# TUTORIAL 8  KRAKEN

BY CODY SCULLY & CHRISTOPHER UZOKWE

# TAXONOMIC CLASSIFICATION
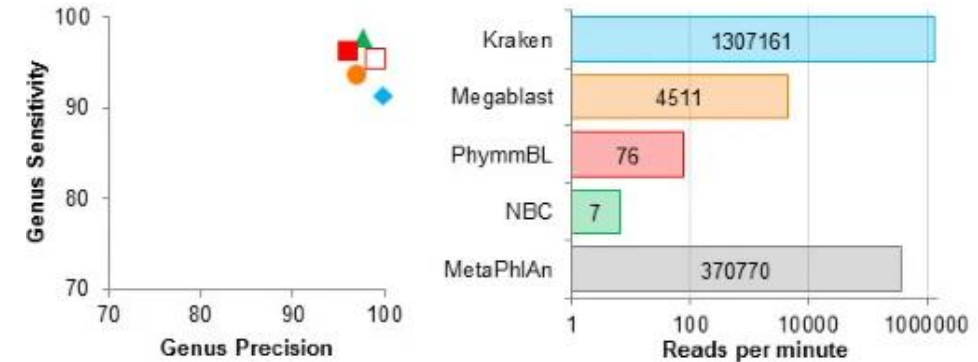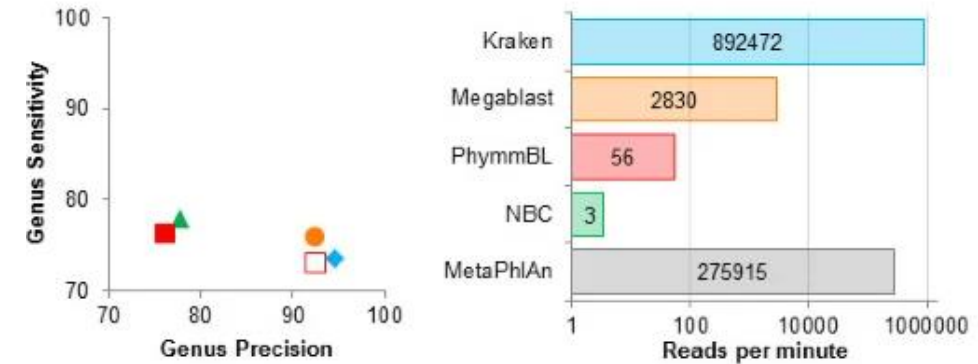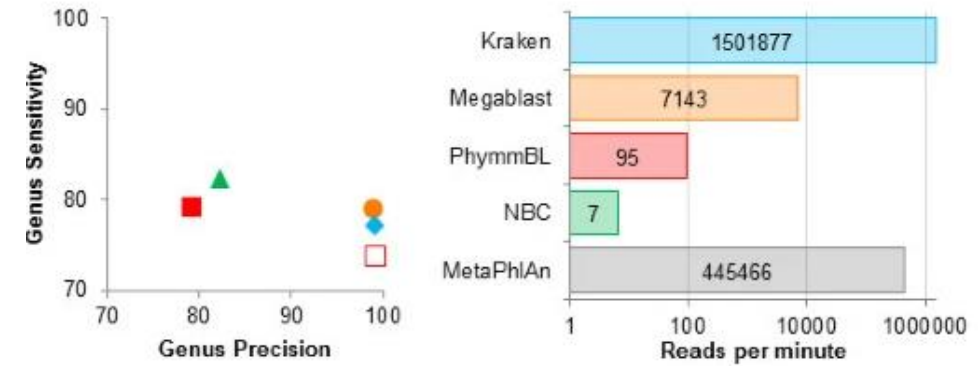
Blast (Magablast)

