# Tutorial 10:

# Run Diamond Docker on the datasets for Kraken2

# Background

# What is DIAMOND?

- A sequence aligner for protein and translated DNA searches

- Designed for high performance analysis of big sequence data

- Has an alignment sensitivity that matches BLAST

- Features various sensitivity modes:
  - Default (fast)
  - Sensitive
  - Very-sensitive
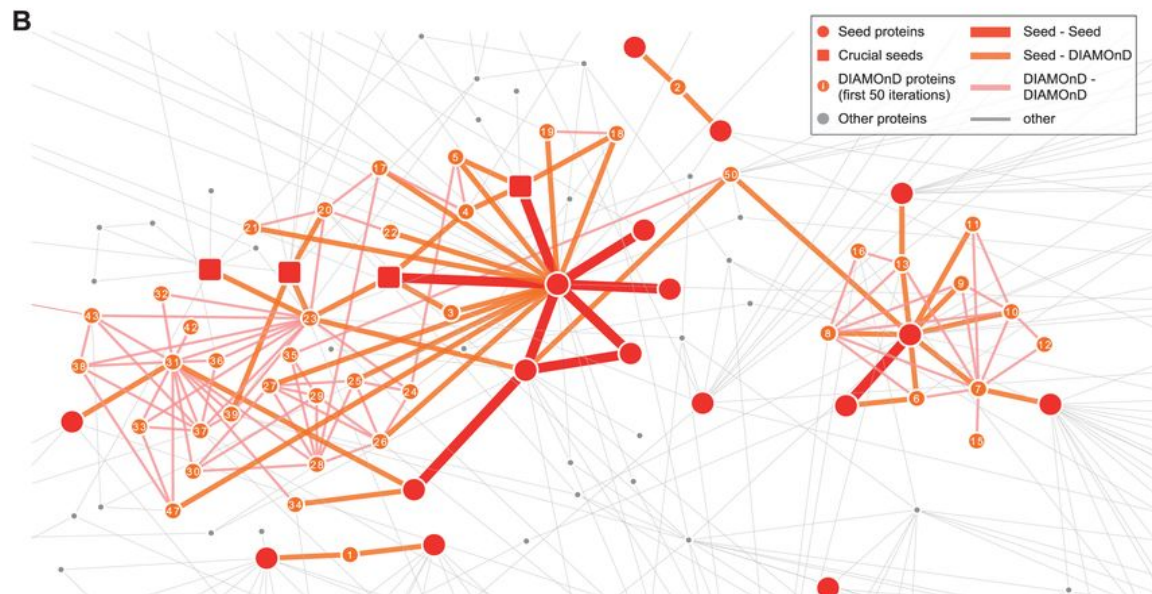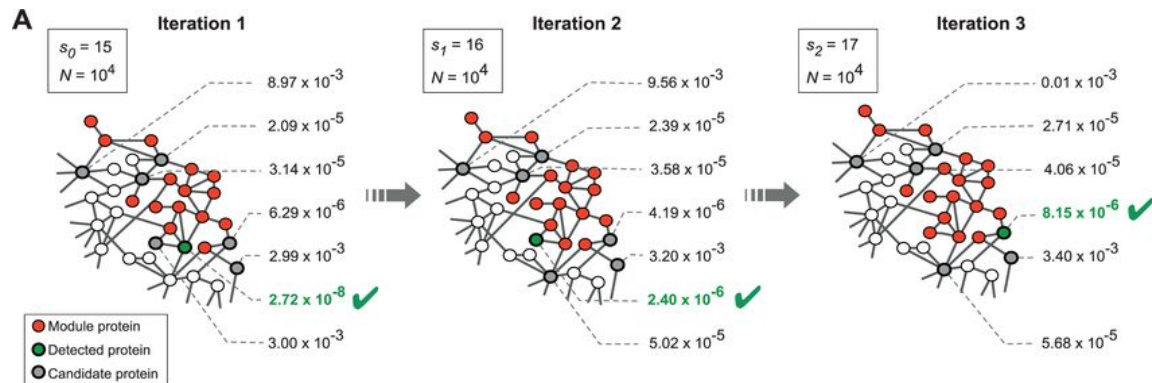  - Ultra-sensitive

# How Does DIAMOND Work?

Double Indexing Approach:

- Database of seed-location pairs are built for query and reference sequences.
  - First index organizes and stores information about the sequences.
  - Second index is for querying sequences provided by the user.
- Matching seeds are paired through a hash join technique.

