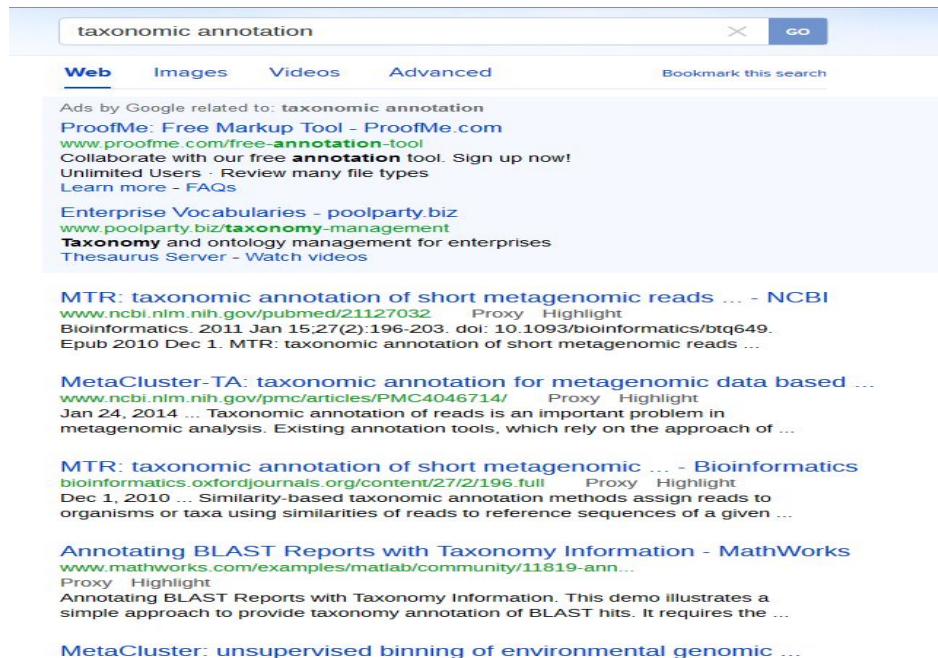# Database searching- Sequence alignment based methods

Nidhi Shah
M$^3$ Workshop
January 10, 2019

# Google: "taxonomic annotation"



- Database of known pages
- Report all that contain keyword

- Ranking important (which of the thousands is most relevant)

# Our "keywords"

>F4BT0V001CZSIM rank=0000138 x=1110.0 y=2700.0 length=5
ACTGCTCTCATGCTGCCTCCCGTAGGAGTGCCTCCCTGAGCCAGGATCAAACGT
>F4BT0V001BBJQS rank=0000155 x=424.0 y=1826.0 length=4
ACTGACTGCATGCTGCCTCCCGTAGGAGTGCCTCCCTGCGCCATCAA
>F4BT0V001EDG35 rank=0000182 x=1676.0 y=2387.0 length=4
ACTGACTGCATGCTGCCTCCCGTAGGAGTCGCCGTCCTCGACNC
>F4BT0V001D2HQQ rank=0000196 x=1551.0 y=1984.0 length=4
ACTGACTGCATGCTGCCTCCCGTAGGAGTGCCGTCCCTCGAC
>F4BT0V001CM392 rank=0000206 x=966.0 y=1240.0 length=82
AANCAGCTCTCATGCTCGCCCTGACTTGGCATGTGTTAAGCCTGTAGGCTAGCGT
>F4BT0V001EIMFX rank=0000250 x=1735.0 y=907.0 length=46
ACTGACTGCATGCTGCCTCCCGTAGGAGTGTCGCGCCATCAGACTG
>F4BT0V001ENDKR rank=0000262 x=1789.0 y=1513.0 length=56
GACACTGTCATGCTGCCTCCCGTAGGAGTGCCTCCCTGAGCCAGGATCAAACTCTG
>F4BT0V001D91MI rank=0000288 x=1637.0 y=2088.0 length=56
ACTGCTCTCATGCTGCCTCCCGTAGGAGTGCCTCCCTGAGCCAGGATCAAACTCTG
>F4BT0V001D0Y5G rank=0000341 x=1534.0 y=866.0 length=75
GTCTGTGACATGCTGCCTCCCGTAGGAGTCTACACAAGTTGTGGCCCAGAACCACTGAGCCAGGATCAAACTCTG
>F4BT0V001EMLE1 rank=0000365 x=1780.0 y=1883.0 length=84
ACTGACTGCATGCTGCCTCCCGTAGGAGTGCCTCCCTGCGCCATCAATGCTGCATGCTGCTCCCTGAGCCAGGATCAAACTCTG

