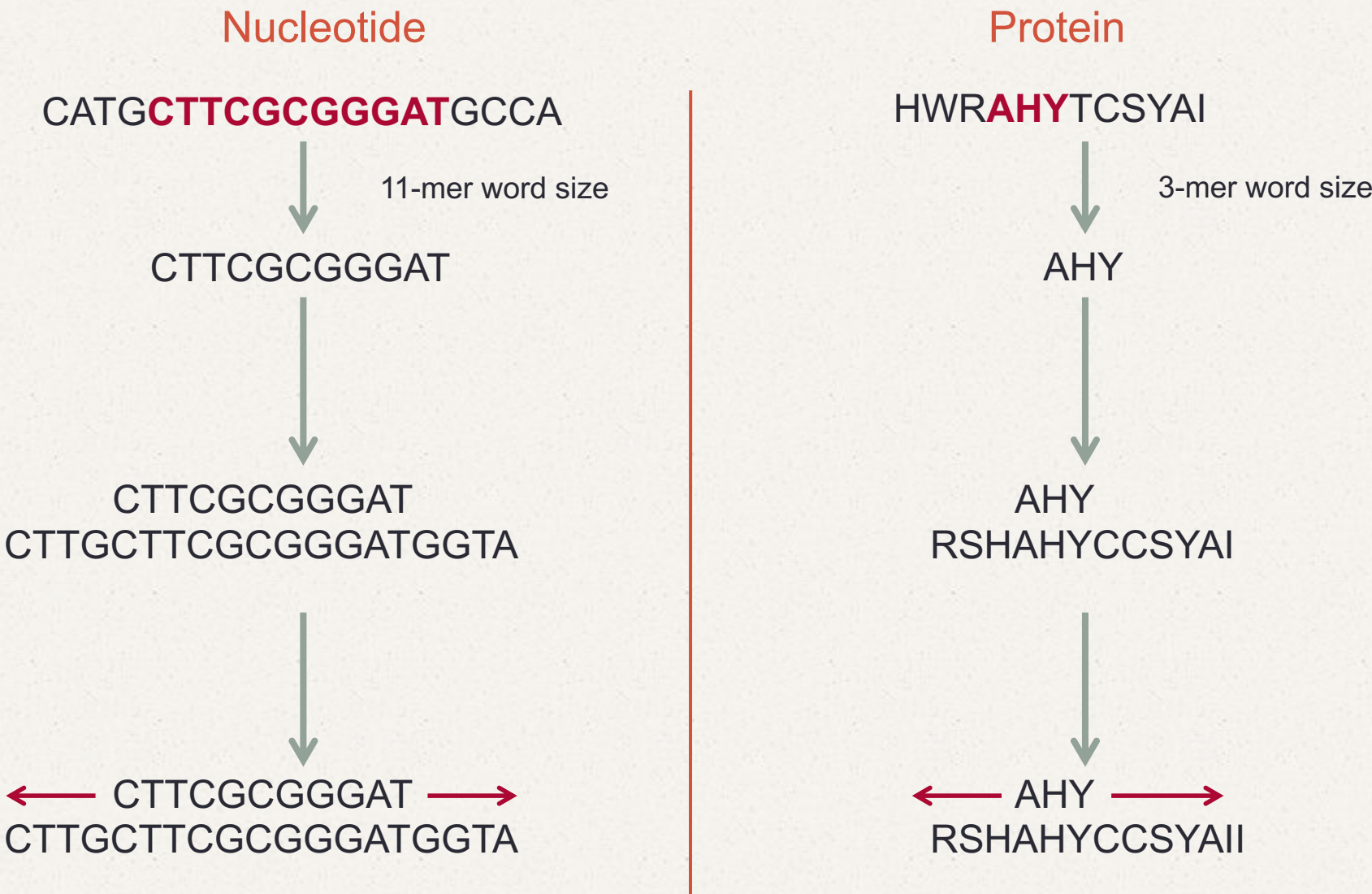**Nucleotide**

- 11-letter words (seeds)
  - ACTACGTGCTATGC
    - ACTACGTGCTA
    - CTACGTGCTAT
    - TACGTGCTATG
    - ACGTGCTATGC

