



Inria



Pres

Optimization Framework For Phylogenetic Compression

Comité de Suivi Individuelle (1st year)

Tam Truong, Dominique Lavenier, Pierre Peterlongo, Karel Břinda

25 June 2025

I. E
II. T
III. P
I.
II.
II.
IV.
V.
V.
IV. C

My Background: International & Multidisciplinary

2018-2021: Licence MIAGE – Rennes University

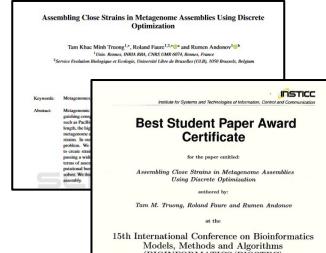
2021: Full-stack developer – Enedis, Paris

2021-2023: Dbl-deg Master Data Science & Business:
Rennes & Aalto Univ (Finland)

2022 & 2023: 2 Internships at INRIA:
Optimization for strains separation - Supervisor: R. Andonov

Start: Nov. 2024
PhD: Optimization for Phylogenetic Compression

Multidisciplinary competences
Big data
Method development & research
Master's internships in bioinformatics



My PhD Thesis

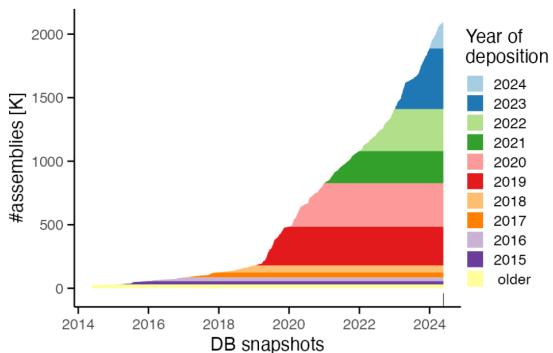
Mot

5

Motivation: Rapidly Growing Bacteria Genome Data & Collections

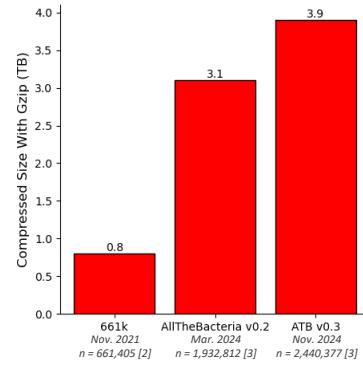
Phyl

Fast Growth Of Bacterial Genomes Data¹ (NCBI)



Karel Brinda, 2024 (CC) <https://doi.org/10.6084/m9.figshare.25879256>

Multi-Terabytes Microbial Genome Collections



18 Jun 2025: ATB v0.4 with 330k new genomes
Next decade: Even larger collections ($n = \sim 10^7$), higher diversity, ...

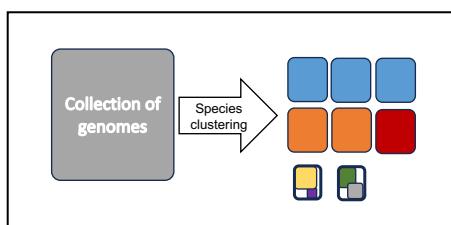
Challenging: Efficient storage & analysis (indexing, searching)
Standard compression protocol is not sufficient

[1] Brinda et al. 2025, [2] Blackwell et al. 2021, [3] Hunt et al. 2024

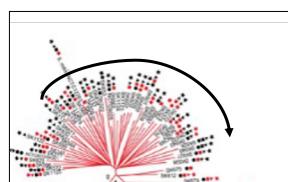
7

Phylogenetic Compression

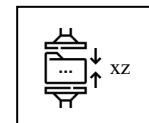
Phylogenetic Compression: Key Steps – Implemented in **MiniPhy**