# Taxonomic annotation algorithm

- Look through database for sequence

- Report all organisms that contain it

- Rank list by ?? (most relevant to least relevant)

- Handle sequencing errors

- What if the sequence is not in the database (Should handle evolutionary divergence)

  - Can we say anything about the data?

    e.g., google "taxnomic anntoation"

# Taxonomic annotation algorithm

Solution:

- Organize the database (taxonomy)

  Kingdom;Phylum;Class;Order;Family;Genus;Species;Strain

- Use search procedure that can generalize from existing knowledge - ??

# Taxonomic annotation algorithm

Solution:

- Organize the database (taxonomy)

  Kingdom;Phylum;Class;Order;Family;Genus;Species;Strain

- Use search procedure that can generalize from existing knowledge

  - Sequence similarity search

  - Assume taxonomically related sequences are similar

# Similarity search

Query
DB

```
ACCATAG-GCCGTCAGACCTAGACTAGA
AC-ATAGAGCCGTCAGA-CTATACTAGA
```

- Finds exact matches
- Handles sequencing errors
- May handle evolutionary divergence
- May provide statistical guarantees (is this a random hit?) -- can help with ranking results!

- MANY tools exist for doing the search! (e.g. BLAST, Diamond, Usearch, BLAT, etc)
- Differ by
  - assumptions about data
  - similarity cutoff
  - heuristics to speed up search (incl. memory/speed trade-off)

# Basic Local Alignment Search Tool

**BLAST** finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance.          Learn more

## Web BLAST

**Nucleotide BLAST**
nucleotide ▶ nucleotide

**blastx**
translated nucleotide ▶ protein

**tblastn**
protein ▶ translated nucleotide

**Protein BLAST**
protein ▶ protein

### BLAST Genomes

Enter organism common name, scientific name, or tax id          **Search**

Human          Mouse          Rat          Microbes

# BLAST web result example

RIH U.S. National Library of Medicine    NCBI National Center for Biotechnology Information    nidhirshah1992@gmail.com   My NCBI   Sign Out

BLAST® » blastn suite » RID-3DKBXPYV014    Home   Recent Results   Saved Strategies   Help

BLAST Results

ⓘ Your search is limited to records that exclude: models (XM/XP), uncultured/environmental sample sequences »Full Entrez Query

Edit and Resubmit   Save Search Strategies   »Formatting options   »Download    How to read this page   Blast report description

Job title: OTU_97.10029 (511 letters)

| | |
|---|---|
| **RID** 3DKBXPYV014 (Expires on 01-11 11:43 am) | **Database Name** nr |
| **Query ID** lcl|Query_240057 | **Description** Nucleotide collection (nt) |
| **Description** OTU_97.10029 | **Program** BLASTN 2.8.1+ »Citation |
| **Molecule type** nucleic acid | |
| **Query Length** 511 | |

Other reports: »Search Summary [Taxonomy reports] [Distance tree of results] [MSA viewer]

⊟ Graphic Summary

Distribution of the top 120 Blast Hits on 100 subject sequences ⓘ

Mouse over to see the title, click to show alignments

Color key for alignment scores
■ <40   ■ 40-50   ■ 50-80   ■ 80-200   ■ >=200

Query
1    100    200    300    400    500

⊟ Descriptions

Sequences producing significant alignments:

Select: All None    Selected:0

⊞ Alignments  ⬇Download ⌄  GenBank  Graphics  Distance tree of results

| Description | Max score | Total score | Query cover | E value | Ident | Accession |
|---|---|---|---|---|---|---|
| Faecalibacterium prausnitzii strain APC918/95b chromosome, complete genome | 859 | 5036 | 95% | 0.0 | 99% | CP030777.1 |
| Faecalibacterium prausnitzii isolate S9G3 16S ribosomal RNA gene, partial sequence | 854 | 854 | 95% | 0.0 | 98% | MF186592.1 |
| Faecalibacterium prausnitzii strain S3L/3 16S ribosomal RNA gene, partial sequence | 845 | 845 | 95% | 0.0 | 98% | HQ457024.1 |
| Faecalibacterium prausnitzii SL3/3 draft genome | 841 | 841 | 95% | 0.0 | 98% | FP929046.1 |
| Butyrate-producing bacterium M21/2 16S ribosomal RNA gene, partial sequence | 839 | 839 | 90% | 0.0 | 99% | AY305307.1 |
| Faecalibacterium prausnitzii isolate S9D8 16S ribosomal RNA gene, partial sequence | 837 | 837 | 95% | 0.0 | 98% | MF186232.1 |
| Faecalibacterium prausnitzii strain CNCM_I_4541 16S ribosomal RNA gene, partial sequence | 837 | 837 | 95% | 0.0 | 98% | MF185659.1 |

# BLAST top hit

5467_464          HM038000.1.1446   E-value: 6e-96        Bit score: 350
Bacteria;Cyanobacteria;Melainabacteria;Vampirovibrionales;Vampirovibrio chlorellavorus

Bit score - the "information" contained in the alignment

E-value - how many random alignments one expects for the same bit score

The lower the E-value, the more significant the alignment score of the sequence match

# Run BLAST on command line

# BLAST top hit

5467_464          HM038000.1.1446  E-value: 6e-96        Bit score: 350
Bacteria;Cyanobacteria;Melainabacteria;Vampirovibrionales;Vampirovibrio chlorellavorus

Bit score - the "information" contained in the alignment
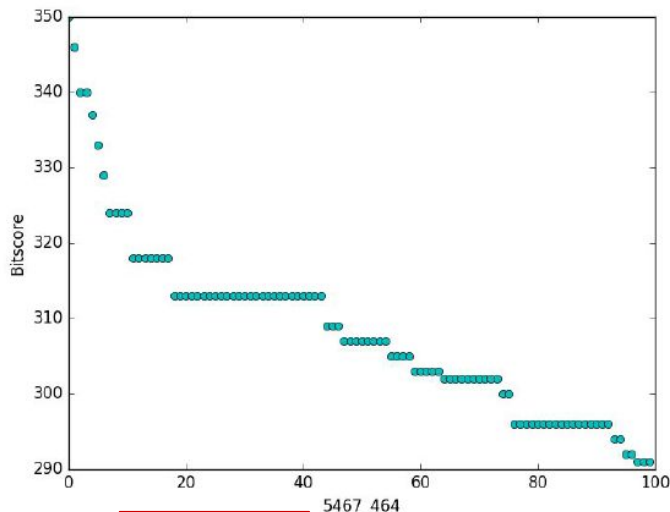
E-value - how many random alignments one expects for the same bit score

The lower the E-value, the more significant the alignment score of the sequence match

# BLAST...more hits

5467_464          HM038000.1.1446 Identity: 80.00%  E-value: 6e-96 Bit score: 350



top 100 hits sorted by bit score

Bacteria;Cyanobacteria;Melainabacteria;Vampirovibrionales;Vampirovibrio chlorellavorus

Bacteria;Proteobacteria;Alphaproteobacteria;Caulobacterales;Caulobacteraceae;Brevundimonas;Brevundimonas mediterranea

Bacteria;Proteobacteria;Alphaproteobacteria;Caulobacterales;Caulobacteraceae;Brevundimonas; Brevundimonas bacteroides

Bacteria;Firmicutes;Clostridia;Clostridiales;Ruminococcaceae;Butyricicoccus;Butyricicoccus pullicaecorum

# How to interpret these results?

- Why does this happen??

 

- Should we use the best hit and transfer annotation to the query sequence?

- Cutoffs on percent identity? Bit score? E-value?

- For e.g. Megan, PhymBl

- What about query coverage?