Seed-and-Extend Algorithm:

- Identifies exact matches between sequences from the indexed database and the query index.
- Once seeds are identified, they are mapped to find longer, high-scoring alignments.

# Seed-and-Extend Algorithm

# Comparison

# Comparison of DIAMOND and Kraken2

**DIAMOND**

- Uses a protein database and can query both DNA or amino acid sequences.
- Query algorithm: alignment using spaced seeds
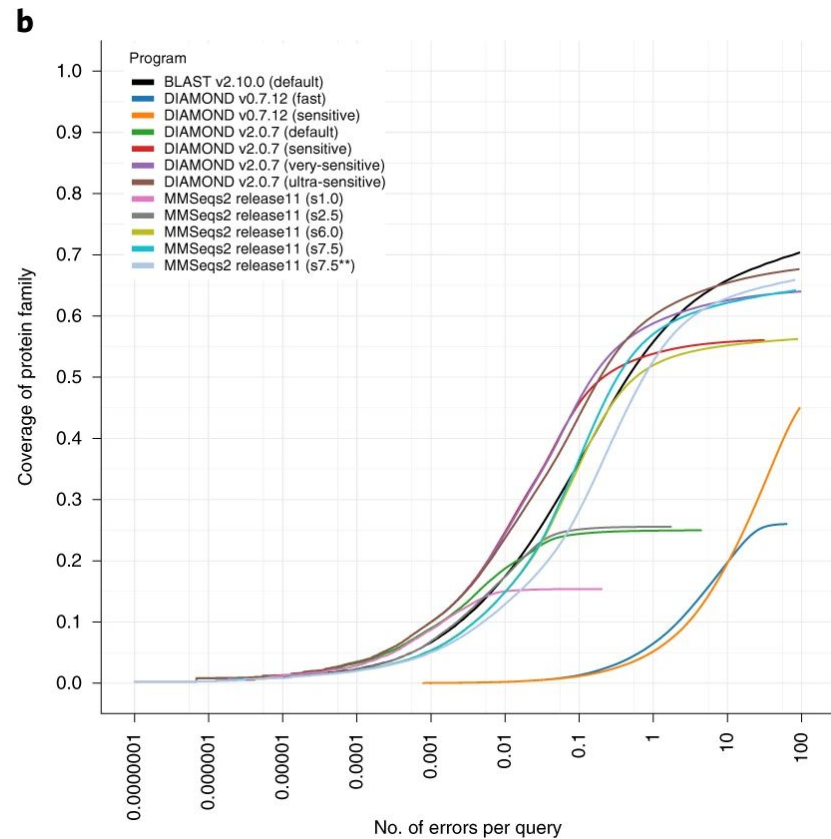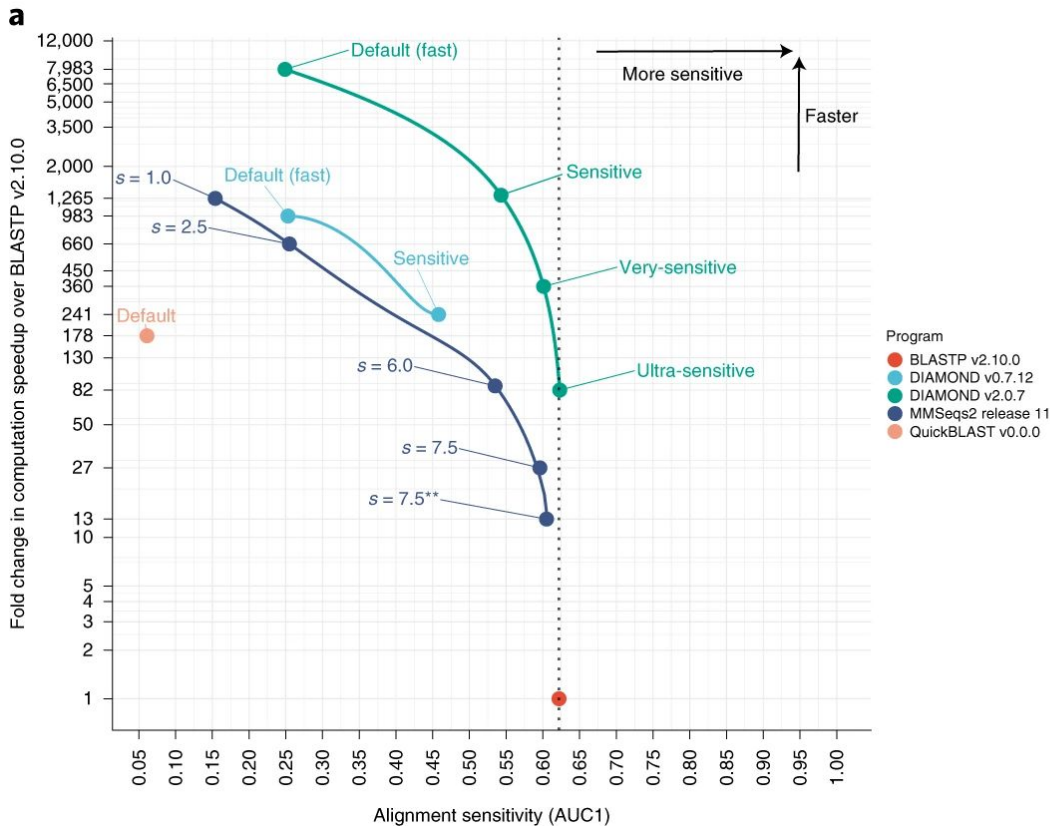- Taxonomic binning strategy: Lowest common ancestor