</div>

<div>

**Protein**

- 3-letter words (seeds)
  - PQGDEF
    - PQG
    - QGD
    - GDE
    - DEF

</div>

# How it works

## Nucleotide

CATG**CTTCGCGGGAT**GCCA

↓ 11-mer word size

CTTCGCGGGAT

↓

CTTCGCGGGAT
CTTGCTTCGCGGGATGGTA

↓

←— CTTCGCGGGAT —→
CTTGCTTCGCGGGATGGTA

## Protein

HWR**AHY**TCSYAI

↓ 3-mer word size

AHY

↓

AHY
RSHAHYCCSYAI

↓

←— AHY —→
RSHAHYCCSYAII

# Blast Scoring and E-values

- Nucleotide sequences search for 11-letter matches
  - (4^11  =  4,194,304 combinations)
  - Match = +5, mismatch = -4
  - Only scores above a threshold (T) are kept

ACTACGTGCTA
ACTACGTGCTA
5+5+5+5+5+5+5+5+5+5+5 = 55

ACTACGTGCTA
ACAAGATGGTA
5+5-4+5-4-4+5+5-4+5+5 = 19

# Blast Scoring and E-values

- ## Proteins use a BLOSUM62 scoring matrix
  - ### 20x20x20 = 8,000 possible 3-letter words
  - ### All possible amino acid pairs are given a score
  - ### All combinations above a threshold (T) are kept
    - #### Minimizes search space



P Q G        P Q G
P E G        E Q R
7+2+6 = 15   -1+5+-2 = 2

# Extending Matches

- ## Match = HSP (High-scoring Sequence Pair)
  - ### Match is found and extended as long as score stays above a threshold value
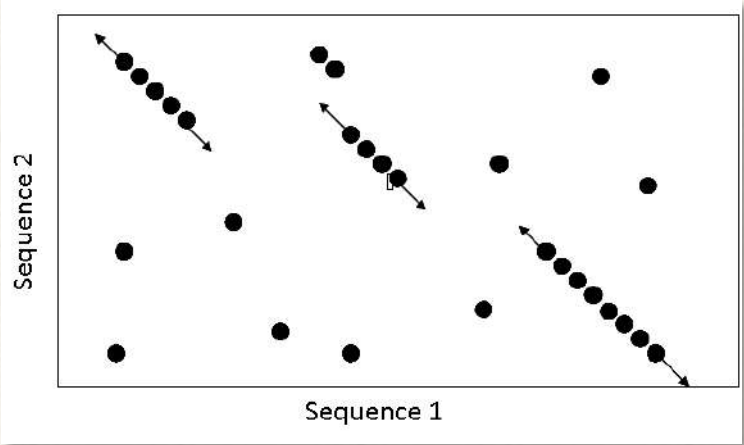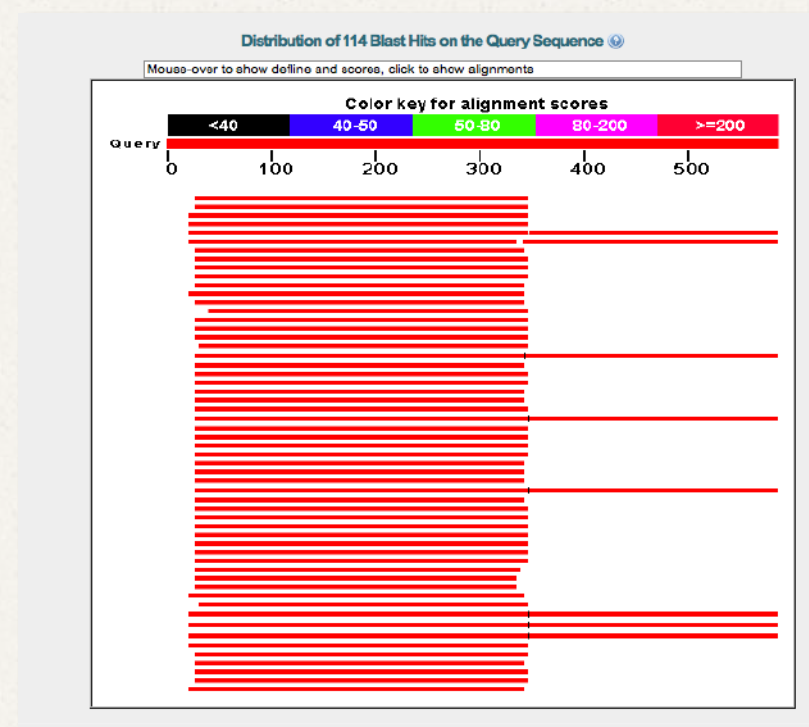    - #### After finished extending, the HSP is kept if above the cutoff score (S)

