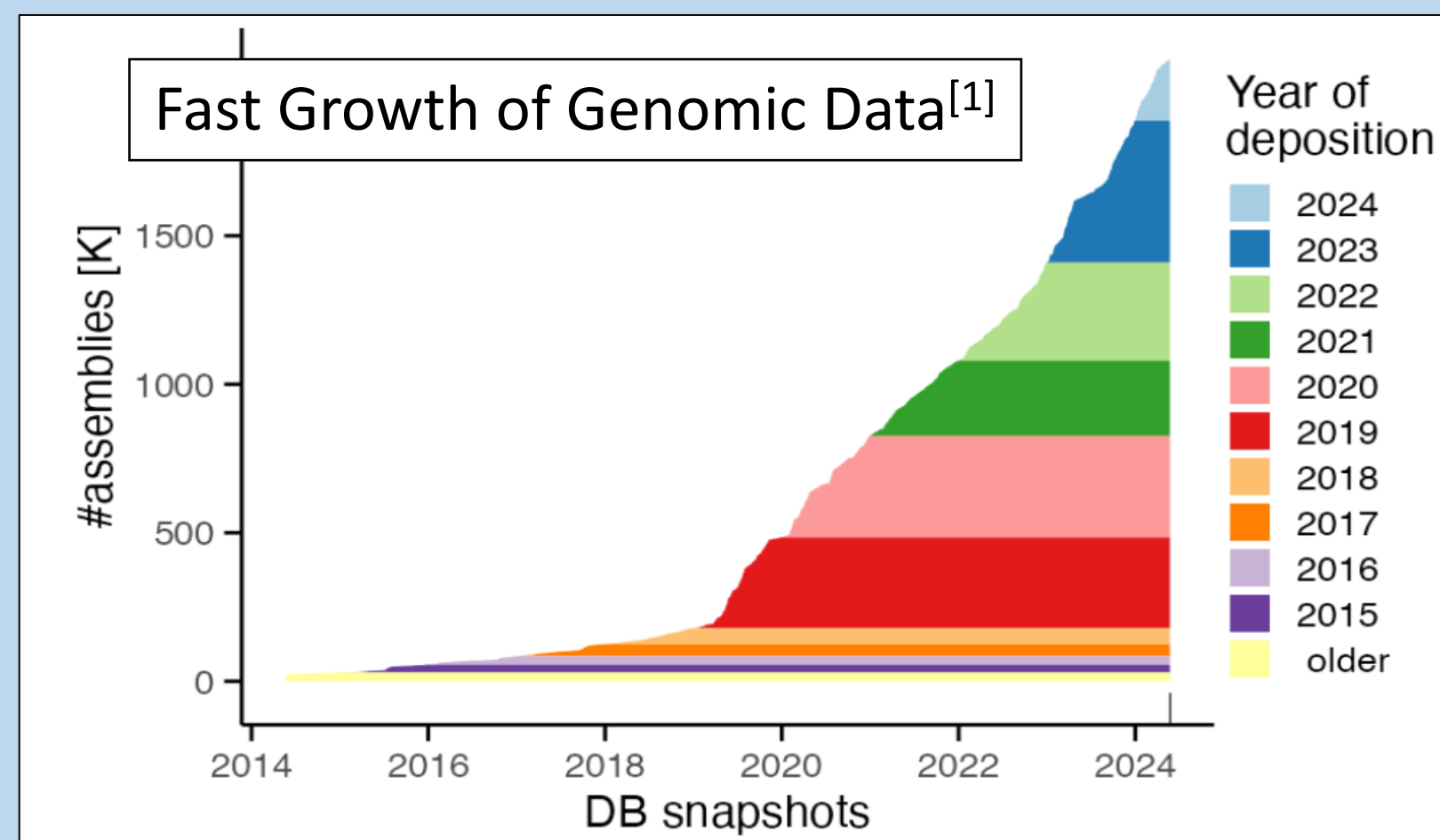


HyperLogLog-Based Load Balancing and Bin Packing for Efficient Compression and Querying of Bacterial Genomes

Motivation	Method	Conclusion and Perspectives
<p>Bacterial genomes are growing rapidly</p> <p>Increasing Availability of Larger Bacterial Genome Collections</p> <ul style="list-style-type: none">661k (Blackwell et al., PLOS Biology, 2021) $n = 661,405$AllTheBacteria (Hunt et al., bioRxiv, 2024) $n = 2,440,377$FUTURE : collection of ten of millions bacteria genomes <p>Efficient storage and query on bacteria collections</p> <p>805 G (standard protocol) vs 29 G (miniPhy-xz) for 661k</p> <p>High compression ratio Compressed using MiniPhy: https://github.com/karel-brinda/miniphy</p> <p>Phylogn: tool for search of across the compressed 661 collection on a standard laptop. https://github.com/karel-brinda/Phylogn</p> <p>Genome compression: clustering/batching</p> <p>Constraint (i.e. memory) while minimizing the number of batches (number of genomes in the batches)</p> <p>Bin Packing Optimization Problem</p> <p>Instance:</p> <ul style="list-style-type: none">Finite set I of items, a size $s(i) \in \mathbb{Z}^+$ for each $i \in I$.Given a capacity K for each batch <p>Objective: Partition the items in set I into batches and minimize the number of batches used</p> <p>Genomes are extensively studied that there are efficient heuristic algorithms to solve this problem for DNA data</p>	<p>Compressed size of batches can be estimated using biological property</p> <p>Fast distinct kmer estimation using probabilistic counting – HyperLogLog sketching</p> <p>Sketching: generating an approximate, compact summary of data (a sketch)</p> <p>TCGATCGATCGATCGA</p> <p>TCGAT</p> <p>CGATC</p> <p>GATCG</p> <p>...</p> <p>ATCG</p> <p>TCGA</p> <p>Probabilistic counting/approximate counting:</p> <p>Sketches</p> <p>Hash(item_1) = 0000... $\Rightarrow \text{Prob} = 1 / 2^4$</p> <p>Hash(item_2) = 0001...</p> <p>Hash(item_3) = 0010...</p> <p>Hash(item_4) = 0011...</p> <p>Hash(item_5) = 0100...</p> <p>Hash(item_6) = 0101...</p> <p>...</p> <p>Longest prefix of leading zeros in the hash values is 4 Estimated cardinality = 2^4</p> <p>Set cardinality estimation with Dashing:</p> <p>Hyperloglog is a sketching algorithm based on advanced approximate counting.</p> <p>Baker, D.N., Langmead, B. implemented it in the tool dashing. https://github.com/dnbaker/dashing</p> <p>Hyperloglog based bin packing and load balancing for genomes</p> <p>HLL-binning: Put genomes into batches that fit within a user-defined memory constraint</p> <p>HLL-balancing: Put genomes into balanced, user-defined number of batches</p> <p>INPUT: GENOMES SKETCHES, CAPACITY (max_distinct_kmers)</p> <p>Initialize: Sort the genomes based on their accession number.</p> <p>Place items in BATCH: For each genome: Try to place it in the first batch that can accommodate it (based on capacity). If it fits, update the batch and move on to the next genome.</p> <p>Create New BATCHES if Necessary: If no batch can accommodate the current genome, create a new batch and place the genome there.</p> <p>https://github.com/sam-km-truong/HLL-Binning</p> <p>INPUT: GENOMES SKETCHES, NUMBER OF BATCH b</p> <p>Initialize: Sort the genomes based on their accession number.</p> <p>Initial Assignment: Assign the first b genomes directly to b batches.</p> <p>Greedy Assignment: For each remaining genomes, place it in the batches with the smallest number of distinct kmers.</p> <p>Update the batch's contents and size.</p> <p>https://github.com/sam-km-truong/HLL-Balancing</p> <p>EXPERIMENT: batching of Mycobacterium Tuberculosis</p> <p>Max distinct kmers for batches: 153000000 kmers so that the compressed size stays under 64MB (calculated by using the correlation of compressed size and distinct kmers count)</p> <p>Using the number off batches found by HLL-Binning, run HLL-Balancing with numbins $b = 24$</p> <p>Eventhough the baches are balances, HLL-Balancing batches have higher distinct kmers count than the capacity set in HLL-Binning</p>	<p>Conclusion</p> <ul style="list-style-type: none">Genomic data benefits greatly from compression when guided by evolutionary characteristicsClustering and batching of genome collections improve compression ratio and facilitate parallel processingThe compressed size of batches (using HyperLogLog) is correlated with the distinct k-mer countThe distinct k-mer count can be efficiently estimated using HyperLogLog sketching (implemented in dashing)The Bin Packing heuristic algorithm can be used to define capacity, i.e., distinct k-mer countThe Load Balancing heuristic algorithm can be used to create balanced batches. <p>Perspectives:</p> <ul style="list-style-type: none">Combining HLL-Binning and HLL-BalancingStudying the collection as a whole instead of individual genomesExtending the scope to include different types of genomic data (e.g., metagenomes, transcriptomes, etc.) <p>Bibliography</p> <p>Grace A. Blackwell, Martin Hunt, Kerri M. Malone, Leandro Lima, Zamin Iqbal, and Zamin Iqbal. 2021. Exploring bacterial diversity in archived DNA sequences. <i>PLOS Biology</i> 19, 11 (November 2021).</p> <p>Martin Hunt, Leandro Lima, Daniel Anderson, Jane Hawkey, and Zamin Iqbal. 2024. <i>AllTheBacteria - all bacterial genomes assembled, available</i>. bioRxiv.</p> <p>Karel Břinda, Leandro Lima, Simone Pignotti, Natalia Quinones, Zamin Iqbal, and Michael Baym. 2024. <i>Efficient and Scalable Phylogenetic Compression</i>. bioRxiv.</p> <p>Po-Ru Loh, Michael Baym, and Bonnie Berger. 2012. Compressed genomic data. <i>Genome Biology</i> 13, 11 (2012).</p> <p>Will P. M. Rowe. 2019. When the levee breaks: a practical guide to the flood of genomic data. <i>Genome Biology</i> 20, 1 (September 2019).</p> <p>Jessica K. Bonnie, Omar Y. Ahmed, and Ben Langmead. 2024. Genomic growth and similarity. <i>iScience</i> 27, 3 (March 2024).</p>