Algorithms for
Molecular Biology

**RESEARCH**

**Open Access**

CrossMark

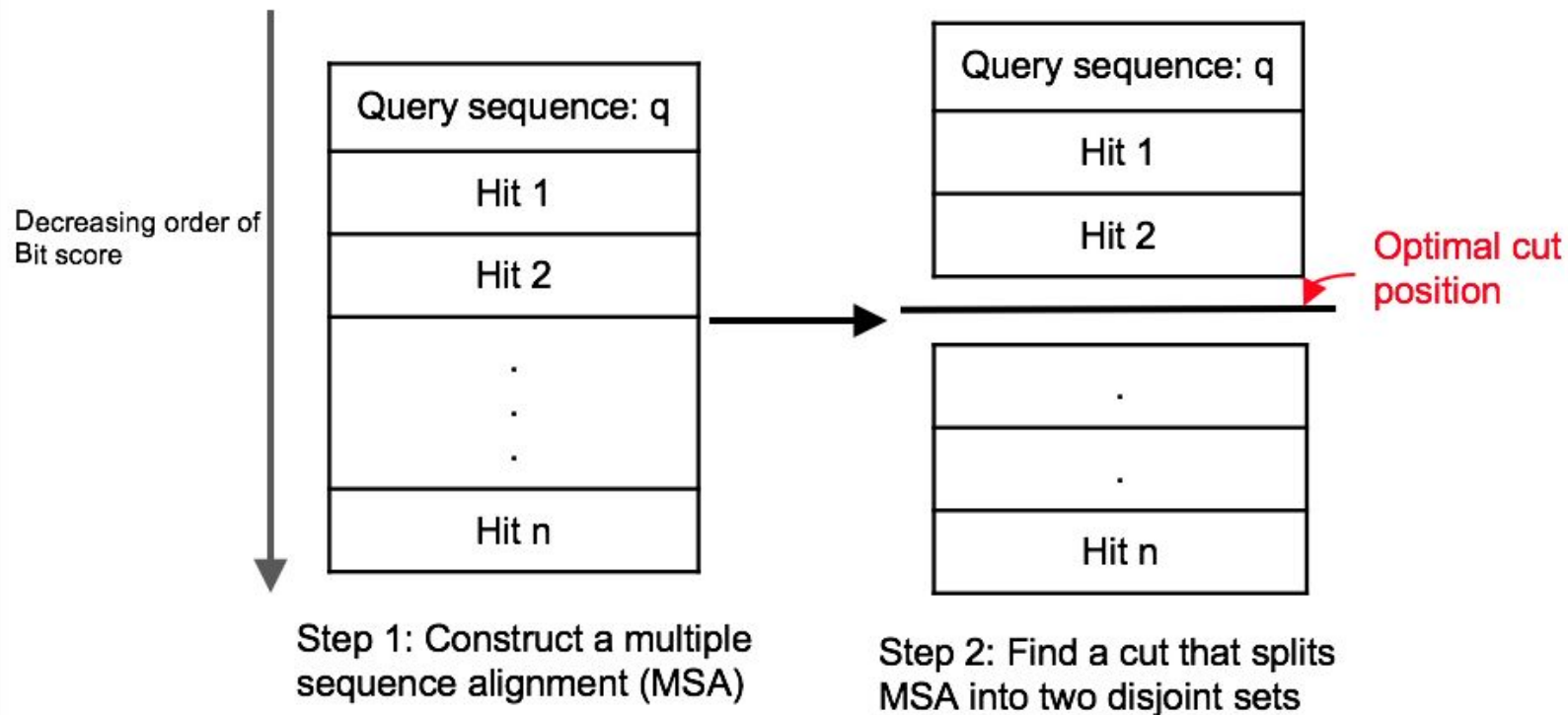# Outlier detection in BLAST hits

Nidhi Shah[1], Stephen F. Altschul[2] and Mihai Pop[1*]

**Abstract**

**Background:** An important task in a metagenomic analysis is the assignment of taxonomic labels to sequences in a sample. Most widely used methods for taxonomy assignment compare a sequence in the sample to a database of known sequences. Many approaches use the best BLAST hit(s) to assign the taxonomic label. However, it is known that the best BLAST hit may not always correspond to the best taxonomic match. An alternative approach involves phylogenetic methods, which take into account alignments and a model of evolution in order to more accurately define the taxonomic origin of sequences. Similarity-search based methods typically run faster than phylogenetic methods and work well when the organisms in the sample are well represented in the database. In contrast, phylogenetic methods have the capability to identify new organisms in a sample but are computationally quite expensive.