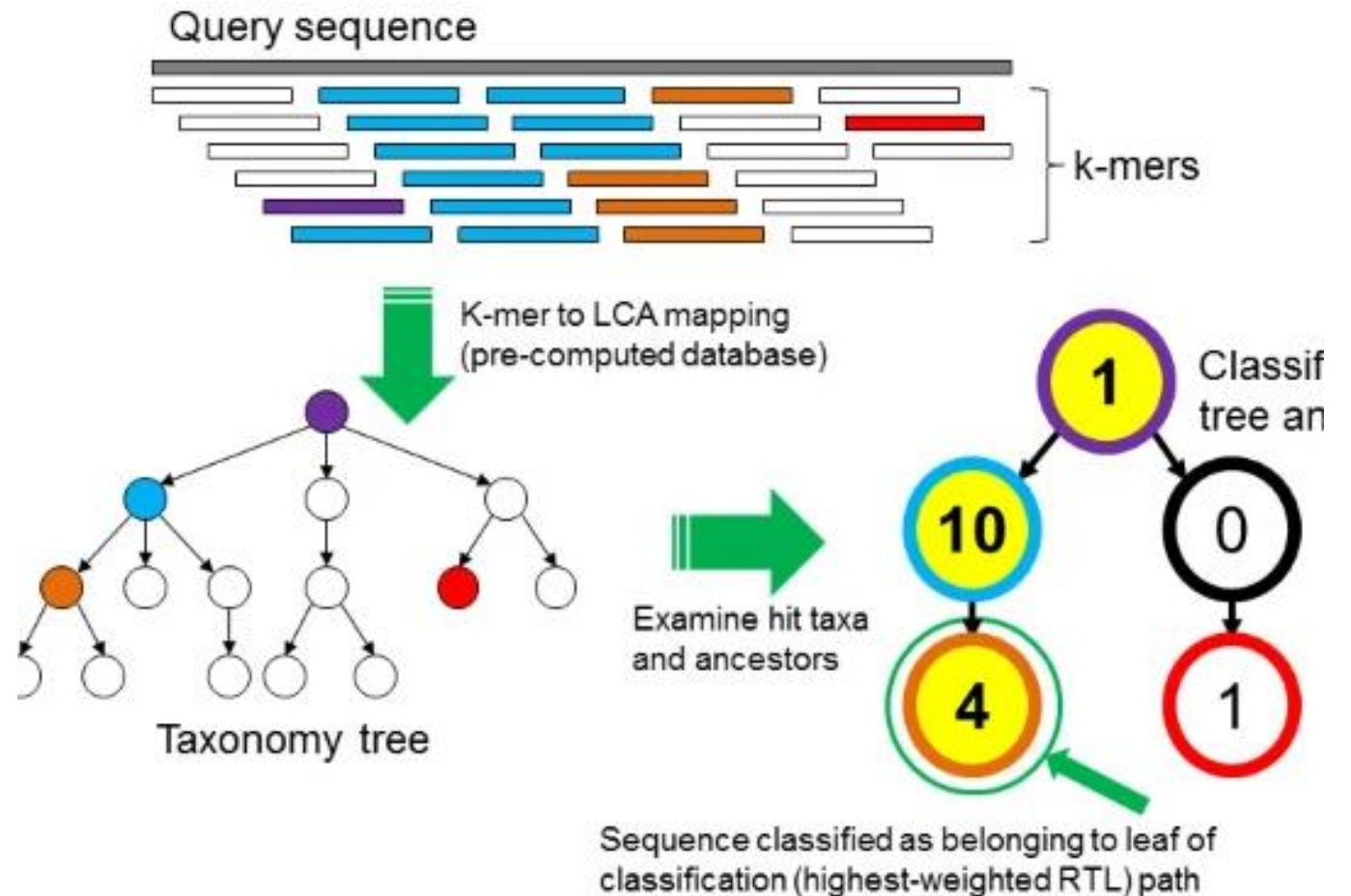The Naïve Bayes Classifier (NBC)

PhymmBL

MatePhlAn

Kraken

# OVERVIEW OF KRAKEN

- This is a Program that reads in Group of unmapped data

- Classifies the data using k-mers to map the LCA

- Creates a Taxonmy tree with the Classifications

- Removes any unclassified leaves from the tree



Query sequence

k-mers

K-mer to LCA mapping (pre-computed database)

Taxonomy tree

Examine hit taxa and ancestors

Classif tree an

Sequence classified as belonging to leaf of classification (highest-weighted RTL) path