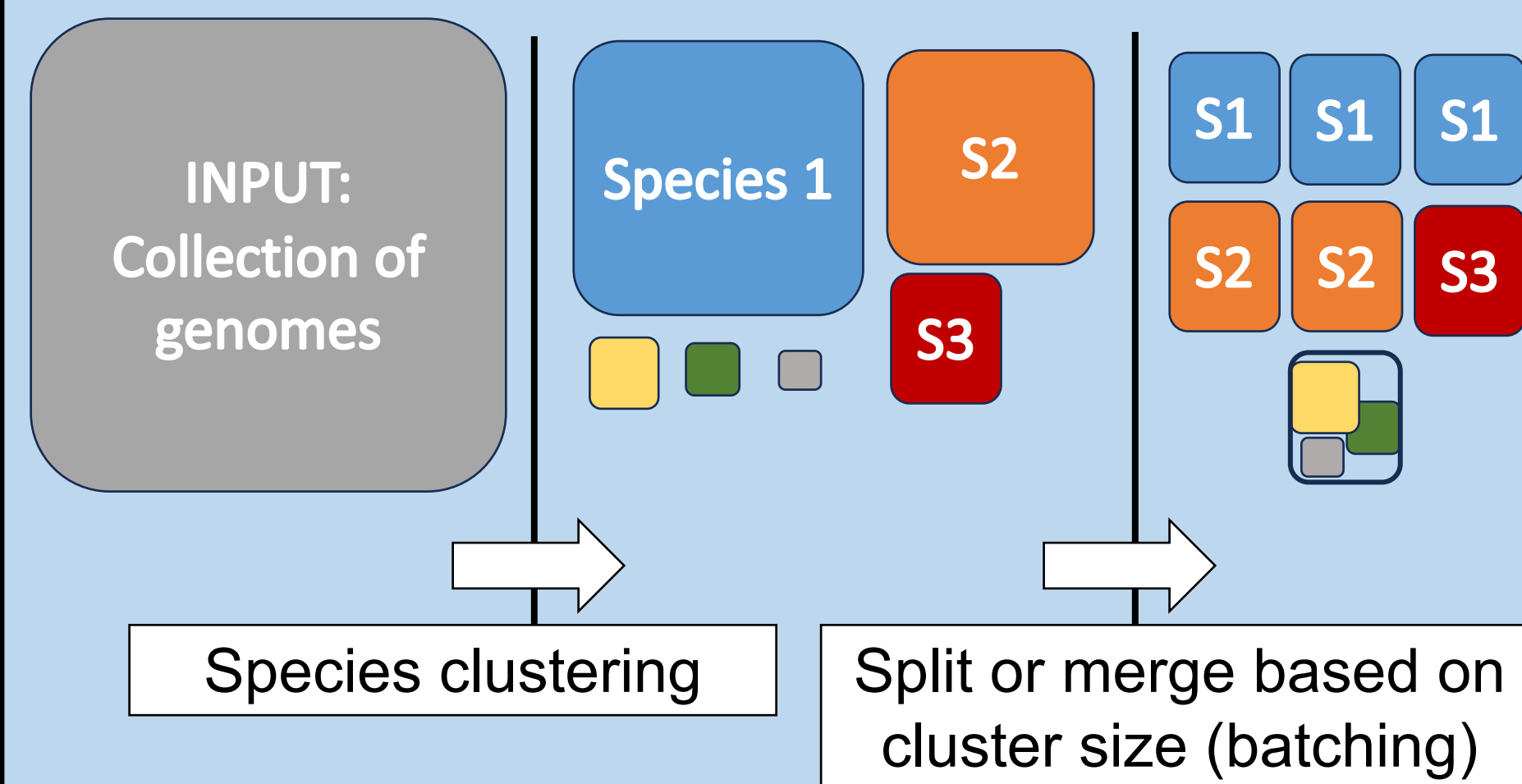
MOTIVATION



Large Bacterial Genome Collections:
661k collection^[2] (2021) n = 661,405
AllTheBacteria^[3] (2024) n = 2,440,377
Future collections could have 10 of millions genomes

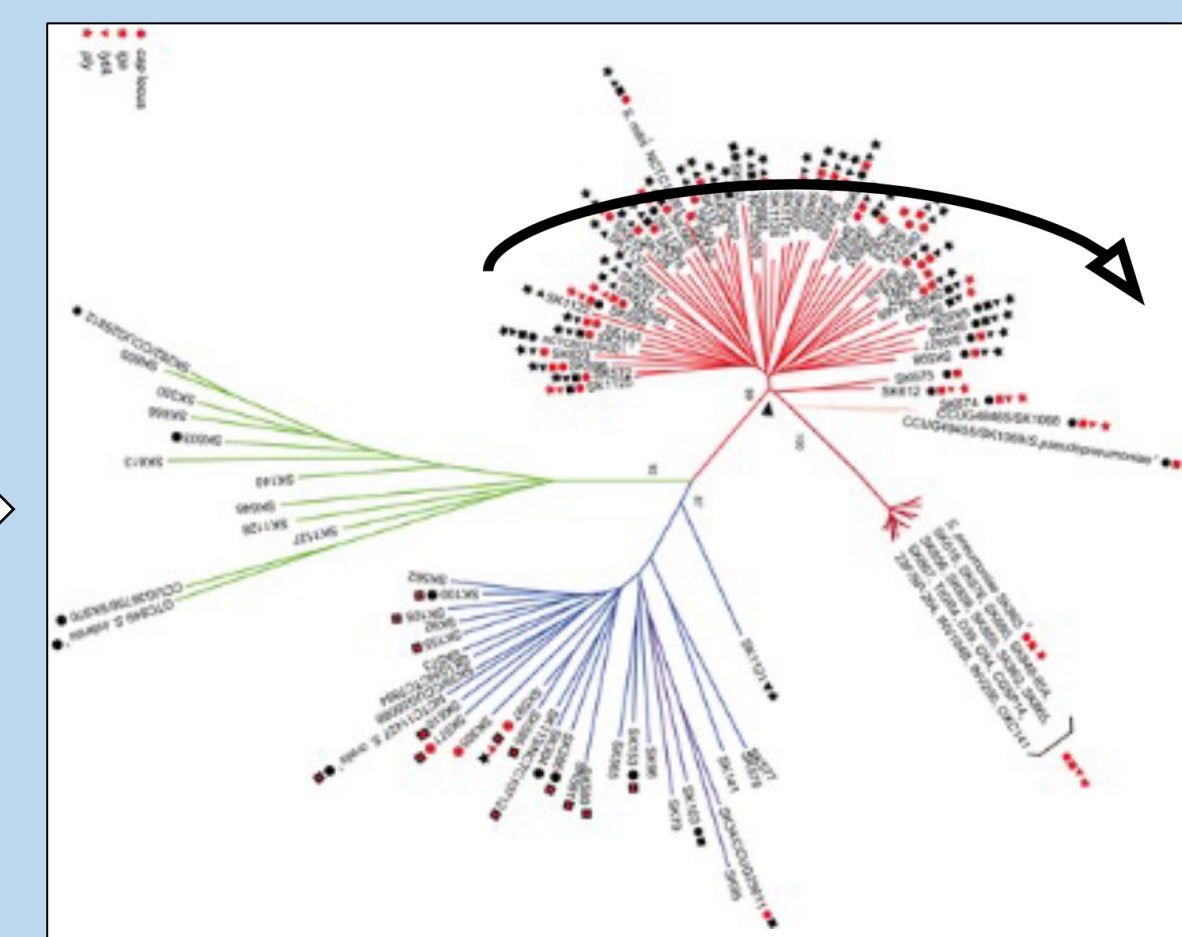
PHYLOGENETIC COMPRESSION: REORDERING STEPS

STEP 1 : PHYLOGENETIC CLUSTERING/BATCHING



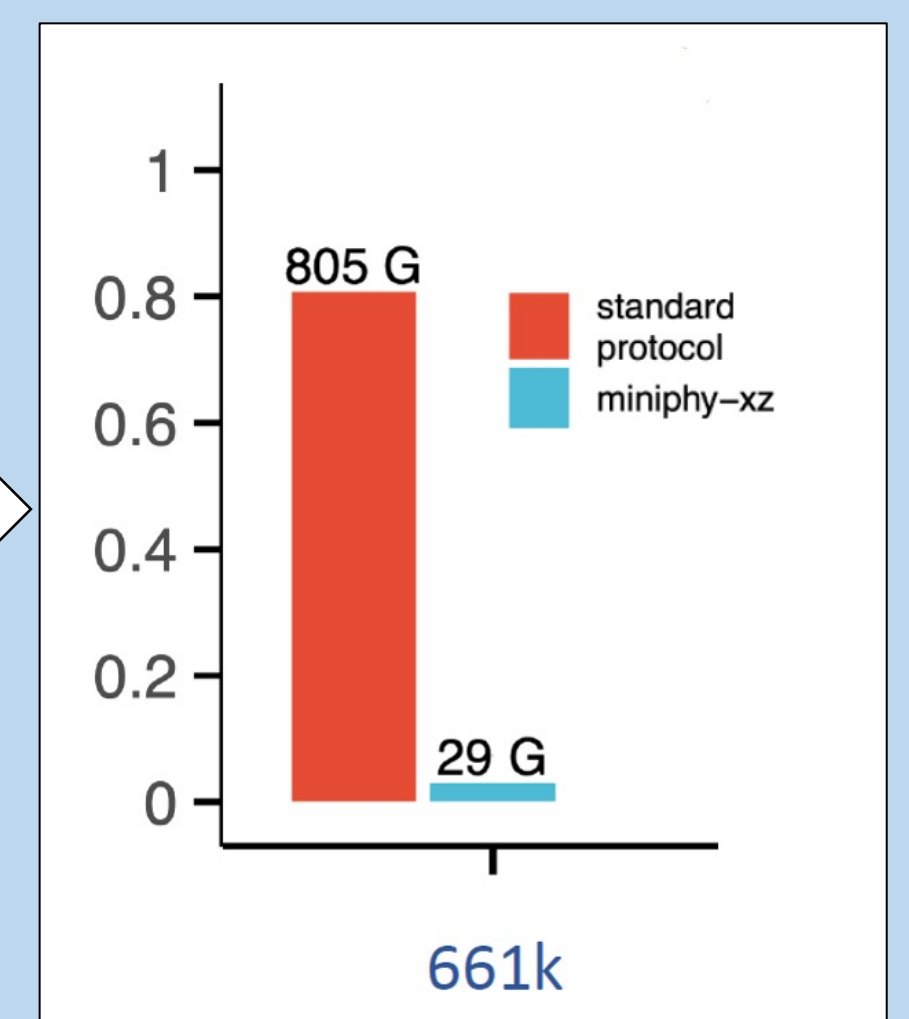
The collection is first partitioned into species clusters. These clusters are then split or merged into batches based on the number of genomes.