Github: https://github.com/shahnidhi/outlier_in_BLAST_hits
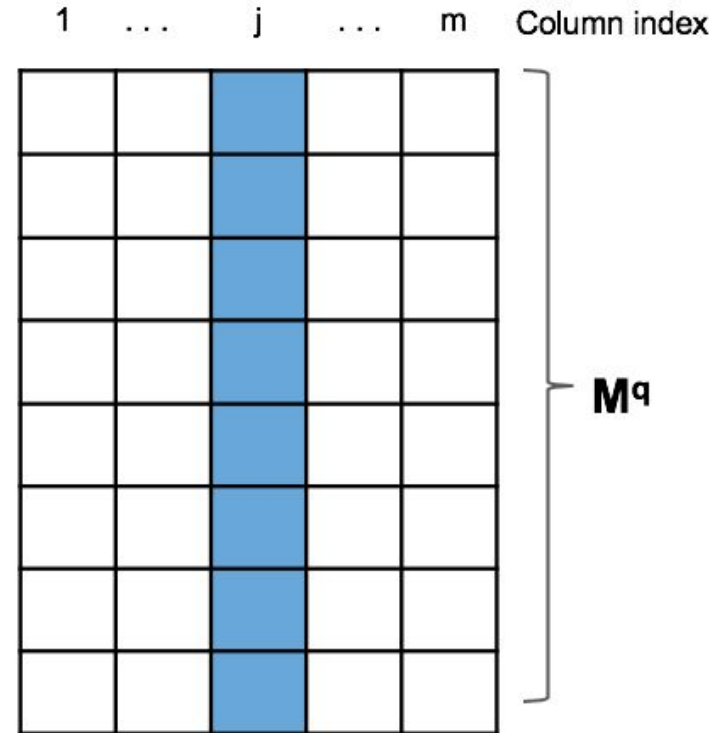
# Methods

# Methods - Score of a Multiple Sequence Alignment

- Bayesian Integral Log Odds (BILD) score

- Dirichlet distribution for nucleotide probabilities prior

$$L(M_j^q) = \log \left[ \frac{\Gamma(\alpha^*)}{\Gamma(\alpha^* + c_j^*)} \prod_{k=1}^{4} \frac{\Gamma(\alpha_k + c_{jk})}{\Gamma(\alpha_k)} \right]$$

- Jeffrey's prior: $\alpha^* = 2$, $\alpha^k = 0.5$

- $C^{ik}$ = total count of base k (A, C, T, G)

- $C^{i*} = C^{iA} + C^{iC} + C^{iT} + C^{iG}$

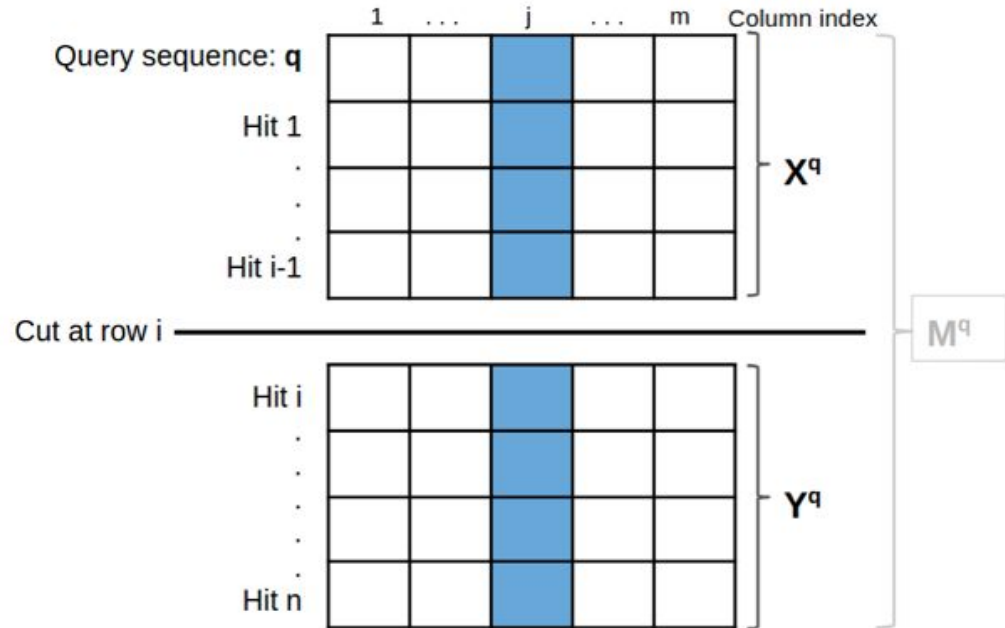| 1 | . . . | j | . . . | m | Column index |