Step 1 : Phylogenetic Species-based
Pre-ordering & Batching



Step 2 : Phylogenetic Reordering Per Batch



Step 3 : Compression

Břinda et al. 2025, Blackwell et al. 2021, Hunt et al. 2024
MiniPhy: github.com/karel-brinda/miniphy

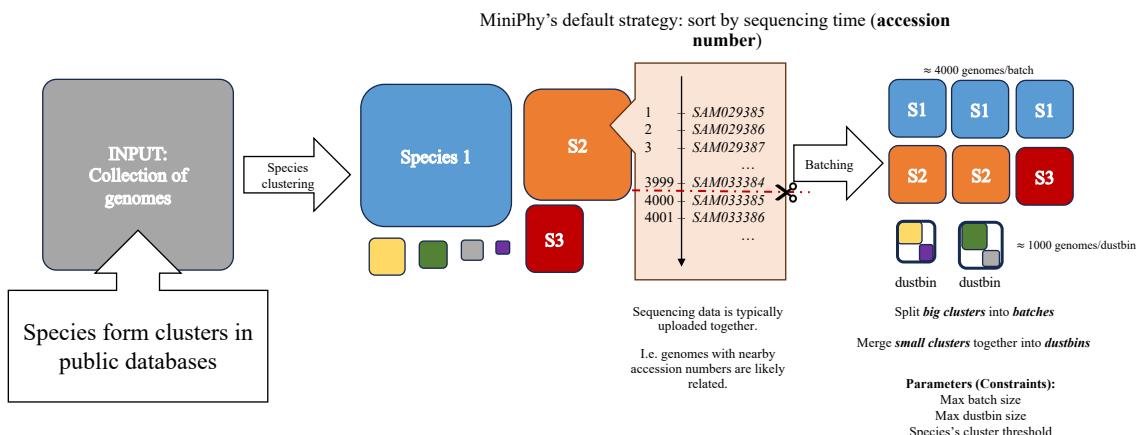
9

Order of Operations

Key Concepts

Inference

Batching In Phylogenetic Compression (MiniPhy)



MiniPhy: github.com/karel-brinda/miniphy

13

Lack Of Suitable Methods For The Phylogenetic Batching Step

No Formalization Of Batching As An Optimization Problem.

Batching Are Not Suitable For Hardware Specific Target Application.

Heavily Dependent On Metadata For An Approximative Input Order

Order-awareness

Constraints on:
Uncompressed sizes
Compressed sizes
Number of genomes
Bounds on size of the used search indexes

Search on GPUs or for processing-in-memory (PIM) architectures

Accession number
Species labels

Opt

Optimization Problem Formulation For Phylogenetic Compression

Exam

Inputs:	Objective function:	
<ul style="list-style-type: none"> $G = \{g_1, g_2, \dots, g_n\}$: set of genomes We want to split G into m ordered batches, $m \leq n$ $B = \{b_1, b_2, \dots, b_m\}$: set of <u>ordered batches</u> 	<ul style="list-style-type: none"> Minimize total compressed batch sizes : $\min \sum_{j=1}^m C(b_j) \cdot y_j$	
Parameters:	<ul style="list-style-type: none"> Find the order to achieve the best compression 	
<ul style="list-style-type: none"> u : bound on uncompressed size c : bound on compressed size e : bound on number of genomes 	<ul style="list-style-type: none"> Minimize the number of batches : <ul style="list-style-type: none"> Fewer batches \Rightarrow better compression 	
	<ul style="list-style-type: none"> Combined: $\min \sum_{j=1}^n C(b_j) \cdot y_j + \sum_{j=1}^n y_j$	
Decision Variables:	Subjects to (possible) constraints:	Subjects
<ul style="list-style-type: none"> $x_{ij} \in \{0,1\}$: 1 if genome g_i is in batch b_j, 0 otherwise $y_j \in \{0,1\}$: 1 if batch b_j is used, 0 otherwise 	<ul style="list-style-type: none"> 1) $\sum_{j=1}^m x_{ij} = 1 \quad \forall i \in \{1, \dots, n\}$: a genome must be assigned 2) $x_{ij} \leq y_j \quad \forall i, j$: no genomes in unselected batches 3) $U(b_j) \leq u \cdot y_j \quad \forall j \in \{1, \dots, m\}$: bound on uncompressed size 4) $C(b_j) \leq c \cdot y_j \quad \forall j \in \{1, \dots, m\}$: bound on compressed size 5) $\sum_{i=1}^n x_{ij} \leq e \cdot y_j \quad \forall j \in \{1, \dots, m\}$: bound on genomes count per batch 	<ul style="list-style-type: none"> 1) 2) 3)
Functions:		
<ul style="list-style-type: none"> $U(b_j) = f_{\text{uncompressed}}(\text{ordered } \{g_i : x_{ij} = 1\})$ Exp: disk size of a batch $C(b_j) = f_{\text{compressed}}(\text{ordered } \{g_i : x_{ij} = 1\})$ Exp: xz compression of a batch (param: level 9 compression,...) 		