STEP 2: COMPRESSIVE PHYLOGENY



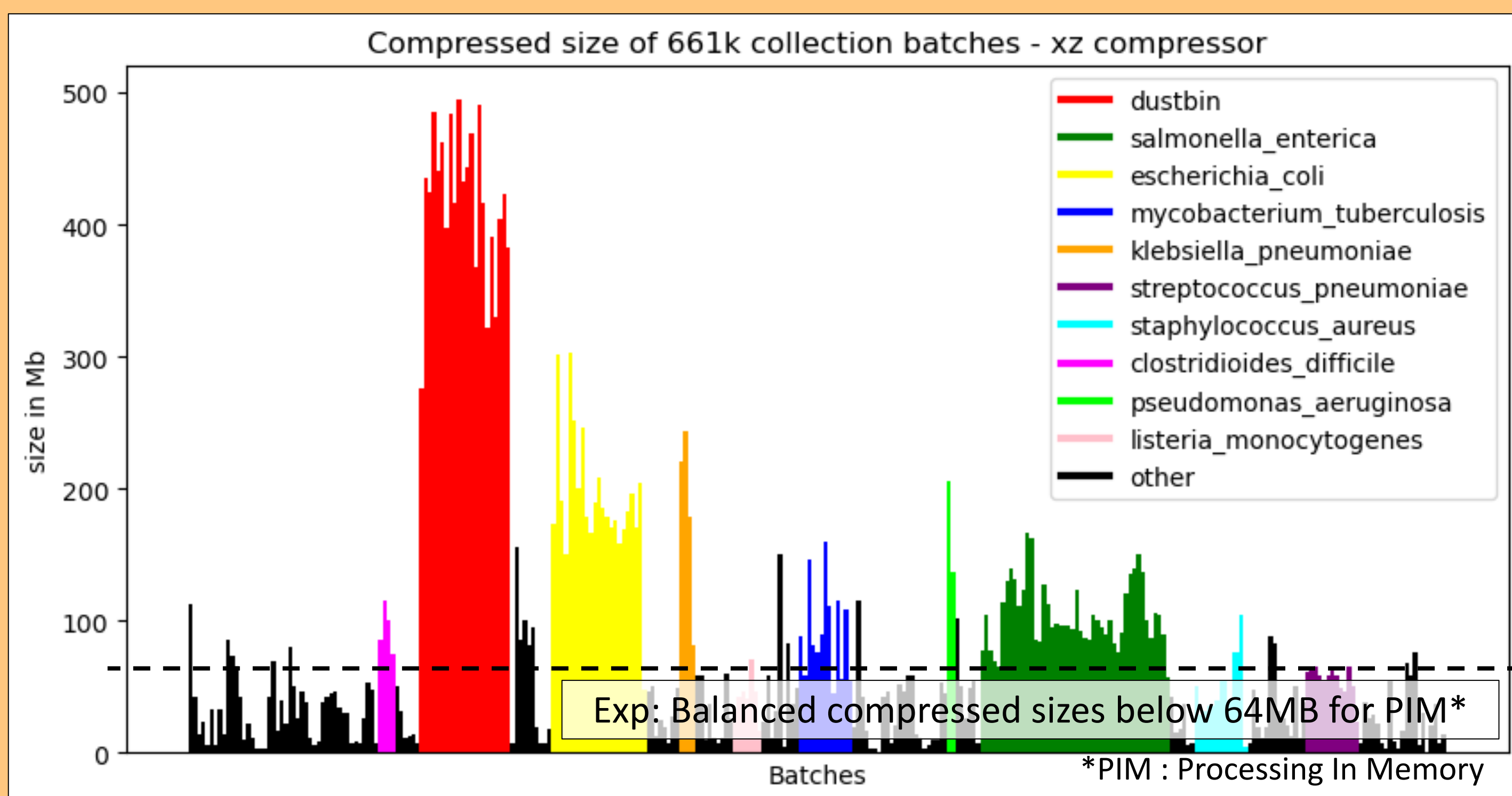
COMPRESSION

STEP 3: COMPRESSION



Lossless compression of 1-3 orders of magnitude

CURRENT LIMITATION: Batching Results In Non-uniform Compressed Sizes



PROBLEM

Unreliable data transmission over bad network

Hinder Parallelization

ULTIMATE OBJECTIVE

Given
Clusters of genomes
Hardware platform

Objective :

$$\text{minimize } \sum \text{ressource}(\text{batch})$$

Per-batch Constraints :

- Bounds on compressed size
- Bounds on decompressed size
- Bounds on number of genomes

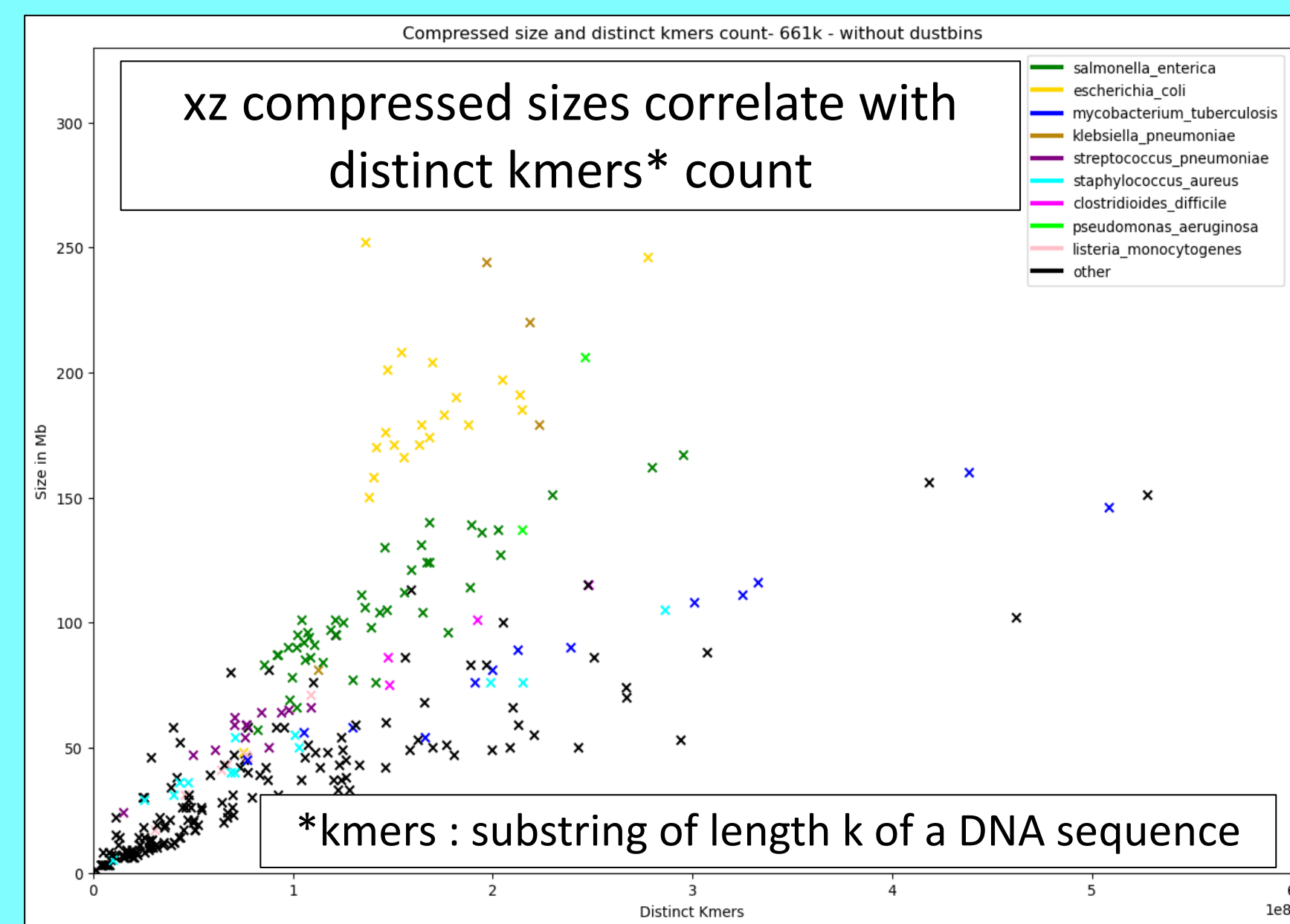
For:

- Data Transmission
- Parallelization

CURRENT GOAL

Create balanced batches after compression that stay below a size constraint and minimizing the number of batches needed

METHODS



Cardinality estimation using HyperLogLog sketching

Sketches : approximate data structures.
HyperLogLog sketches : bit patterns,
i.e. $\text{hash(ATGCG)} \rightarrow 00010100, \text{hash(CGTAC)} \rightarrow 00000010$.
Fast and efficient UNION operation for sketches.

Implemented in the Dashing^[5] tool.
Average relative error ^[4] : 6.537×10^{-4}

MAIN IDEA:

Prediction of Genome Batch's Compression Size Using Distinct Kmers Estimation

HyperLogLog Based Load Balancing^[6] and Bin Packing^[7]

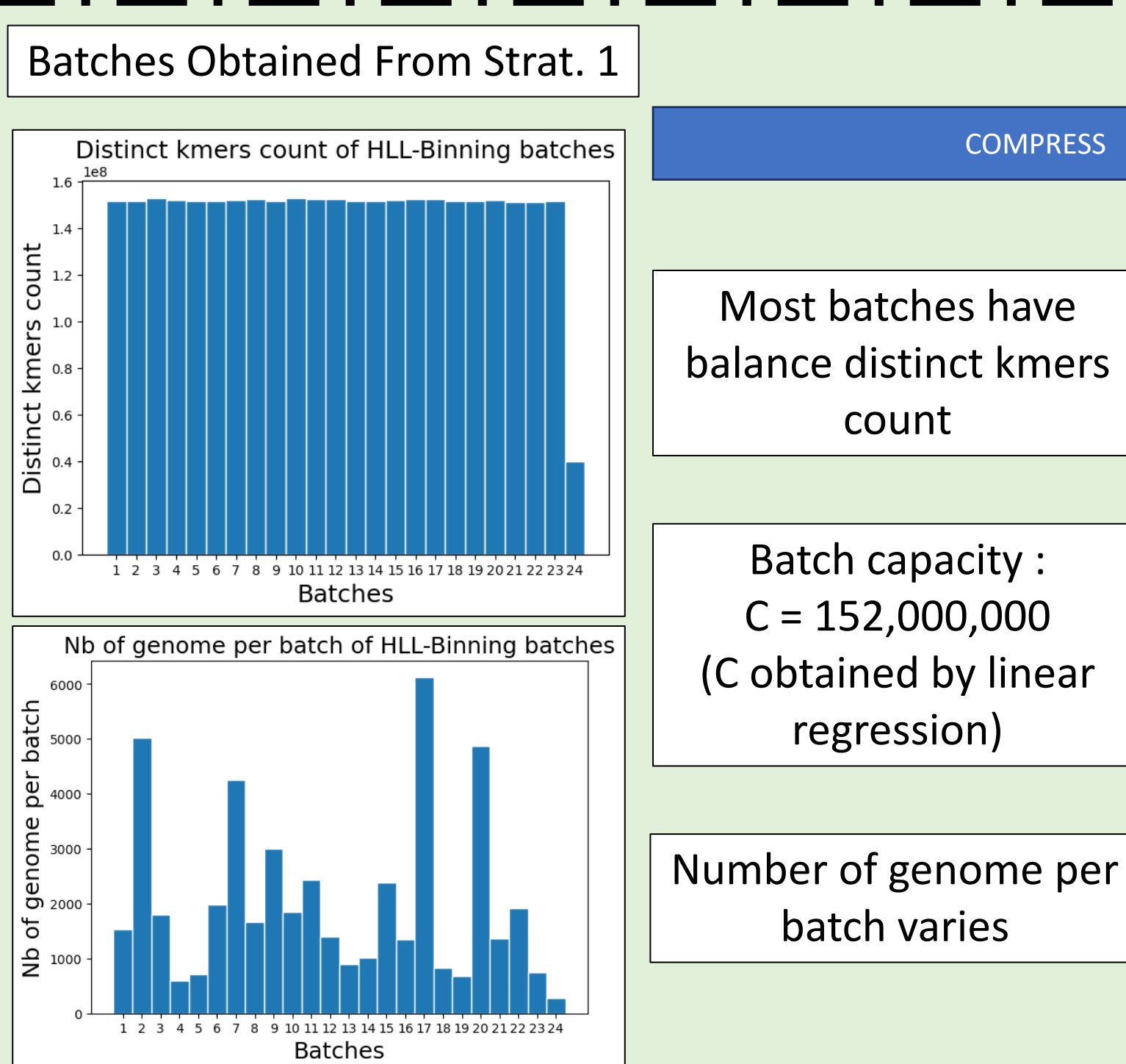
HLL-binning:
Given unlimited batches with capacity C, m genomes, put genomes into batches:

$$\text{Minimize nb_of_batch } B = \sum_{j=1}^n b_j, \quad \text{for } (j = 1, \dots, n)$$

s.t. $\text{distinct_kmers}(b_j) < C$, for $(j = 1, \dots, n)$
Greedy implementation : first-fit bin packing
<https://github.com/tam-km-truong/HLL-Binning>