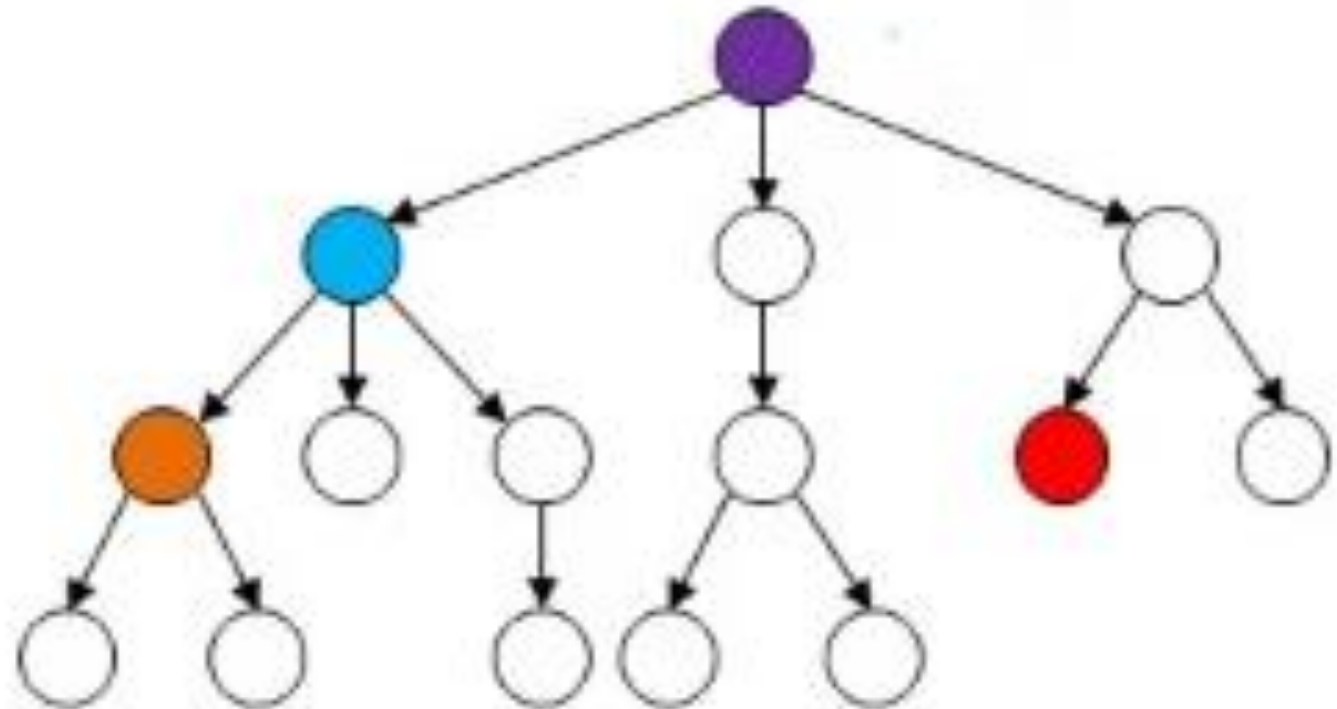# OVERVIEW OF METAGENOMICS WORKFLOW

# QUERY SEQUENCE

- The First step that Kraken takes is to sequence the data set
  - This is done by a library of metagenome data
- Every k-Mer is either
  - **Classified** with the data from the library
