

# Application of FracMinHash to analyse the phylogenetic context of *Phytophthora*

---

Felix Seidel

[felix.seidel@student.uni-tuebingen.de](mailto:felix.seidel@student.uni-tuebingen.de)

June 6, 2024

# Overview

A short story

Background

Phylogenetic Context of *Phytophthora*

Conclusion

## A short story

---

# The setting

Picture of idyllic farms, potato, maybe avocado?

## The villain: *Phytophthora*

Picture of affected plants, maybe three bullet points (yearly damage)

## The hero?

"Unfortunately, we cannot just walk to some mountain and throw jewelrey in it to defeat the villain" Research on the modes of operation, effector genes, and **phylogeny** of that species

## Background

---

# Phylogenetic Outline<sup>1</sup>

---

<sup>1</sup>Bagci et al., “Microbial Phylogenetic Context Using Phylogenetic Outlines”; Bryant and Huson, “NeighborNet”.



---

<sup>2</sup>Bagci et al., “Microbial Phylogenetic Context Using Phylogenetic Outlines”.

# FracMinHash<sup>3</sup> (1)

Idea: Hash  $k$ -mers and keep values smaller or equal to a threshold

- Sequence: *ATGCATGATG*

---

<sup>3</sup>Irber et al., *Lightweight Compositional Analysis of Metagenomes with FracMinHash and Minimum Metagenome Covers*.

Idea: Hash  $k$ -mers and keep values smaller or equal to a threshold

- Sequence: ATGCATGATG
- $k$ -mers: ATG, TGC, GCA, CAT, ATG, TGA, GAT, ATG

---

<sup>3</sup>Irber et al., *Lightweight Compositional Analysis of Metagenomes with FracMinHash and Minimum Metagenome Covers*.

Idea: Hash  $k$ -mers and keep values smaller or equal to a threshold

- Sequence: *ATGCATGATG*
- $k$ -mers: *ATG, TGC, GCA, CAT, ATG, TGA, GAT, ATG*
- hash values: 1, 8, 3, 4, 1, 10, 6, 1

---

<sup>3</sup>Irber et al., *Lightweight Compositional Analysis of Metagenomes with FracMinHash and Minimum Metagenome Covers*.

Idea: Hash  $k$ -mers and keep values smaller or equal to a threshold

- Sequence: *ATGCATGATG*
- $k$ -mers: *ATG, TGC, GCA, CAT, ATG, TGA, GAT, ATG*
- hash values: 1, 8, 3, 4, 1, 10, 6, 1
- define threshold: 3

---

<sup>3</sup>Irber et al., *Lightweight Compositional Analysis of Metagenomes with FracMinHash and Minimum Metagenome Covers*.

# FracMinHash<sup>3</sup> (1)

Idea: Hash  $k$ -mers and keep values smaller or equal to a threshold

- Sequence: *ATGCATGATG*
- $k$ -mers: *ATG, TGC, GCA, CAT, ATG, TGA, GAT, ATG*
- hash values: 1, 8, 3, 4, 1, 10, 6, 1
- define threshold: 3
- FracMinHash sketch: {1, 3}

---

<sup>3</sup>Irber et al., *Lightweight Compositional Analysis of Metagenomes with FracMinHash and Minimum Metagenome Covers*.

## FracMinHash<sup>4</sup> (2)

- $W$ : input sequence
- $h$ : hash function producing values in  $[0, H]$
- $s$ : scaling parameter  $0 < s \leq H$
- $k(W)$ : set of  $k$ -mers of  $W$

$$\text{FRAC}_s(W) = \{h(w) \leq \frac{H}{s} \mid \forall w \in k(W)\} \quad (1)$$

---

<sup>4</sup>Irber et al., *Lightweight Compositional Analysis of Metagenomes with FracMinHash and Minimum Metagenome Covers*.

$$J_{frac}(A, B) = \frac{1}{1 - (1 - \frac{1}{s})^{|A \cup B|}} \frac{|\text{FRAC}_s(A) \cap \text{FRAC}_s(B)|}{|\text{FRAC}_s(A) \cup \text{FRAC}_s(B)|} \quad (2)$$

$$D_{frac}(A, B) = 1 - \left( \frac{2J_{frac}(A, B)}{1 + J_{frac}(A, B)} \right)^{\frac{1}{k}} \quad (3)$$

---

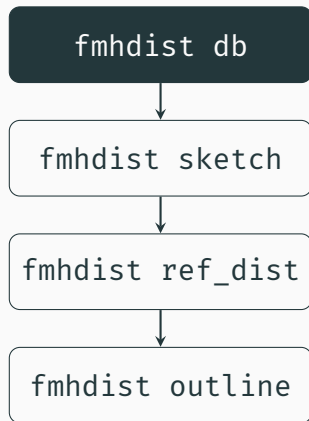
<sup>5</sup>Hera et al., “Deriving Confidence Intervals for Mutation Rates across a Wide Range of Evolutionary Distances Using FracMinHash”.