HLL-balancing:
Given n batches, m genomes, put genomes into batches:
Minimize $\max(\text{distinct_kmers}(b_j))$, for $(j = 1, \dots, n)$
Greedy partitioning algorithm :
<https://github.com/tam-km-truong/HLL-Balancing>

STRATEGY 1 HLL-Binning



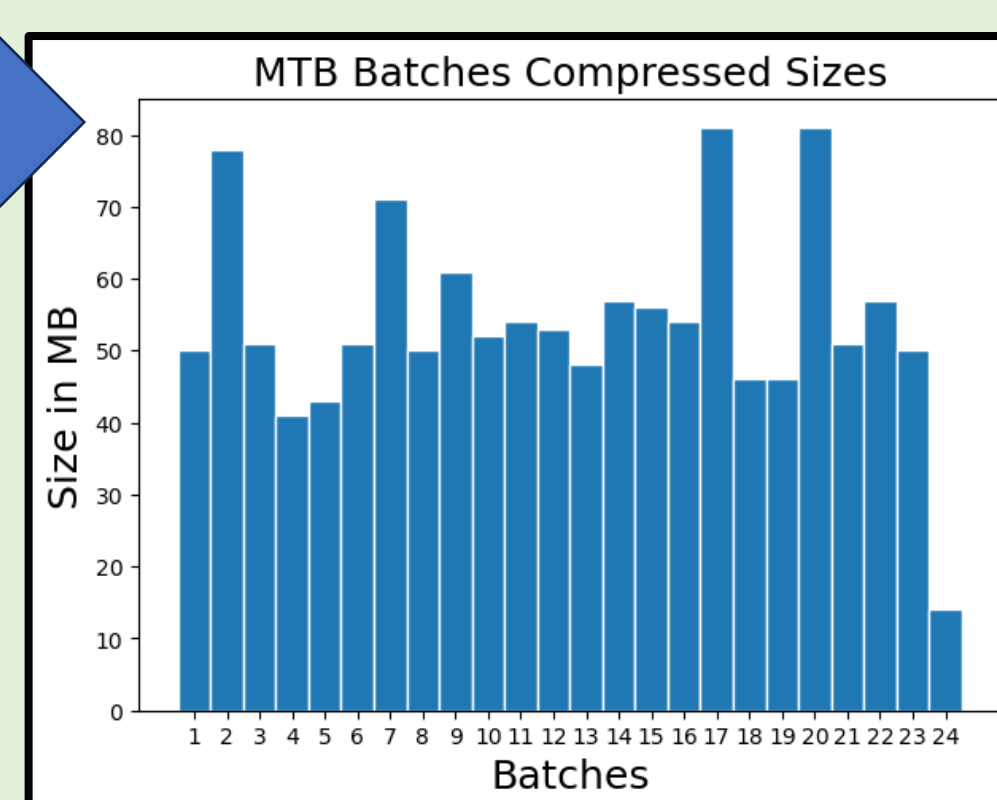
Most batches have balance distinct kmers count

Batch capacity : C = 152,000,000 (C obtained by linear regression)

Number of genome per batch varies

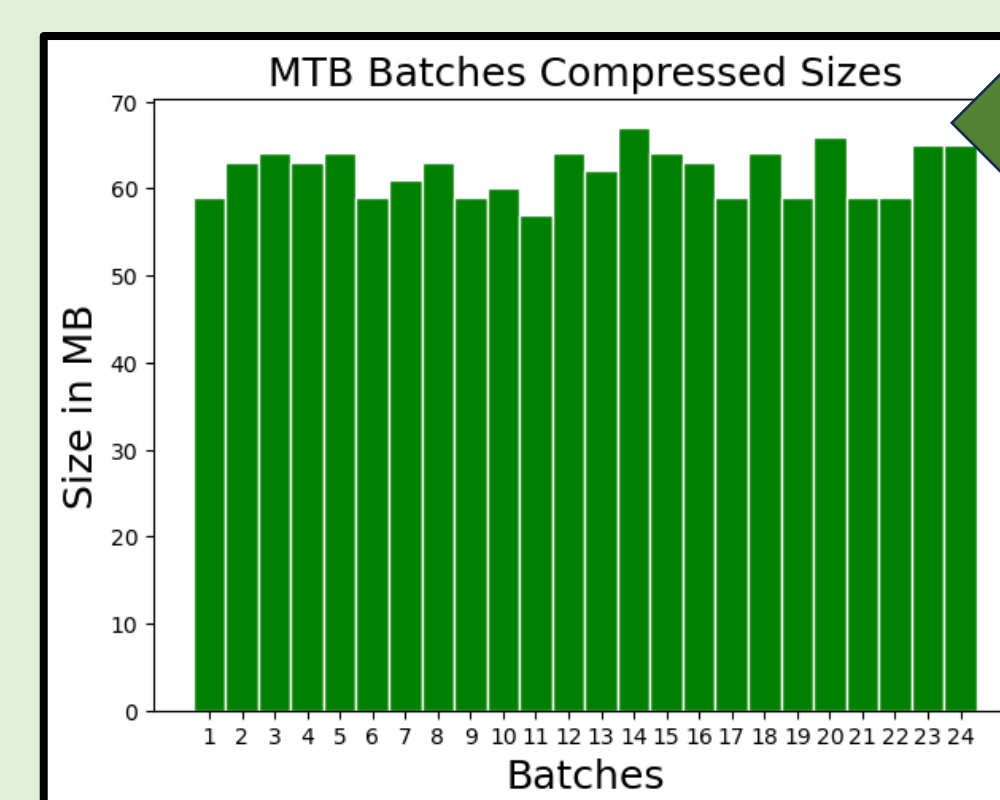
PRELIMINARY RESULTS

DATA : Genomes of Mycobacterium Tuberculosis (MTB) from the 661k Collection^[2]



Most of the batches are balanced (between 40-50MB)

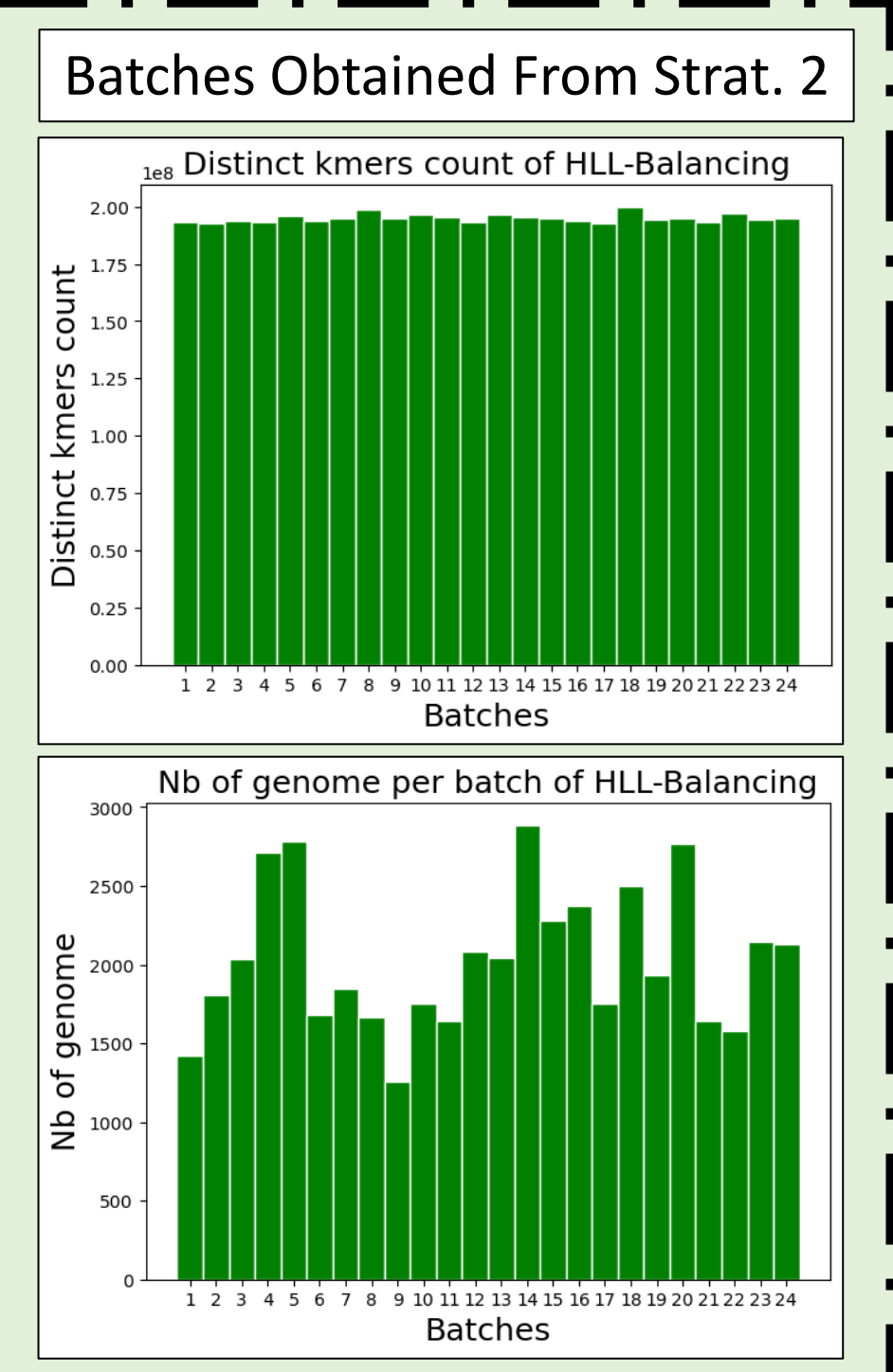
Evaluation strat. 1:
Allowing a capacity on distinct kmers. The result remain somewhat imbalanced.



All Batches are well balanced (max size is 69MB)

Evaluation strat. 2:
Producing more balanced batches. No control over the maximum distinct k-mer count per batch.