  - **Unclassified** if nothing is found for them.
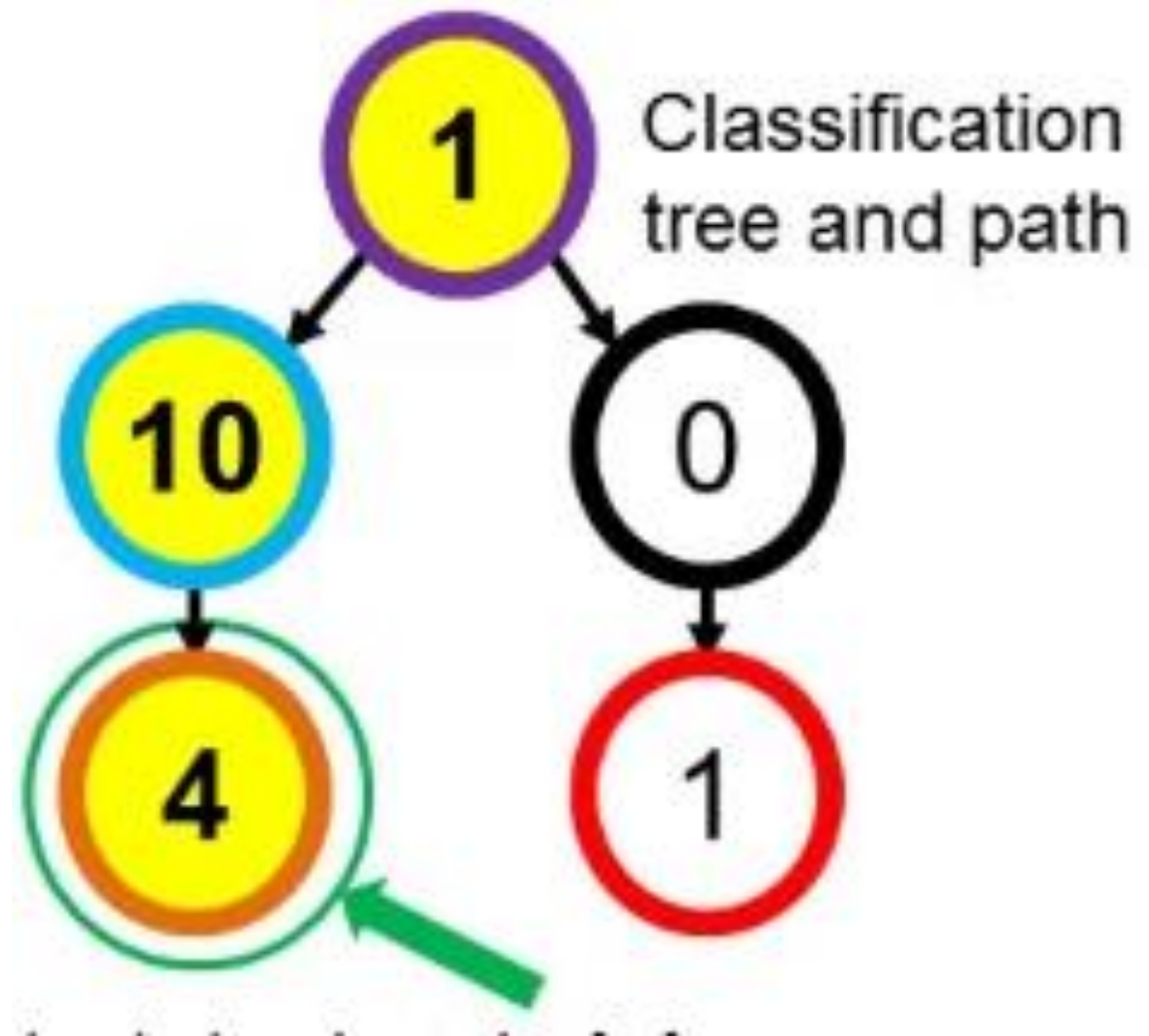


Query sequence

# TAXONOMY TREE

- Each K-mer is mapped to the lowest common ancestor (LCA) of the genomes that contains that k-mer
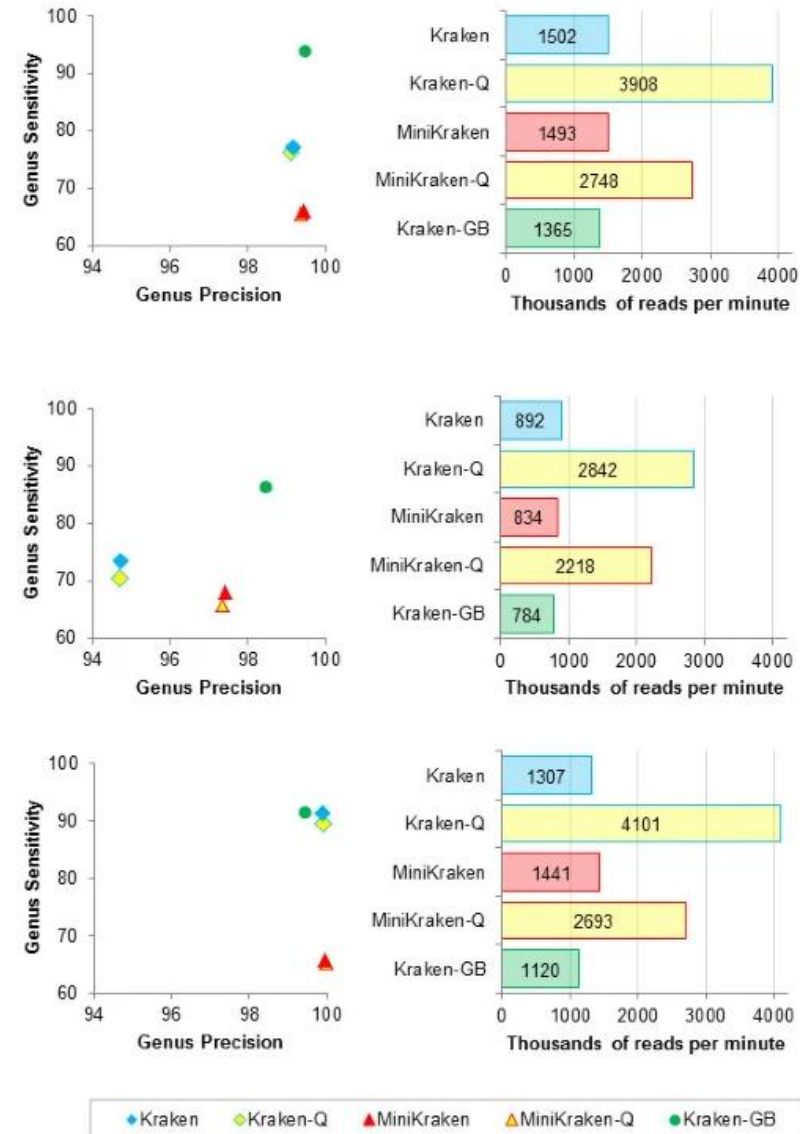