# Phylogenetic Context of *Phytophthora*

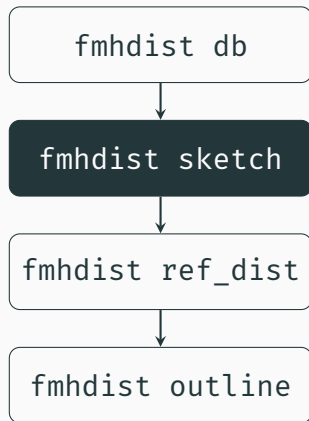
---

## The implementation with fmhdist (1)



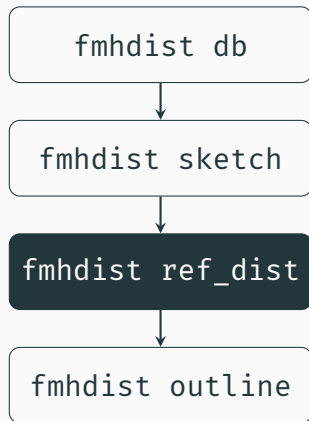
- Input: List of NCBI accession codes
- Sketching parameters:  $s$ ,  $k$ ,  $h$  and random seed for  $h$
- Output: Reference database

## The implementation with fmhdist (2)



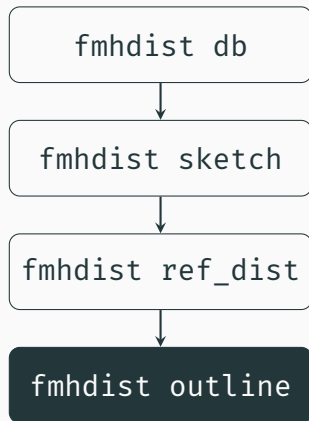
- Input: List of paths to FASTA files
- Sketching parameters:  $s$ ,  $k$ ,  $h$  and random seed for  $h$
- Output: sketches and their coordinates

## The implementation with `fmhdist` (3)

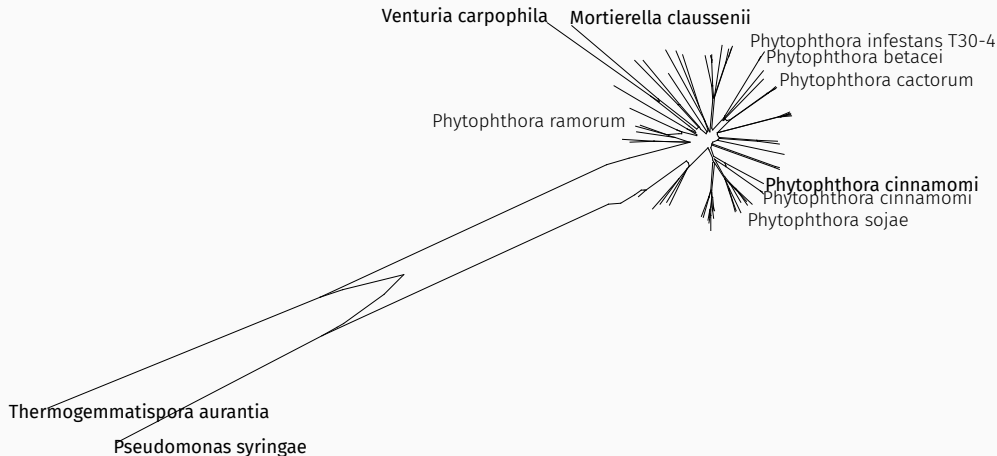


- Input: reference database (`fmhdist db` or `fmhdist sketch`), query sketches (`fmhdist sketch`)
- Distance threshold
- Output: distance matrix

## The implementation with fmhdist (4)

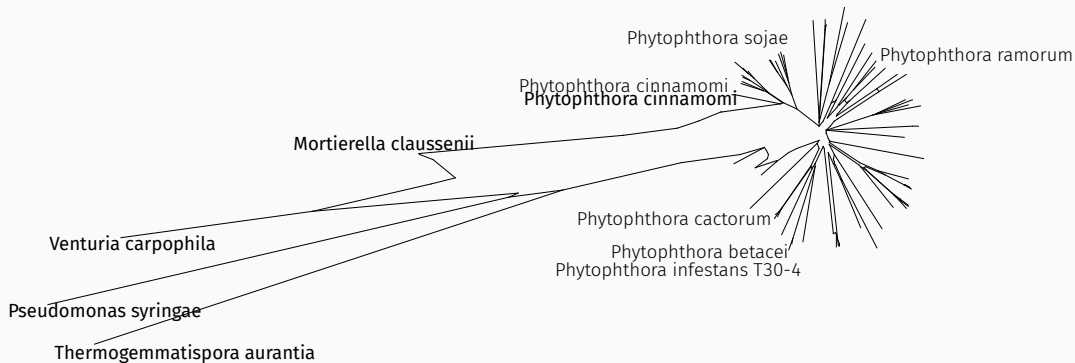


- Input: distance matrix
- Image parameters (width, height, scaling, offset)
- Output: phylogenetic outline as SVG