STRATEGY 2 : HLL-Balancing



Most batches have balance distinct kmers count

Nb of batch B = 24

Nb of genomes per batch varies but to a lesser extent compared to Strat 1

CONCLUSION & PERSPECTIVES

Batching by Predicting Compression Size Using HyperLogLog Distinct K-mer Estimation:

Improves balancing of the final compressed sizes.
Allows for better control over compression sizes.

Perspectives:

Extend the results to the whole collection.
Enable control over the number of genomes in each batch.

BIBLIOGRAPHY

- [1] Karel Brinda et al., 2024. Efficient and Robust Search of Microbial Genomes via Phylogenetic Compression. *Nature Methods*.
- [2] Grace A. Blackwell et al., 2021. Exploring bacterial diversity via a curated and searchable snapshot of archived DNA sequences. *PLOS Biology* 19, 11
- [3] Martin Hunt et al., 2024. AllTheBacteria - all bacterial genomes assembled, available and searchable. *bioRxiv*.
- [4] Jessica K. Bonnie et al., 2024. DandD: Efficient measurement of sequence growth and similarity. *iScience* 27, 3
- [5] Daniel N Baker, Ben Langmead, 2019. Dashing: Fast and Accurate Genomic Distances with HyperLogLog. *bioRxiv*
- [6] Mertens, Stephan, 2006, *The Easiest Hard Problem: Number Partitioning*, in Allon Percus; Gabriel Istrate; Cristopher Moore (eds.), *Computational complexity and statistical physics*, Oxford University Press US, p. 125.
- [7] Coffman et al., 2012. *Bin Packing Approximation Algorithms: Survey and Classification*. 10.1007/978-1-4419-7997-1_35.

Motivation

HLL-binning:

s.t.

Minimize $B = \sum_1^n b_j$, for $(j = 1, \dots, n)$

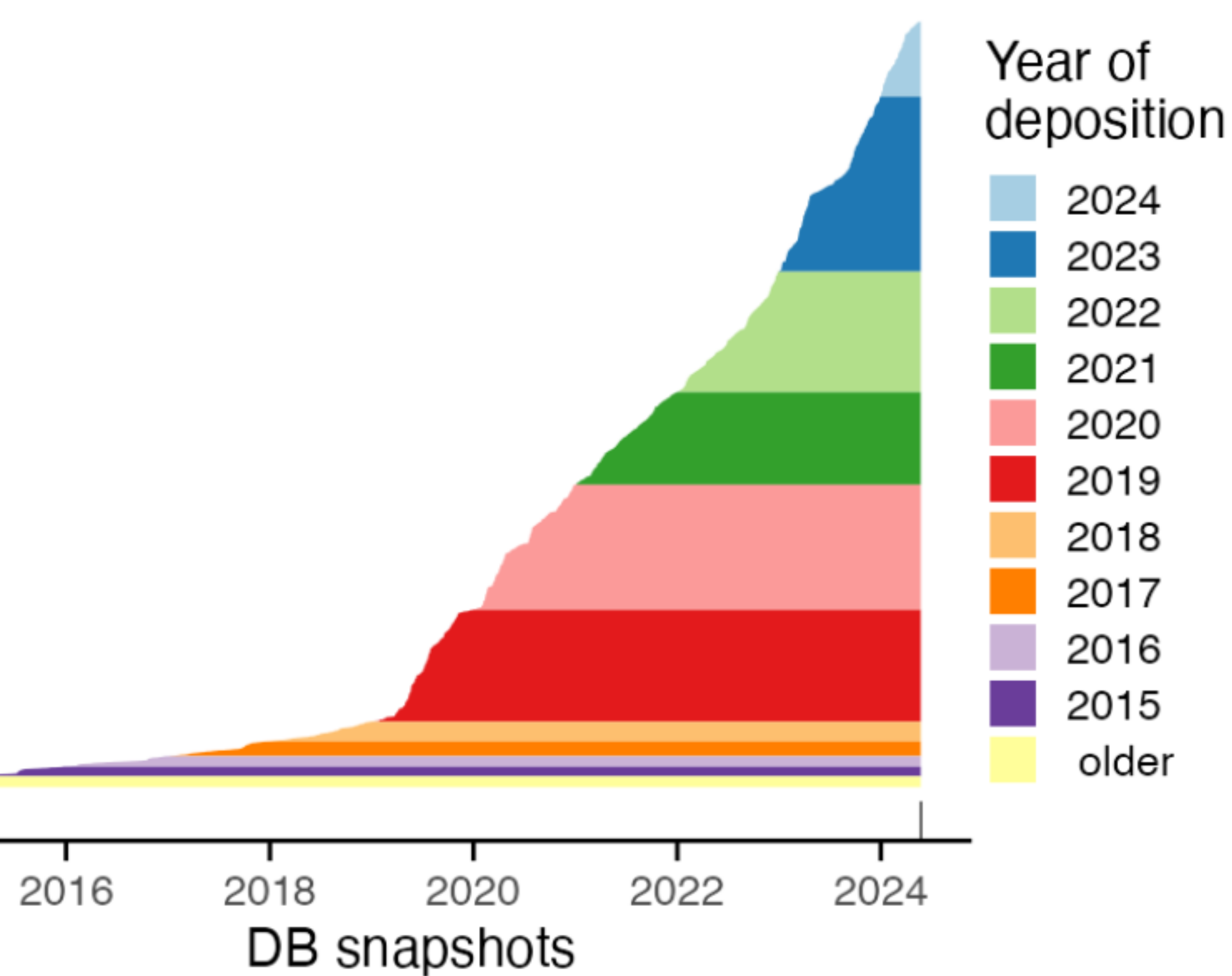
$|b_j| < C$, for $(j = 1, \dots, n)$

$\sum_{i=1}^m x_{ij} = 1$, for $(j = 1, \dots, n)$

$b_j \in \{0,1\}, x_{ij} \in \{0,1\}$

Collection of bacteria genomes is growing rapidly

st growth of bacterial genomes data



[Brinda et al.,bioRxiv, 2024]

Increasing Availability of Larger Bacterial Genome

661k (Blackwell et al., PLOS Biology, 2021) n = 6

AllTheBacteria (Hunt et al., bioRxiv, 2024) n = 2,4

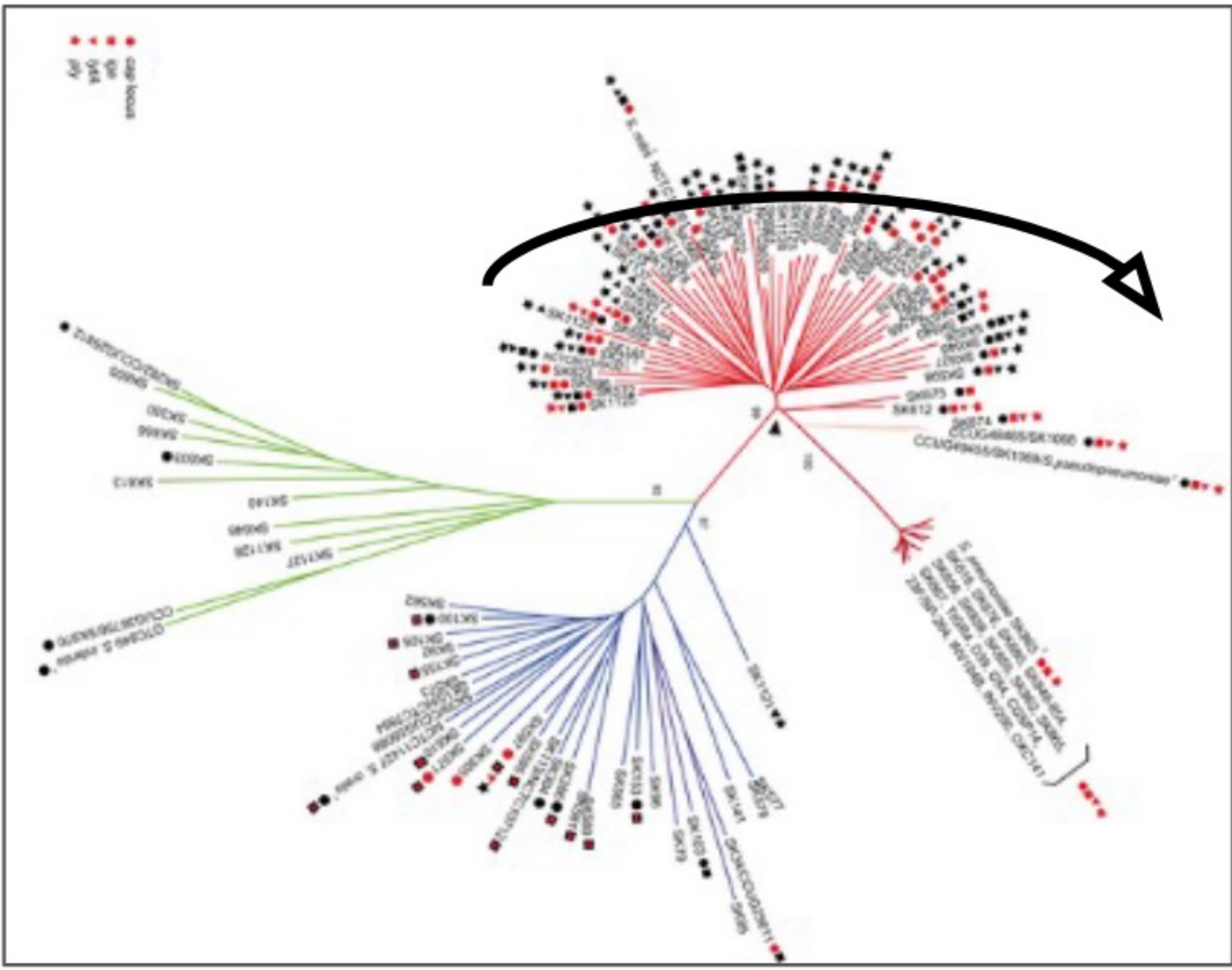
FUTURE : collection of ten of millions bacteria genomes,