**Kraken2**

- Uses protein or DNA, whole genome or single locus databases
- Query algorithm: exact k-mer matching
- Taxonomic binning strategy utilizes the highest number of k-mer matches considering a root-to-leaf path

# Comparison to Other Tools

- DIAMOND's default mode is reported to be up to 20,000 times faster than BLASTX and reports about 80-90% of the matches that BLASTX would find.

- In sensitive mode, DIAMOND is about 2,500 times faster than BLASTX, finding more than 94% of all matches.

# Comparison to Other Tools (Cont.)

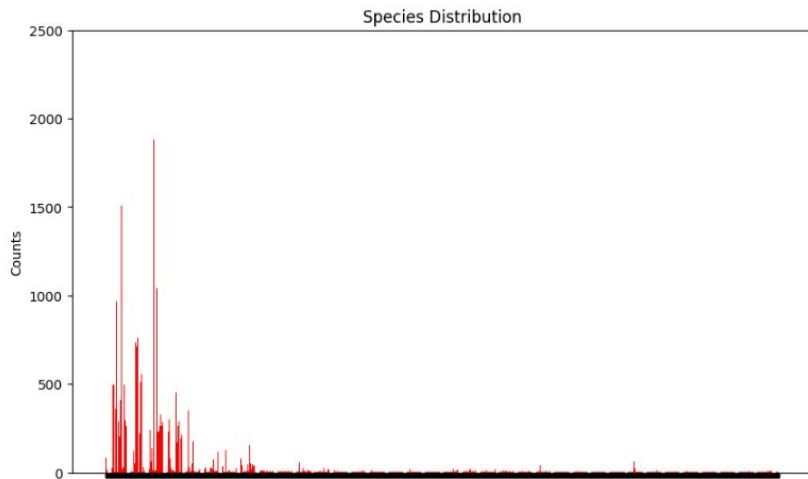# Usage and Experiments

# Usage

- singularity exec --bind /path/to/your/data:/data diamond_latest.sif /bin/bash
- diamond makedb --in /data/combined_uniprot_sprot_2_10.fasta.gz -d /data/sprotdb_2_10
- diamond blastp --ultra-sensitive -d sprotdb_2_10 -q /data/uniprot_sprot_part1.fasta.gz -o /data/output_file.m8 --threads 96