Taxonomy tree

# CLASSIFICATION TREE AND PATH

- Finally the unclassified K-mers are removed and the tree is reduced and weighted/scored



Classification tree and path

# KRAKEN VARIANCE

- **Kraken**
  - Needs 70 GB of RAM

- **Kraken-Q/MiniKraken-Q**
  - Faster then Kraken /MiniKraken
  - Decreased Accuracy

- **MiniKraken**
  - Reduced database
    - Needs 4 GB of RAM
  - Less overall Sensitivity (~11-25%)

- **Kraken-GB**
  - Uses GenBank Database
    - Very high sensitivity

# KRAKEN2 TUTORIAL

# INSTALLATION & DOWNLOADS

- Download Kraken2 Docker:

```
$ conda create --yes -n kraken kraken2 bracken
$ conda activate kraken
```

- (Or) find it on Picotte:

```
biobakery_workflows.sif   edirect_latest.sif    metaspades_latest.sif   nextflow                  qiime2_latest.sif
conda-qiime2_latest.sif   kraken2_latest.sif    miniconda3              nfcore-magbusco-1.2.0.img  README
diamond_latest.sif        metabat_latest.sif    modulefiles             qiime
[cnu25@picotte001 containers]$
```

- Must also download taxo database:

```
$ curl -O ftp://ftp.ccb.jhu.edu/pub/data/kraken2_dbs/minikraken2_v2_8GB_201904_UPDATE.tgz

# alternatively we can use wget
$ wget ftp://ftp.ccb.jhu.edu/pub/data/kraken2_dbs/minikraken2_v2_8GB_201904_UPDATE.tgz
```

- Also available on picotte

```
k2_pluspf_20210127   k2_pluspf_20210127.tar.gz
[cnu25@picotte001 kraken2_db_plus_protozoa_fungi]$
```

# RUN KRAKEN2 CONTAINER IN SINGULARITY

▪ Request node to run on:

```
[cnu25@picotte001 data]$ srun --nodes=1 --ntasks=1 --cpus-per-task=32 --mem=120GB --time=00:30:00 --pty /bin/bash
```

▪ Run the docker using singularity, bind to our groups folder

```
[cnu25@node003 eces450650Grp]$ singularity run --bind /ifs/groups/eces450650Grp/ containers/kraken2_latest.sif
```

▪ We are using the run command to execute through the kraken container

```
Usage:
    singularity [global options...]

Description:
    Singularity containers provide an application virtualization layer enabling
    mobility of compute via both application and environment portability. With
    Singularity one is capable of building a root file system that runs on any
    other Linux system where Singularity is installed.

Options:
    -c, --config string    specify a configuration file (for root or
                           unprivileged installation only) (default
                           "/etc/singularity/singularity.conf")
    -d, --debug            print debugging information (highest verbosity)
    -h, --help             help for singularity
        --nocolor          print without color output (default False)
    -q, --quiet            suppress normal output
    -s, --silent           only print errors
    -v, --verbose          print additional information

Available Commands:
    build       Build a Singularity image
    cache       Manage the local cache
    capability  Manage Linux capabilities for users and groups
    config      Manage various singularity configuration (root user only)
    delete      Deletes requested image from the library
    exec        Run a command within a container
    help        Help about any command
    inspect     Show metadata for an image
    instance    Manage containers running as services
    key         Manage OpenPGP keys
    oci         Manage OCI containers
    plugin      Manage Singularity plugins
    pull        Pull an image from a URI
    push        Upload image to the provided URI
    remote      Manage singularity remote endpoints, keyservers and OCI/Docker registry credentials
    run         Run the user-defined default command within a container
    run-help    Show the user-defined help for an image
    search      Search a Container Library for images
    shell       Run a shell within a container
    sif         siftool is a program for Singularity Image Format (SIF) file manipulation
    sign        Attach digital signature(s) to an image
    test        Run the user-defined tests within a container
    verify      Verify cryptographic signatures attached to an image
    version     Show the version for Singularity

Examples:
    $ singularity help <command> [<subcommand>]
    $ singularity help build
    $ singularity help instance start

For additional help or support, please visit https://www.sylabs.io/docs/
[cnu25@node003 eces450650Grp]$
```