ic compression allows efficient storage and query on bacteria

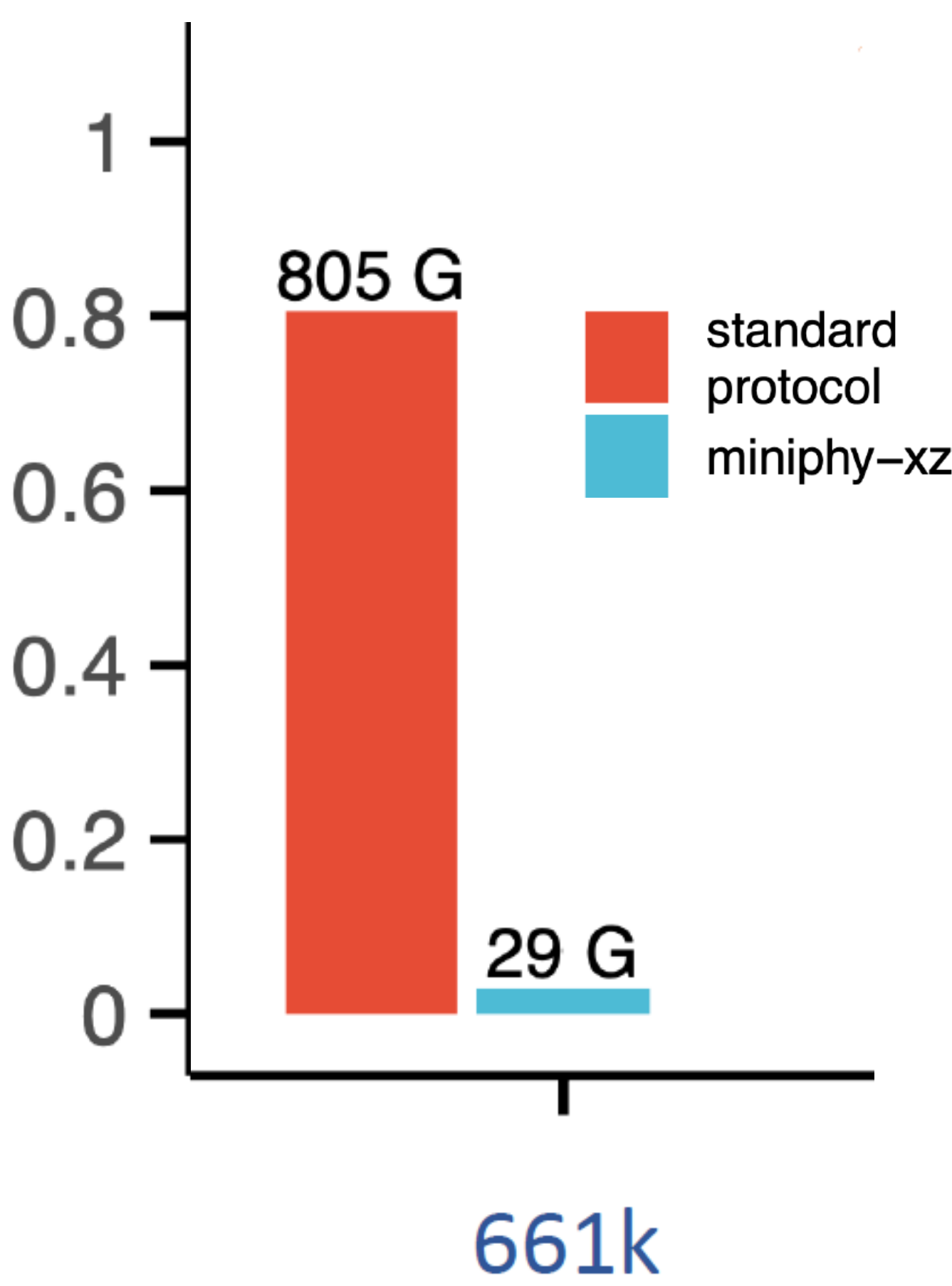
related genomes are extremely compressible.

mpression: Utilizes evolutionary relationships to
on and search

Genomes that are closer to each other on the phylogenetic tree are
more compressible



[Brinda et al.,bioRxiv, 2024]



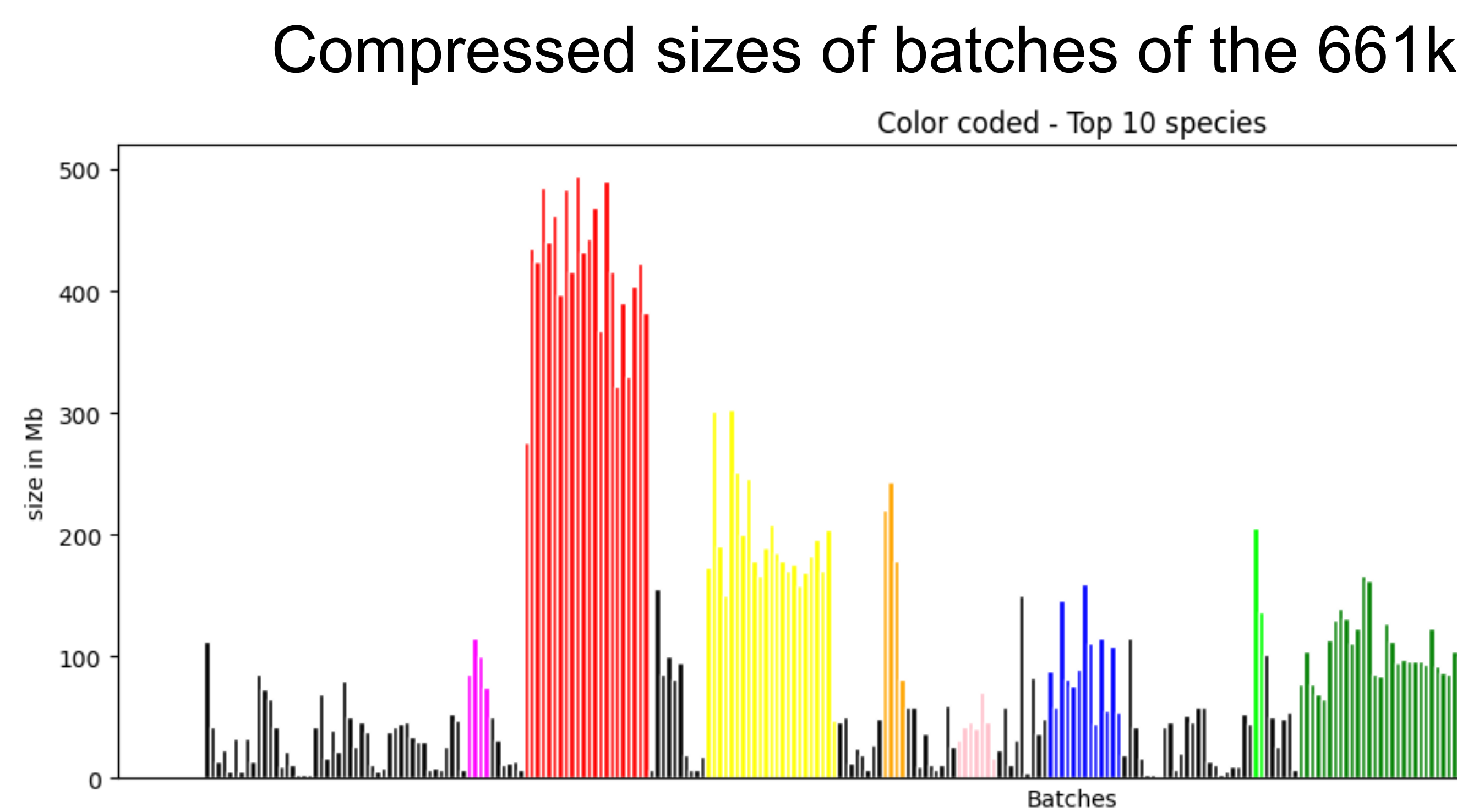
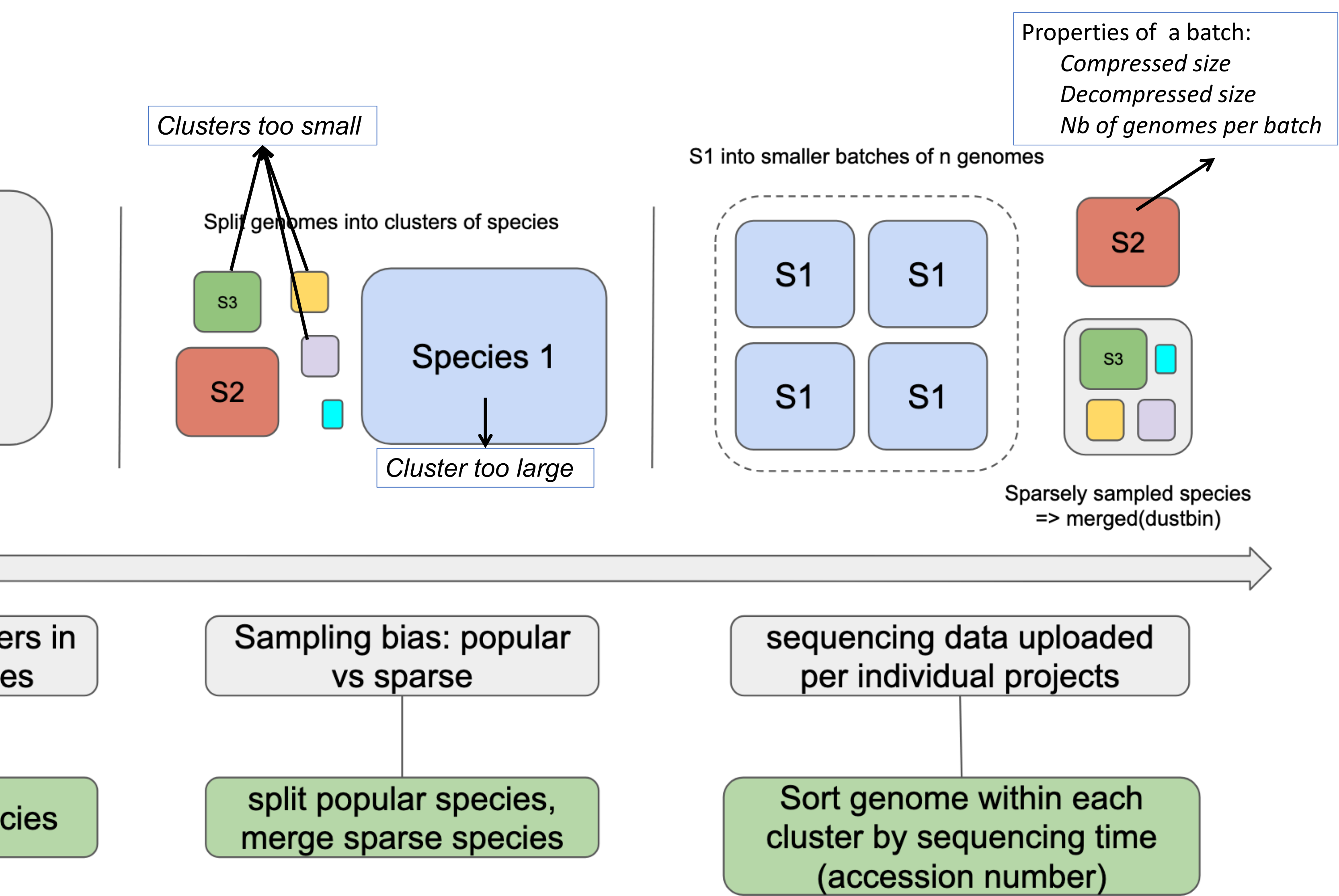
High compression ratio
Compressed using MiniPhy:
<https://github.com/karel-brinda/miniphy>

Phylinn: tool for search
compressed 661 colle
laptop.

<https://github.com/karel-brinda/phylinn>

State of the art

Key idea in phylogenetic compression: clustering/batching



This batching enables compression and search (parallelization) or linear time (using a single processor)

This batching is not adapted for different hardware memory constraint
The compressed size of the batch is not constant

Find batches that fit within a constraint (i.e. memory) while minimizing the number of batches
(maximize the number of genomes in the batches)

MultiProcessor Scheduling
ization Problem

Instance:

Items, a size $s(i) \in \mathbb{Z}^+$ for each $i \in I$.