M�q

# Methods - Scoring a split in an MSA

Score of a cut at row i : $V_i^q$

$$V_{ji}^q = L(X_{ji}^q) + L(Y_{ji}^q) - L(M_j^q)$$

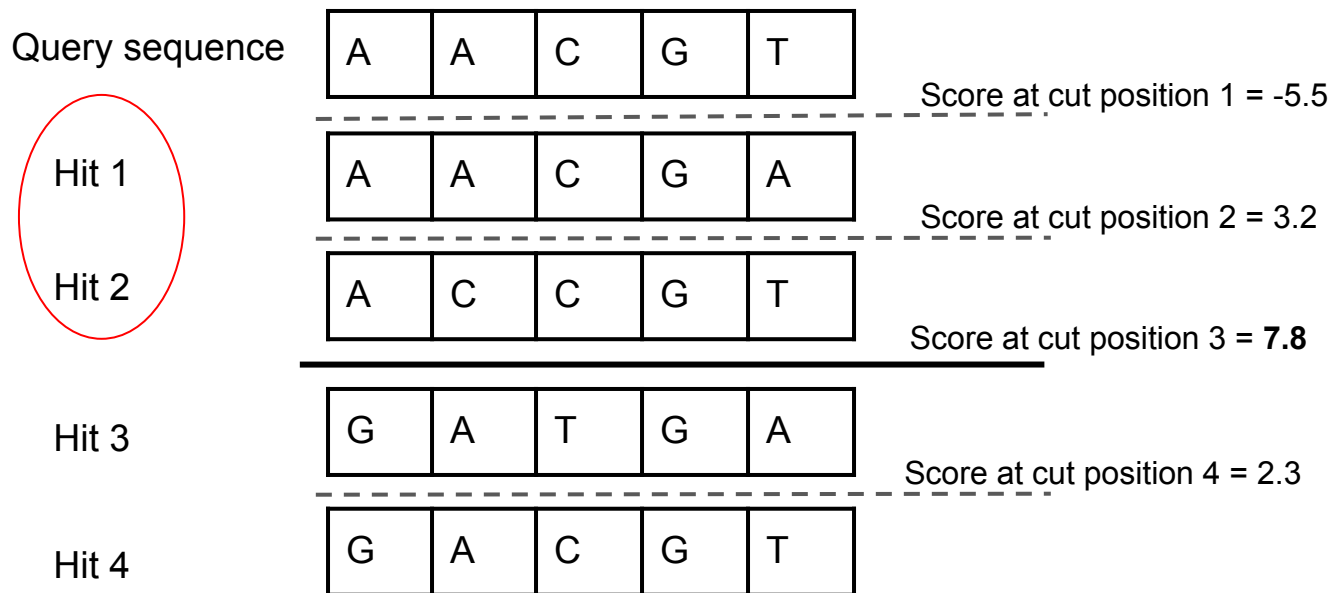$$V_i^q = \sum_{j=1}^{m} e_j^a V_{ji}^q$$

