# Fast estimation of pangenome openness using prefix-free parsing

MIGGS - Novel Pangenome Paradigms JOBIM 2025, July 10, Bordeaux

Université de Lille

Igor Martayan aimartayan



### Pangenome Openness

f(n) = # new genes added by the  $n^{th}$  genome ~  $\frac{k n^{-a}}{k n^{-a}}$  (Heap's law)

 $a < 1 \Rightarrow \text{open}$ 



 $a > 1 \Rightarrow closed$ 

Common approach expensive to compute (every genome ordering) [Tettelin et al., 2005]

Faster histogram-based approach for k-mers in  $O(n^2)$  [Parmigiani Wittler Stoye, 2024]

Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: Implications for the microbial "pan-genome"

Hervé Tettelin, Vega Masignani, Michael J. Cieslewicz, +42, and Claire M. Fraser Authors Info & Affiliations

#### PNAS, 2005

Revisiting pangenome openness with *k*-mers

Parmigiani, Luca<sup>1, 2, 3</sup> (b); Wittler, Roland<sup>1, 2</sup> (c); Stoye, Jens<sup>1, 2</sup> (c)  $\uparrow$ 

PCJ, 2024





## Prefix-Free Parsing (PFP)

- Select windows with hash mod p = 0 (p acts as a sampling parameter)
- Selected windows partition the text into phrases



Local consistency: <u>same context  $\Rightarrow$  same selection  $\Rightarrow$  same phrases</u>

The # phrases is a good metric for repetitiveness

Prefix-free parsing for building big BWTs

Christina Boucher 🖾, Travis Gagie, Alan Kuhnle, Ben Langmead, Giovanni Manzini & Taher Mun

#### AMB, May 2019

### Pfp-fm: an accelerated FM-index

Aaron Hong<sup>1</sup>, Marco Oliva<sup>1</sup>, Dominik Köppl<sup>2,3</sup>, Hideo Bannai<sup>3</sup>, Christina Boucher<sup>1\*</sup> and Travis Gagie<sup>4</sup>

#### AMB, April 2024

#### Prefix-free parsing for merging big BWTs

Diego Diaz-Dominguez, Travis Gagie, Veronica Guerrini, Ben Langmead, Zsuzsanna Liptak, Giovanni Manzini, Francesco Masillo, Vikram Shivakumar

Preprint, June 2025









### **Estimating Openness with PFP**



4

### Take-home messages

- PFP is a simple and high-throughput parsing method
- Provides a good estimation of pangenome openness  $(a \text{ such that # new genes } f(n) \sim k n^{-a})$
- Can be reused for lightweight indexing (ongoing work)

Thank you!

### PFP speed on an Apple M1



Number of threads