Given b batches.

Partition the items in set I into b batches such that we minimize the max size of the batches

- ➔ These problem are extensively studied
- ➔ NP-hard but there are efficient heuristic algorithms
- ➔ How to adapt these problem for DNA data

Bin Packing Optimization Problem

Instance:

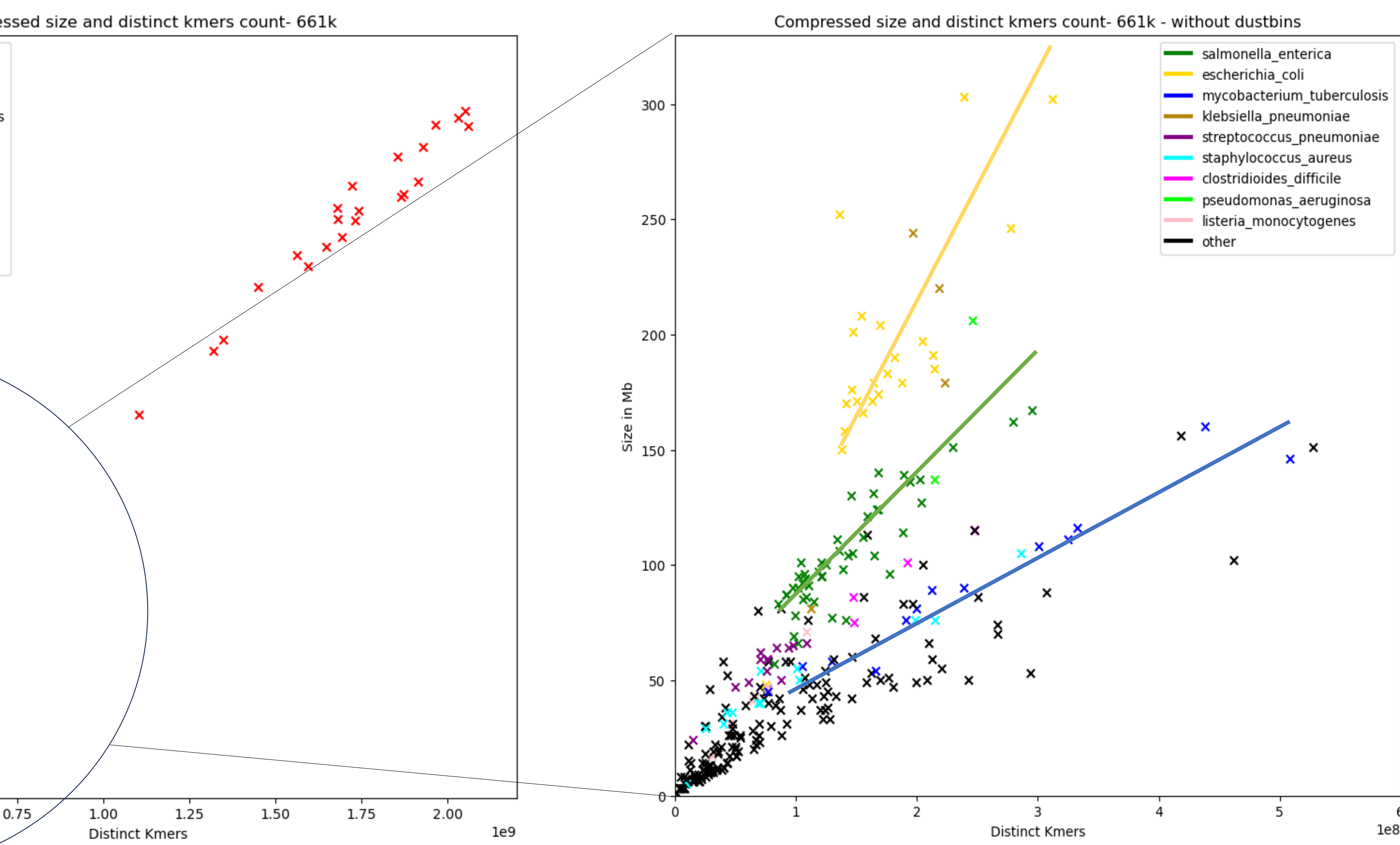
Finite set I of items, a size $s(i) \in \mathbb{Z}^+$ for each $i \in I$.

Given a capacity K for the bins

Objective: Partition the items into b batches and minimize the number of batches used

Method

Compressed size of batches can be estimated using biological properties



k-mers: substring of length k of DNA sequence

Correlation between kmer and compressed size of batches

Different correlation for different species

➔ Distinct kmers count can be used to estimate compressed size of batches

Accurate kmer estimation using probabilistic counting – HyperLogLog

Estimating an approximate, compact sketch of a dataset (e.g., a sketch)

GATCGATCGA

G

...

ATCG

TCGA

Probabilistic counting/approximate counting:

Sketches

Hash(item_1) = 0000... => Prob = $1 / 2^4$

Hash(item_2) = 0001...

Hash(item_3) = 0010...

Hash(item_4) = 0011...

Hash(item_5) = 0100...

Hash(item_6) = 0101...

...

Longest prefix of leading zeros in the hash values is 4

Estimated cardinality $\approx 2^4$

Set cardinality estimation with HyperLogLog

Hyperloglog is a sketching algorithm for advanced approximate counting

Baker, D.N., Langmead, B. improved tool dashing.

<https://github.com/dnbaker/dashing>

Result

Hyperloglog based bin packing and load balancing for genomes

HLL-Binning:
Put genomes into batches that fit within a user-defined capacity constraint

INPUT: GENOMES SKETCHES, CAPACITY (max_distinct_kmers)

Initialize:
Sort genomes based on their accession number.

H:

Assign:
For each genome, assign it to the first batch that can accommodate it (based on capacity).
If it fits, add it to the batch and move on to the next genome.

Overflow if Necessary:

If the current batch cannot accommodate the current genome, create a new batch and place the genome in it.

<https://github.com/tam-km-truong/HLL-Binning>

HLL-balancing:

Put genomes into balanced, user-defined number of batches

INPUT: GENOMES SKETCHES, NUMBER OF BATCH b

Initialize:

Sort the genomes based on their accession number.

Initial Assignment:

Assign the first b genomes directly to b batches.

Greedy Assignment:

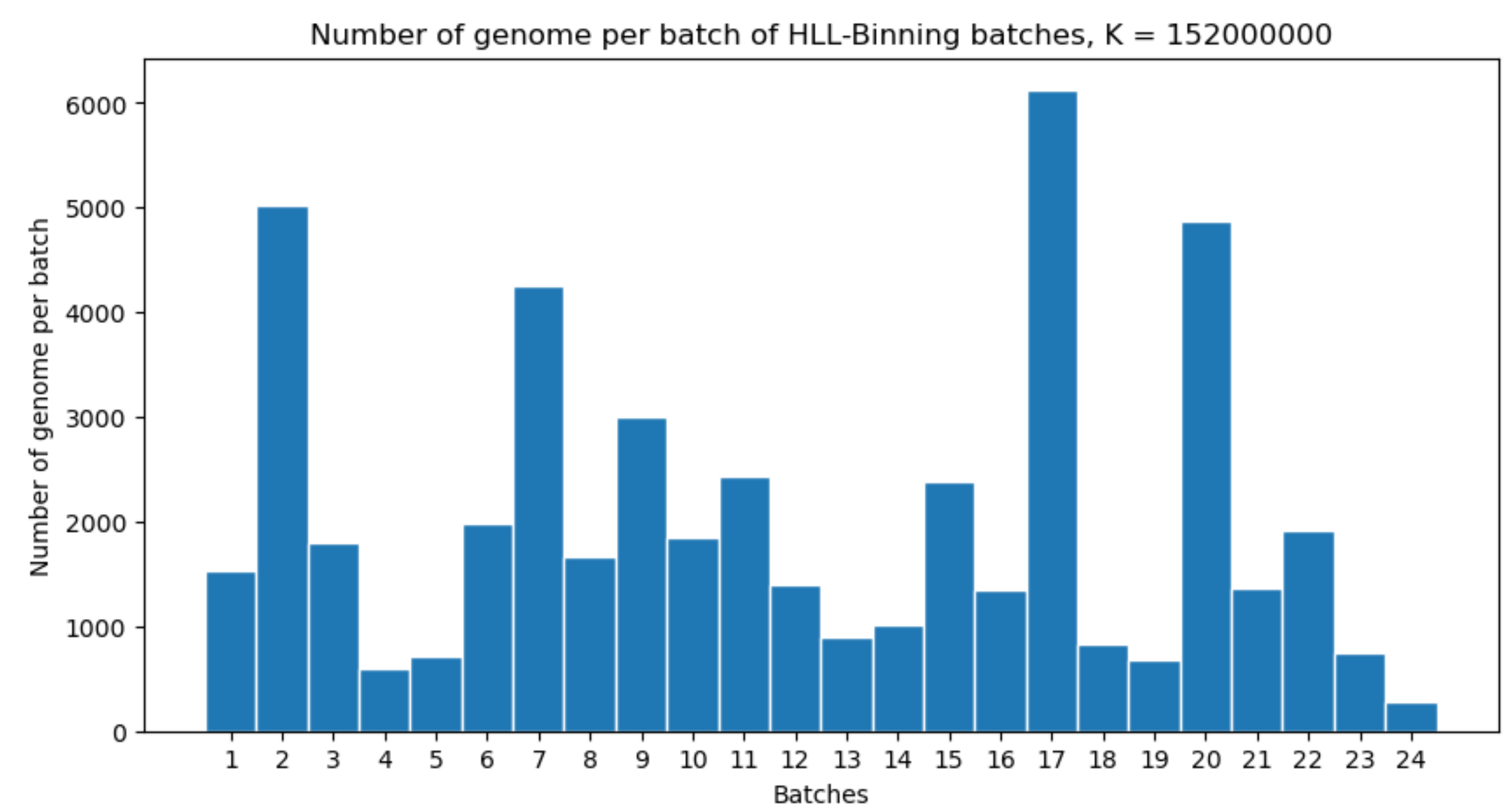
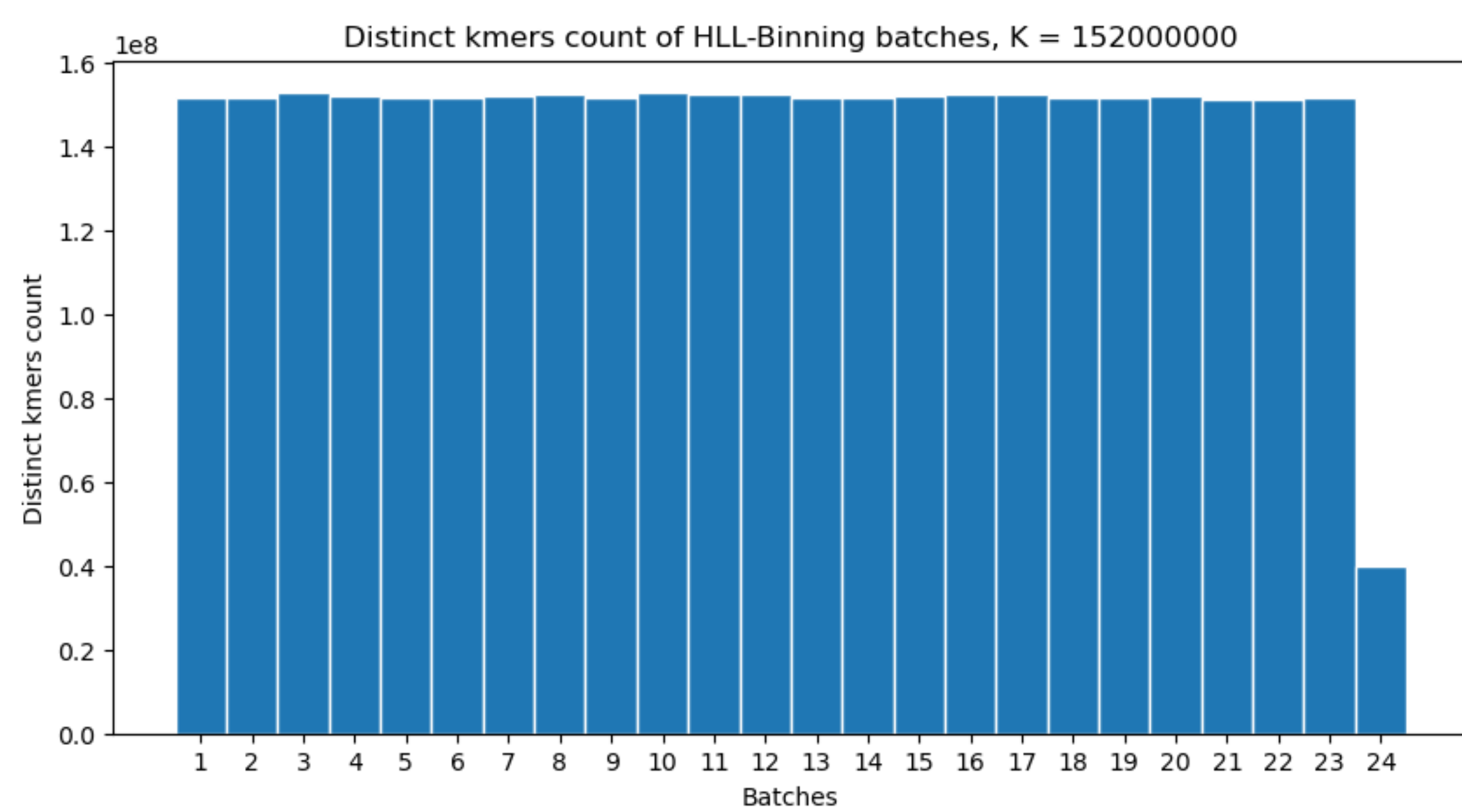
For each remaining genome, place it in the batch with the smallest distinct kmers count.