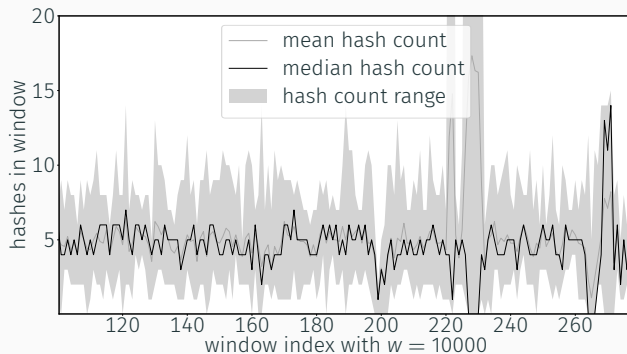
**Figure 1:** Outline based on FracMinHash distances of *Phytophthora* reference sequences and fungal and bacterial query sequences (bold). Only some labels shown.

## Compare this to Mash



**Figure 2:** Outline based on Mash distances of *Phytophthora* reference sequences and fungal and bacterial query sequences (bold). Only some labels shown.

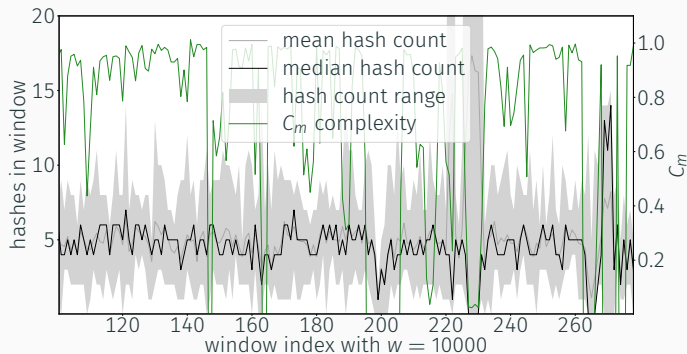
## Origin of the hashes in the sketch of *Phytophthora sojae*



**Figure 3:** Hash counts in windows of the NW\_009258123.1 sequence of the *Phytophthora sojae* reference genome.



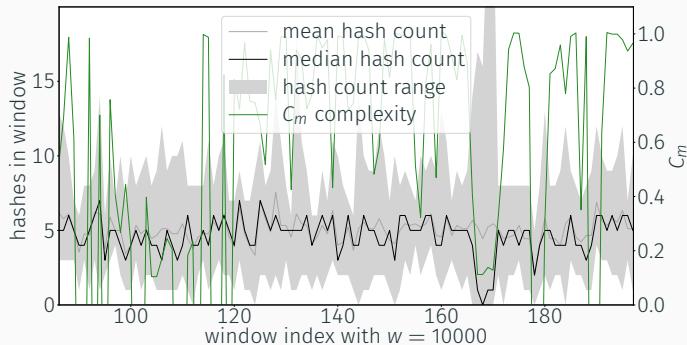
## ...and the corresponding sequence complexity<sup>6</sup>



**Figure 4:** Hash counts and sequence complexity in windows of the NW\_009258123.1 sequence of the *Phytophthora sojae* reference genome.

<sup>6</sup>Pirogov et al., “High-Complexity Regions in Mammalian Genomes Are Enriched for Developmental Genes”.

## How about other genomes?



**Figure 5:** Hash counts and sequence complexity in windows of the NW\_003303751.1 sequence of the *Phytophthora infestans* reference genome.

	$n$	$u$	$r$	$p$
<i>P. cambivora</i>	727	0	-0.0272	
<i>P. betacei</i>	26688	437	0.2446	4.5575e-265
<i>P. cinnamomi</i>	10799	28	-0.0933	4.3767e-19
<i>P. infestans</i>	11173	37	0.1830	2.8352e-24
<i>P. sojae</i>	7339	74	0.1570	6.5134e-48
<i>P. nicotianae</i>	2529	0	-0.0042	

**Table 1:** Excerpt of the statistical analysis of window ( $w = 10000$ ) hash counts and sequence complexity of *Phytophthora* genomes.

# Benchmarks

	mash <sup>7</sup> (1)	mash (6)	sourmash <sup>8</sup> (1)	fmhdist (1)	fmhdist (6)
min (s)	135	44	171	201	75
max (s)	140	58	178	215	91
avg (s)	137	51	174	208	84

**Table 2:** Runtime comparison of three tools calculating sketches for sequences totalling 5.477Gb. The number of threads is in parantheses.

---

<sup>7</sup>Ondov et al., “Mash”.

<sup>8</sup>Irber et al., *Lightweight Compositional Analysis of Metagenomes with FracMinHash and Minimum Metagenome Covers*.

## Conclusion

---

## Method works in general

## Potential next steps

Thanks!