# Dataset

- gene_refseq_uniprotkb_collab.gz
  - 176,513,537 accessions(ids)
  - no sequence
- uniprot_trembl.fasta.gz
  - 249,751,891 sequences
  - predicted data
- **uniprot_sprot.fasta.gz**
  - 570,830 sequences
  - curated data
- ETC
  - evol1.sorted.unmapped.R2.fastq.gz (17,692 sequences)
  - astral-scopedom-seqres-gd-all-2.08-2923-901-06.fa (1/11th, 27,776 sequences)

# Dataset Analysis

- **uniprot_sprot.fasta.gz**
  - 570,830 sequences
  - curated data
  - 12,014 species
  - one sample species: 4,888
  - less than 10 samples species: 9,606



Species Distribution

# Run Diamond with sensitive settings

- default
- --sensitive
- --more-sensitive
- --ultra-sensitive

|  | time(s) | pairwise aligned(#) | Aligned queries(#) |
|---|---|---|---|
| default | 12.555 | 4077 | 403 |
| sensitive | 65.239 | 5828 | 405 |
| more | 65.613 | 5828 | 405 |
| ultra | 64.898 | 5828 | 405 |

Query: evol1.sorted.unmapped.R2.fastq.gz
DB: uniprot_sprot

# Run Diamond with sensitive settings

- default
- --sensitive
- --more-sensitive
- --ultra-sensitive

|  | time(s) | pairwise aligned(#) | Aligned queries(#) |
|---|---|---|---|
| default | 4.226 | 5275 | 211 |
| sensitive | 22.264 | 5369 | 217 |
| more | 22.059 | 5369 | 217 |
| ultra | 69.202 | 5405 | 226 |

Query: 1/11th astral-scopedom-seqres-gd-all-2.08-2923-901-06.fa
DB: another 1/11th astral-scopedom-seqres-gd-all-2.08-2923-901-06.fa

# Run Diamond with sensitive settings

- default
- --sensitive
- --more-sensitive
- --ultra-sensitive

|  | time(s) | pairwise aligned(#) | Aligned queries(#) |
|---|---|---|---|
| default | 33.283 | 429380 | 30195 |
| sensitive | 152.994 | 625449 | 36799 |
| more | 189.506 | 638549 | 37119 |
| ultra | 907.452 | 676317 | 37963 |

Query: 1/10th uniprot_sprot
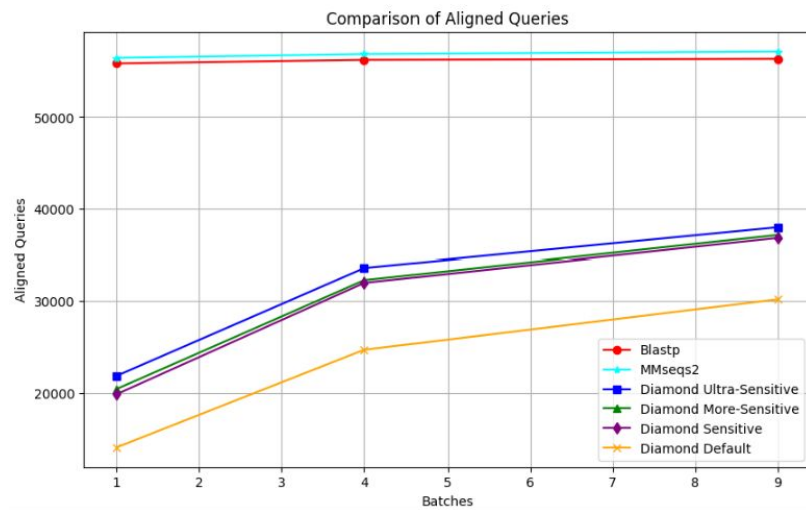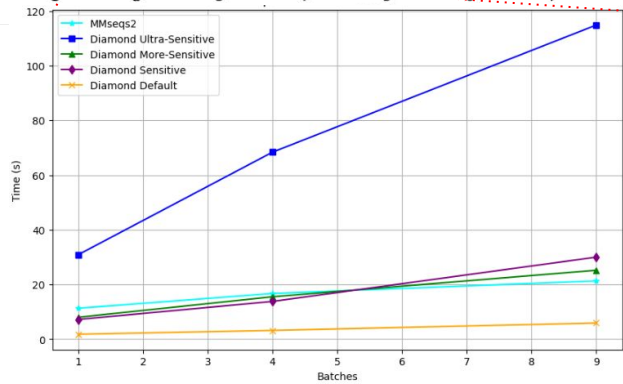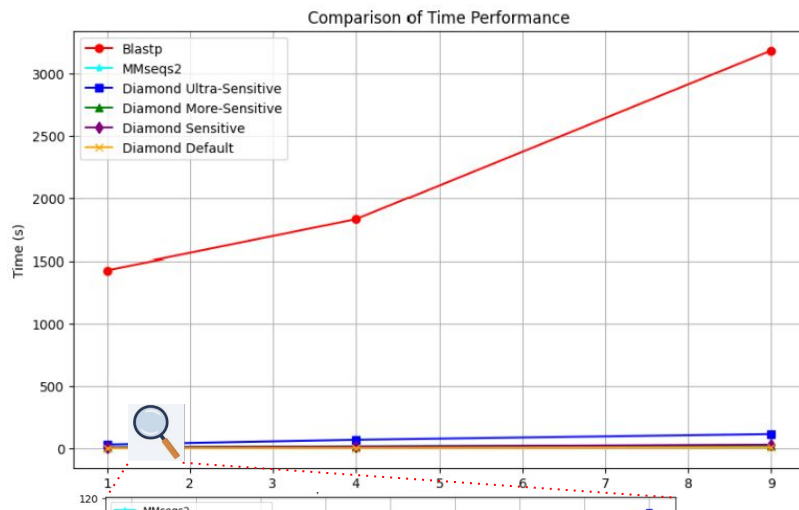DB: another 9/10th uniprot_sprot