Update the batch's contents and size.

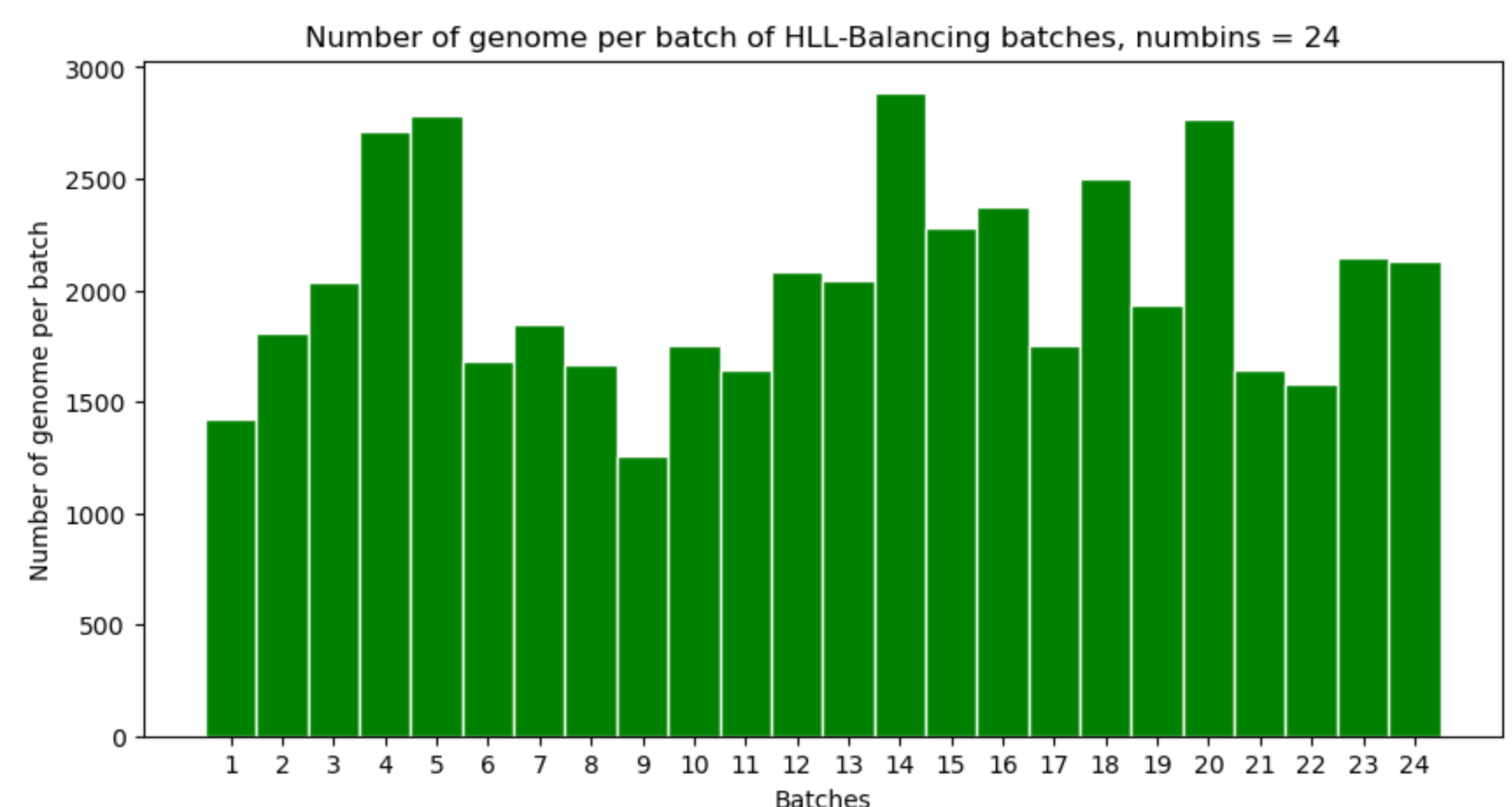
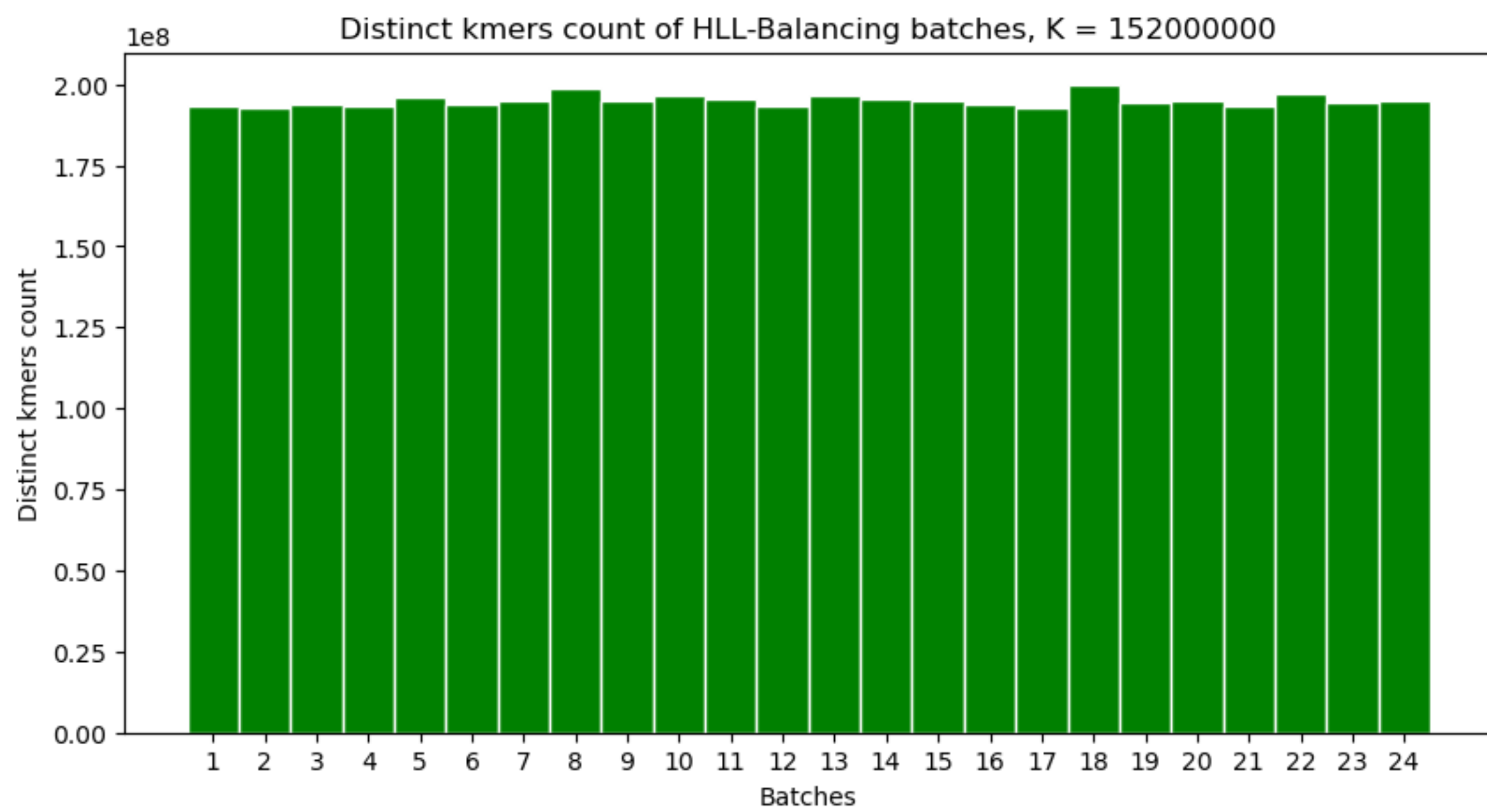
<https://github.com/tam-km-truong/HLL-Balancing>

EXPERIMENT: batching of Mycobacterium Tuberculosis

Capacity constraint for batches: 153000000 kmers so that the compressed size stays under 64MB (calculated by using the correlation of compressed size and kmers count)



Using the number of batches found by HLL-Binning, run HLL-Balancing with numbins b = 24



Even though the batches are balanced, HLL-Balancing batches have higher distinct kmers count than the capacity set in HLL-Binning

Discussion