# USAGE NOTES (KRAKEN)

`--use-names`: print scientific names

`--db`: Database reference file path

`--fastq-input`[deprecated] : We are using fastq formatted files

`--paired`: We are dealing with paired end data

`--report FILE`: This provides us with a sample-wide report

# KRAKEN2 COMMAND, I/O, RUNTIME

```
Singularity> kraken2 --use-names --threads 4 --db ./data/kraken2_db_plus_protozoa_fungi/k2_pluspf_20210127 --fastq-input -report ./ECES450650_SP21/Tutorial8/reporttest --paired ./data/mappings/evol1.sorted.
unmapped.R1.fastq ./data/mappings/evol1.sorted.unmapped.R2.fastq> ./ECES450650_SP21/Tutorial8/evol1.krakentest
```

- Input files: /ifs/groups/eces450650Grp/data/mappings/evol1.sorted.unmapped.R1.fastq
             /ifs/groups/eces450650Grp/data/mappings/evol1.sorted.unmapped.R1.fastq

- Output files: ./ECES450650_SP21/Tutorial8/evol1.kraken ← classifications
             ./ECES450650_SP21/Tutorial8/report ← report

Execution time:

```
17692 sequences (0.85 Mbp) processed in 0.039s (26969.5 Kseq/m, 1291.14 Mbp/m)
  877 sequences classified (4.96%)
  16815 sequences unclassified (95.04%)
```

# RESULTS & UNDERSTANDING - CLASSIFICATIONS

1.C/U: one letter code indicating that the sequence was either classified or unclassified.

2.The sequence ID, obtained from the FASTA/FASTQ header.

3.The taxonomy ID Kraken2 used to label the sequence; this is **0** if the sequence is unclassified and otherwise should be the NCBI Taxonomy identifier.

4.The length of the sequence in bp.

5.A space-delimited list indicating the lowest common ancestor (in the taxonomic tree) mapping of each k-mer in the sequence. For example, 562:13 561:4 A:31 0:1 562:3 would indicate that:
- the first 13 k-mers mapped to taxonomy ID #562
- the next 4 k-mers mapped to taxonomy ID #561
- the next 31 k-mers contained an ambiguous nucleotide
- the next k-mer was not in the database
- the last 3 k-mers mapped to taxonomy ID #562

Source:
https://genomics.readthedocs.io/en/latest/ngs-taxonomic-investigation/index.html#wood2014

# RESULTS & UNDERSTANDING – REPORT SUMMARY

1. **Percentage** of reads covered by the clade rooted at this taxon

2. **Number of reads** covered by the clade rooted at this taxon

3. **Number of reads** assigned directly to this taxon

4. A rank code, indicating **(U)nclassified, (D)omain, (K)ingdom, (P)hylum, (C)lass, (O)rder, (F)amily, (G)enus, or (S)pecies**. All other ranks are simply **"-"**.

5. NCBI Taxonomy ID

6. The indented scientific name

Source:
https://genomics.readthedocs.io/en/latest/ngs-taxonomic-investigation/index.html#wood2014

# SOURCES

- "6. Taxonomic investigation¶," *6. Taxonomic investigation - Genomics Tutorial 2020.2.0 documentation*. [Online]. Available: https://genomics.readthedocs.io/en/latest/ngs-taxonomic-investigation/index.html#kraken2. [Accessed: 20-May-2021].

- D. E. Wood and S. L. Salzberg, "Kraken: ultrafast metagenomic sequence classification using exact alignments," *Genome Biology*, 03-Mar-2014. [Online]. Available: https://genomebiology.biomedcentral.com/articles/10.1186/gb-2014-15-3-r46. [Accessed: 20-May-2021].

- D. Wood, "DerrickWood/kraken2," *GitHub*. [Online]. Available: https://github.com/DerrickWood/kraken2. [Accessed: 20-May-2021].

- M. Lee, "Genomics," *Happy Belly Bioinformatics*, 2021. [Online]. Available: https://astrobiomike.github.io/genomics/. [Accessed: 20-May-2021].