17

The Two Tracks Of My Work: Preordering & Partitionning

Axi

OBJECTIVE:

$$\min \left[\sum_{j=1}^n C(b_j) \cdot y_j \right] + \left[\sum_{j=1}^n y_j \right]$$

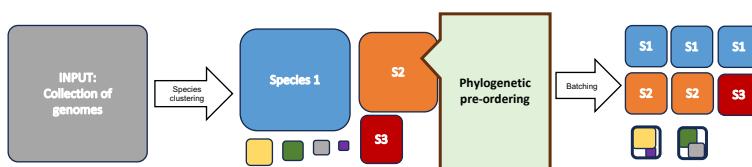
1. Phylogenetic Pre-ordering Of Genomes

2. Order-based Genome Batching

19

Track's Objective: Achieving Global Phylogenetic Pre-order At The Million-genome Scale

Skele



Infere phylogenetic tree for each species : Distance estimation (Mash) + Neighbor joining tree (Quicktree)

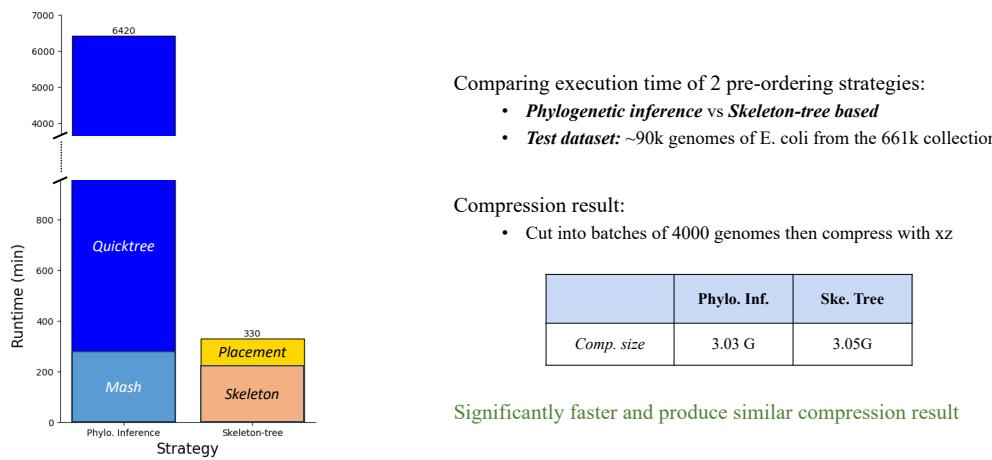
Challenging for highly sampled species (>50k genomes)

Infeasible for large modern collection (million-genome scale)

Ondov et al. 2016, Howe, Bateman, and Durbin 2002
Attotree: github.com/karel-brinda/attotree