# Comparison with Blast

● Use 96 thread for blastp

| | time(s) | pairwise aligned(#) | Aligned queries(#) |
|---|---|---|---|
| blastp | 3181.1 | 9611230 | 56366 |
| default | 5.9 | 429380 | 30195 |
| sensitive | 30.0 | 625531 | 36799 |
| more | 25.2 | 638673 | 37123 |
| ultra | 114.9 | 676055 | 37962 |

Query: 1/10th uniprot_sprot
DB: another 9/10th uniprot_sprot

# Incremental Comparison

# Comparison with Kraken2

- Get texa for uniqrot_sprot file using ncbi api
  - Only able to get about 10%(57,083 seqs) of uniprot_sprot(about 8 hours to take)
- Build Kraken db with 9/10th of the file with texa
- Test 1/10th of the file(5,708 seqs) as a query
- Got unclassified result

# Resources

- About DIAMOND:

  Buchfink, B., Reuter, K., & Drost, H.-G. (2021, April 7). *Sensitive protein alignments at tree-of-life scale using Diamond*. Nature News. https://www.nature.com/articles/s41592-021-01101-x#Sec2

- DIAMOND Algorithm:

  Ghiassian, S. D., Menche, J., & Barabási, A.-L. (n.d.). *A disease module detection (diamond) algorithm derived from a systematic analysis of connectivity patterns of disease proteins in the human interactome*. PLOS Computational Biology. https://journals.plos.org/ploscompbiol/article?id=10.1371%2Fjournal.pcbi.1004120

- Comparison for DIAMOND and other tools:

  McIntyre, A.B.R., Ounit, R., Afshinnekoo, E. *et al.* Comprehensive benchmarking and ensemble approaches for metagenomic classifiers. *Genome Biol* 18, 182 (2017). https://doi.org/10.1186/s13059-017-1299-7

# Resources

- About source code and libraries:

  DerrickWood's Kraken 2 on GitHub: https://github.com/DerrickWood/kraken2

  Official NCBI BLAST+ Docker Image Documentation on GitHub: https://github.com/ncbi/blast_plus_docs

  soedinglab MMseqs2: Ultra-fast and sensitive sequence search and clustering suite on GitHub: https://github.com/soedinglab/MMseqs2

  DIAMOND: A sequence aligner for protein and translated DNA searches, designed for high performance analysis of big sequence data on GitHub by bbuchfink: https://github.com/bbuchfink/diamond