# Tutorial 12
# MetaPhlAn and HUMAnN

Zoe Sucato, Catherine Nguyen, Tobenna Oduah

# MetaPhlAn [5]

MetaPhlAn stands for **Meta**genomic **Ph**ylogenetic **A**nalysis.

- Computational tool for profiling the composition of microbial communities from metagenomic shotgun sequencing data.
- Relies on unique clade-specific marker genes.
    - From ~17,000 reference genomes:
        - 13,500 bacterial and archeal genomes
        - 3,500 viral genomes
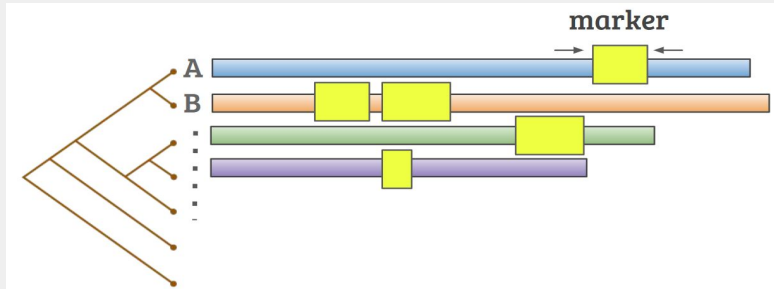        - 110 eukaryotic genomes
- Uses bowtie2.

## Definitions

**Clade:** Group of organisms believed to have evolved from a common ancestor on a phylogenetic tree.
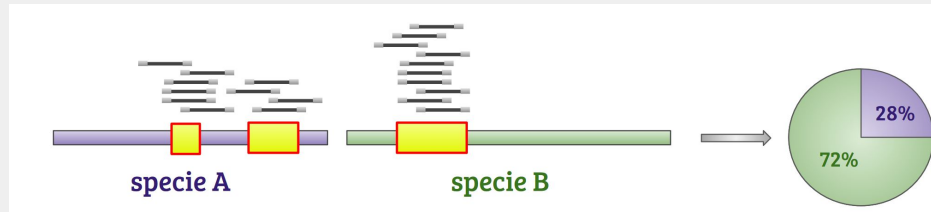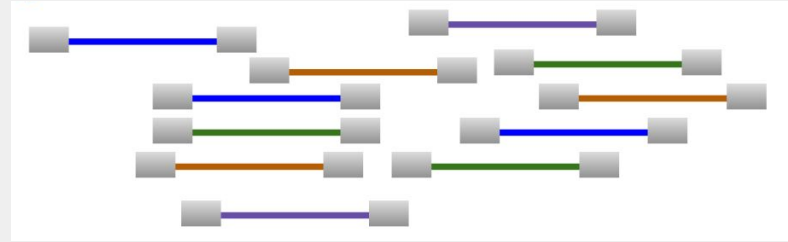
**Clade-specific marker**: Coding sequences that are strongly conserved within the clade's genomes and do not possess substantial local similarity with any sequence outside the clade.

# MetaPhlAn Algorithm [2]



Database of reference genomes and their relationships, with identified clade-specific markers

Sequenced sample

Reads mapped to marker genes

# HUMAnN [3]

HUMAnN stands for the **H**MP **U**nified **M**etabolic **A**nalysis **N**etwork.

- A method for profiling the abundance of microbial metabolic pathways, including other molecular functions from metagenomic/metatranscriptomic sequencing data.

Definitions

**HMP:** Human Microbiome Project. An initiative to research and  understand the microbial components of the human genetic landscape and how that translates to health-related norms and concerns.

# HUMAnN Algorithm [6]

Genes matched to function

Sequenced reads

Species 1    Species 2

Unclassified    Novel

Reads get assigned to species (MetaPhlAn)

Species 1 and 2 marker genes recruit reads

Protein sequence

| Feature | RPK |
|---|---|
| Σ GeneX | 8 |
| GeneX \| Species1 | 2 |
| GeneX \| Species2 | 3 |
| GeneX \| Unclassified | 3 |

# MetaPhlAn

Which microbes are there?