# Assembling HSPs

- HSPs, after extension, are assembled into a longer alignment
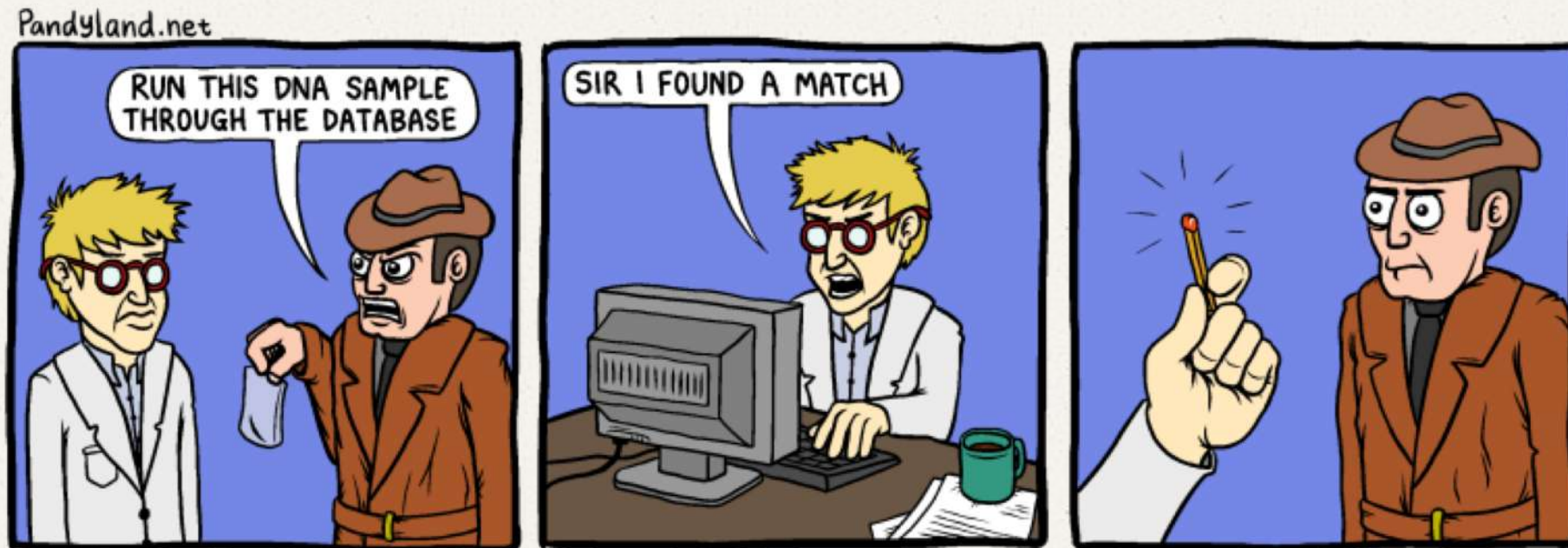
# Output

- <u>Max/Total Score</u>:  quality of the alignment
  - Higher the score the better the match
- <u>Query Coverage</u>: what proportion of the query the particular HSP covers
- <u>E-value</u>:  probability that a match ≥ Max Score occurs by random chance (based on database size)
- <u>Max Identity</u>:  For that HSP, the % of bases that match

| Accession | Total Score | Query Coverage | E-value | Max Ident |
|-----------|-------------|----------------|---------|-----------|
| X56286.1 | 579 | 54% | 7e-162 | 99% |
| AF091629.1 | 573 | 54% | 3e-160 | 99% |
| L48348.1 | 481 | 55% | 2e-132 | 93% |

# Interpretation

- The matches you get are only acceptable matches, not necessarily the optimal match
- Your search is only as good as your database
  - If the optimal match is not in the database, you will not find it.
  - If you have sequences not in the database, SUBMIT THEM!

# Take Home Points

- Blast is a powerful tool for database searching

- Very fast, but at the expense of sensitivity

- Flexible (types, databases)

- Interpret results carefully

- Help make it grow!

# Several examples

- Example 1: SRA Blast (https://www.ncbi.nlm.nih.gov/sra)
  - Query: M55627.1
    - *Coccidiodes immitis* (Valley fever fungus*)* ssuRNA
  - Project: SRX633288
    - Puma 454 transcriptome reads


- Example 2: Blast an assembly (https://blast.ncbi.nlm.nih.gov/)
  - Query:
    - TruSeq Universal Adapter
    - AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT
  - Database: nt
  - Organism: *Cyprinus carpio (taxid:7962)*