- $e_j$ is the entropy of column j

# Methods - Detecting outliers

| Query sequence | A | A | C | G | T |
|---|---|---|---|---|---|

Score at cut position 1 = -5.5

| Hit 1 | A | A | C | G | A |
|---|---|---|---|---|---|

Score at cut position 2 = 3.2

| Hit 2 | A | C | C | G | T |
|---|---|---|---|---|---|

Score at cut position 3 = **7.8**

| Hit 3 | G | A | T | G | A |
|---|---|---|---|---|---|

Score at cut position 4 = 2.3

| Hit 4 | G | A | C | G | T |
|---|---|---|---|---|---|

Note: We are looking for the first positive peak in the score

# Methods - Detecting outliers

| Query sequence | A | A | C | G | T |

Score at cut position 1 = -5.5

| Hit 1 | A | A | C | G | A |

Score at cut position 2 = 3.2

| Hit 2 | A | C | C | G | T |

Score at cut position 3 = **7.8**

| Hit 3 | G | A | T | G | A |

Score at cut position 4 = 2.3

| Hit 4 | G | A | C | G | T |

Note: We are looking for the first positive peak in the score

# Run outlier detection pipeline on command line

# How to get taxonomic annotation?

- Use Most recent common ancestor of candidate DB sequences (outliers)
- Look at output files
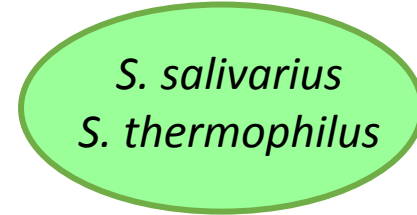

- Can we do something better?

# Taxonomic resolution

| Family | Genus | Species |
| --- | --- | --- |
| Streptococcaceae | Streptococcus | salivarius |
| Streptococcaceae | Streptococcus | |
| Streptococcaceae | Streptococcus | salivarius |
| Streptococcaceae | Streptococcus | |
| Streptococcaceae | Streptococcus | thermophilus |

↑ Lowest Common Ancestor

**Non-Pathogenic**

*S. salivarius*
*S. thermophilus*

**Pathogenic**

*S. mitis*
*S. pyogenes*

# Creating taxonomic-agnostic clusters



Query
sequences

Database
sequences

"Confusion graph"
on database
sequences

Graph partitioning

Database
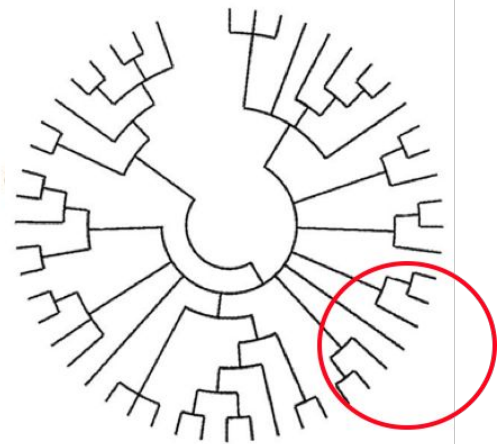partitions/clusters
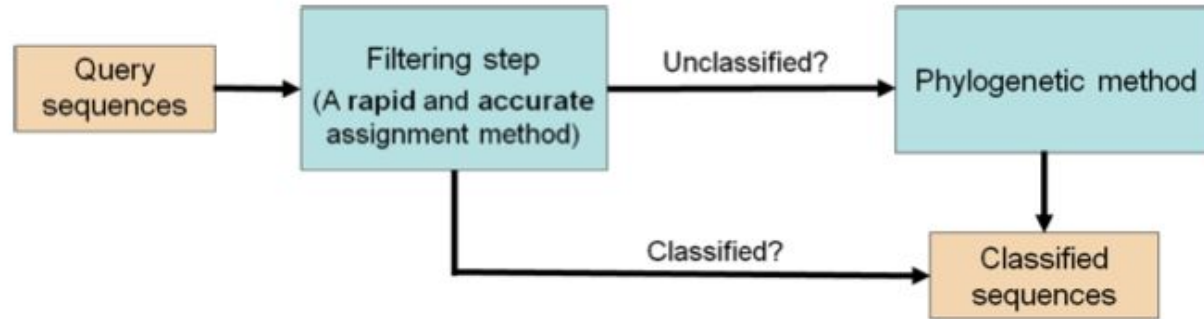
# Look at partition files

# Why certain query sequences have no outliers?

# Why certain query sequences have no outliers?

- Query sequence is very different from everything in the database

- The algorithm cannot find any cut position with positive score (peak)

Thanks!

# A two-step approach

# Suggestions for WMS

Assembly

Variant detection

Look at assembly graphs

Read level classification

Abundance profiling - marker gene based methods

# Landscape of Taxonomic Approaches



runtime
match to "truth"

**k-mer counting**
kraken

**k-mer stat model**
RDP

**Sequence Alignment**
blast
usearch
diamond

**Phylogenetic**
pplacer
TIPP

Assume taxonomically related sequences have similar k-mer frequencies

Assume taxonomically related sequences are similar

Simulate evolution