# HUMAnN

What can the microbes do?

# Demo

# Original Plans... Changed

Originally meant to use Biobakery Workflows.

Could not use Docker directly, so we used Singularity to run a Docker image.

Many issues were encountered:

- Pre-existing Singularity Container did not have the tools (wmgz) installed.
- Creating a new Singularity container with the Docker image failed due to incorrect dependencies.
- Using pip install was not viable since many modules were installed and some failed to install.

# New Plan

Use Galaxy/Hutlab! [4]



All online; no package installs.

However, it has limitations.

Example: Doesn't make the output from one process nice for the next process

# Go to the Galaxy/Hutlab Website.



https://huttenhower.sph.harvard.edu/galaxy/

# To upload your fasta file, click "load your own data".

# Drag your fasta file into this upload box.



https://huttenhower.sph.harvard.edu/galaxy/

# Click start to upload.

# Wait until status is 100%.

# The uploaded file should appear in the right panel.

# Navigate to the MetaPhlAn2 profile from the left panel.

# Make sure the correct input is set.



https://huttenhower.sph.harvard.edu/galaxy/

# Keep default settings and click "Execute".

# Successful execution! Job is now added to the queue.

# Result will appear green when it is completed.

https://huttenhower.sph.harvard.edu/galaxy/

# Click the eye icon to download and view result.



https://huttenhower.sph.harvard.edu/galaxy/

# Click the eye icon to download and view result.

# MetaphlAn Output Snippet

Galaxy3-MetaPhlAn2_on_data_2.metaphlan - Notepad

File    Edit    View

```
#SampleID    Metaphlan2_Analysis
k__Viruses  72.88282
k__Bacteria 27.11718
k__Viruses|p__Viruses_noname  72.88282
k__Bacteria|p__Proteobacteria 24.16031
k__Bacteria|p__Bacteroidetes  2.95687
k__Viruses|p__Viruses_noname|c__Viruses_noname  72.88282
k__Bacteria|p__Proteobacteria|c__Gammaproteobacteria  20.08052
k__Bacteria|p__Proteobacteria|c__Alphaproteobacteria  4.0798
k__Bacteria|p__Bacteroidetes|c__Flavobacteriia  2.95687
k__Viruses|p__Viruses_noname|c__Viruses_noname|o__Viruses_noname  54.15491
k__Bacteria|p__Proteobacteria|c__Gammaproteobacteria|o__Enterobacteriales    20.08052
```

# Visualized Output

# Species Detected

# References

[1] "biobakery/biobakery_workflows." n.d. GitHub. Accessed May 20, 2022.
    https://github.com/biobakery/biobakery_workflows.
[2] Borenstein Lab. n.d. "MetaPhlAn." Accessed May 20, 2022. PowerPoint.
    http://borensteinlab.com/courses/TAU_CS_3116_B_19/presentations/7_MetaPhlan.pdf.
[3] The Huttenhower Lab. n.d. "biobakeryWorkflows." The Huttenhower Lab. Accessed May 20, 2022.
    https://huttenhower.sph.harvard.edu/biobakery_workflows/.
[4] The Huttenhower Lab. n.d. "Galaxy / Hutlab." Accessed May 20, 2022.
    https://huttenhower.sph.harvard.edu/galaxy/.
[5] The Huttenhower Lab. n.d. "MetaPhlAn2." The Huttenhower Lab. Accessed May 20, 2022.
    https://huttenhower.sph.harvard.edu/metaphlan2/.
[6] Mehta, Subina, Pratik Jagtap, and Saskia Hiltemann. 2021. "Introduction to metatranscriptomics." Galaxy Training!
    https://training.galaxyproject.org/archive/2021-10-01/topics/metagenomics/tutorials/metatranscriptomics/slides-
    plain.html.

# Questions?