21

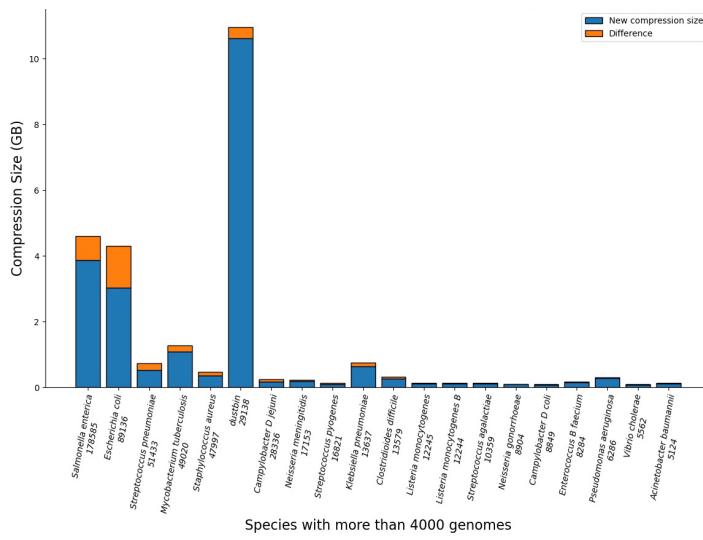
Run Time Comparison: >10x Improvement Compared To Standard Strategy



Resu

23

Result (661k) Species-wise: Absolute Compressed Size Reduction P. Species

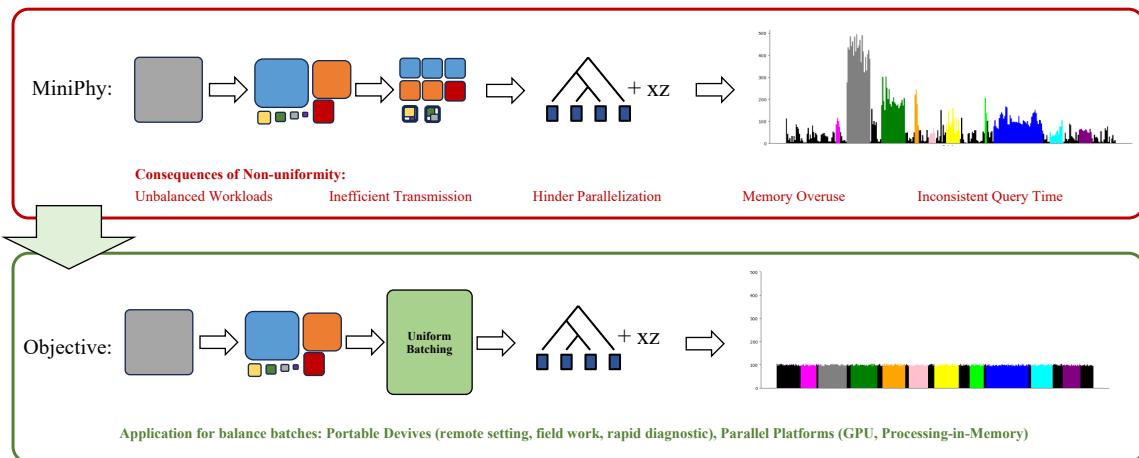


25

Axis 2: Order-based Genome Batching



Consequences and Applications of Uniform Batches



Bin 1

Inputs:
•
•
•
•
•

Parameter:
•
•
•
•
•

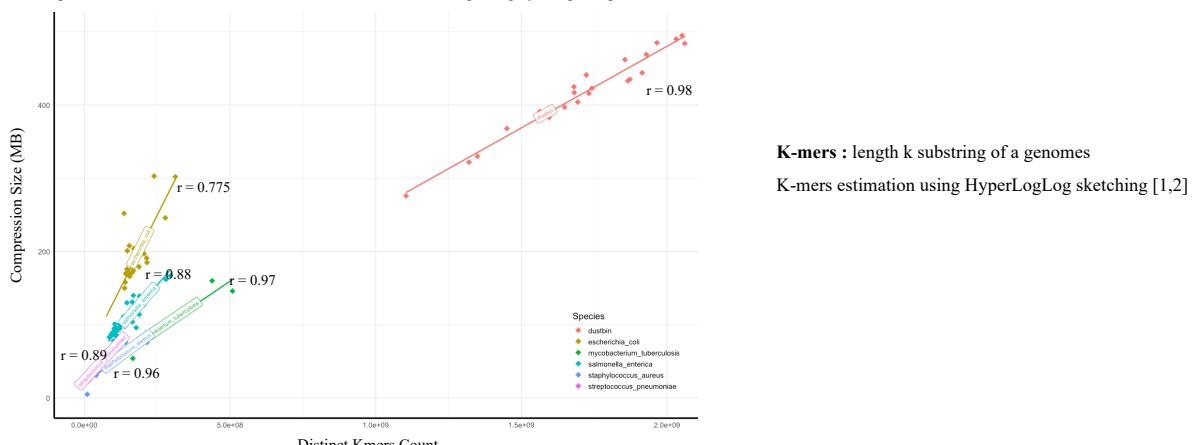
Decision V:
•
•
•
•
•

Functions:
•
•
•
•
•

Subjects to:
1)
2)
3)

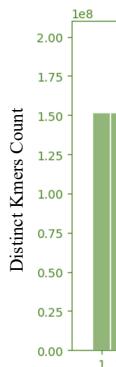
Compression Size Estimation Using Proxy: Distinct K-mers Counts

Compression Size Vs Distinct Kmers Count – 661k Collections – Top 5 Highly Sampled Species & Dustbin



Axis 2 Experimental Result:

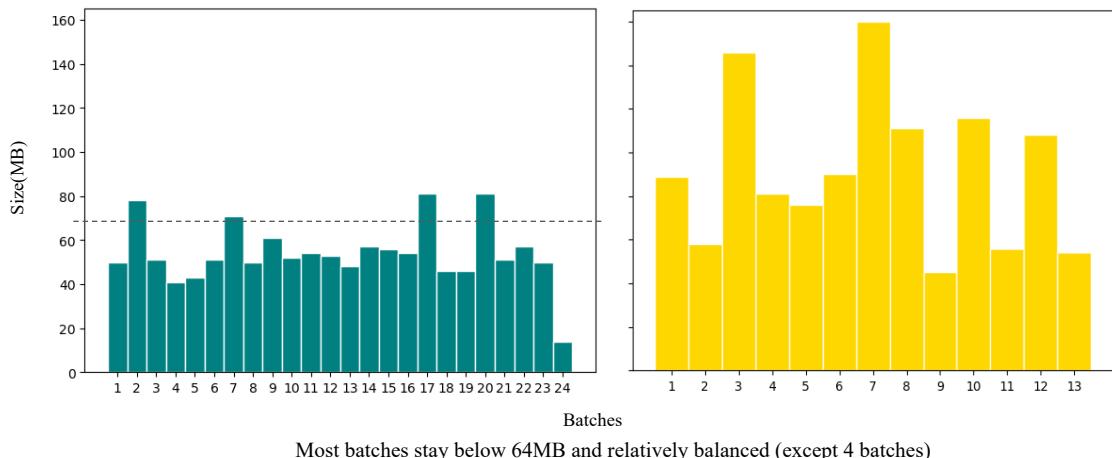
- Dataset: Assemblies of *Mycobacterium tuberculosis* from 661k collection
 - Number of Genomes: around 49,000
 - We want the compression size stay below 64MB (DRAM for PIM [1,2])
 - CAPACITY (distinct kmers count) of batches: 152,000,000



Ghose et al. 2019; Mutlu et al. 2019

33

Compression size per batch



- Post
 - Present
 - **Next**
 - Future