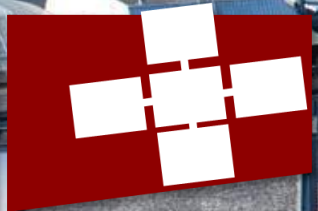
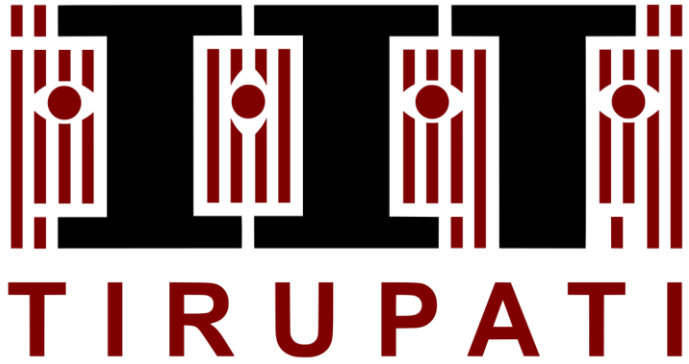


Communication-Efficient Jaccard Similarity for High-Performance Distributed Genome Comparisons

MACIEJ BESTA, RAGHAVENDRA KANAKAGIRI, HARUN MUSTAFA, MIKHAIL KARASIKOV, GUNNAR RÄTSCH, TORSTEN HOEFLER, EDGAR SOLOMONIK

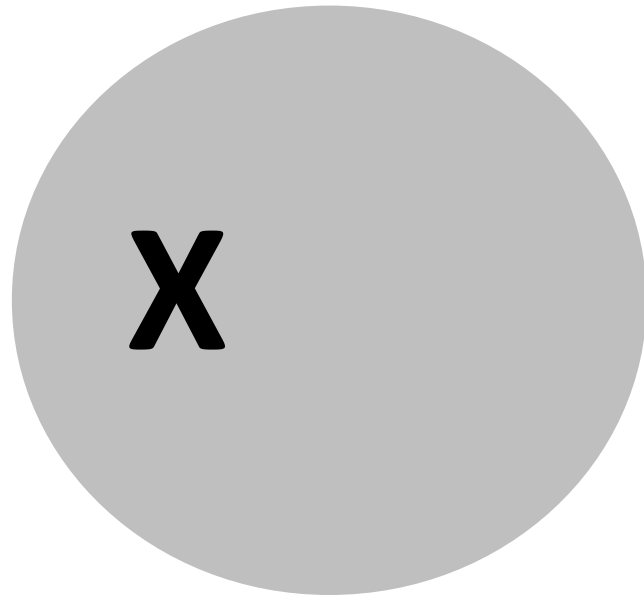
भारतीय प्रौद्योगिकी संस्थान तिरुपति



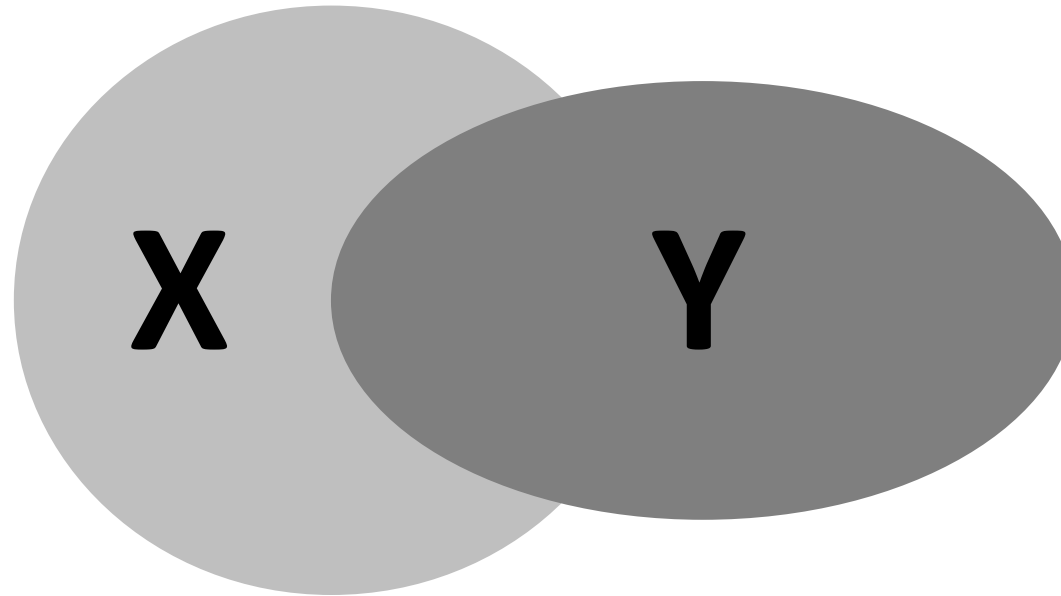
BIOMEDICAL
INFORMATICS

SET SIMILARITY

SET SIMILARITY

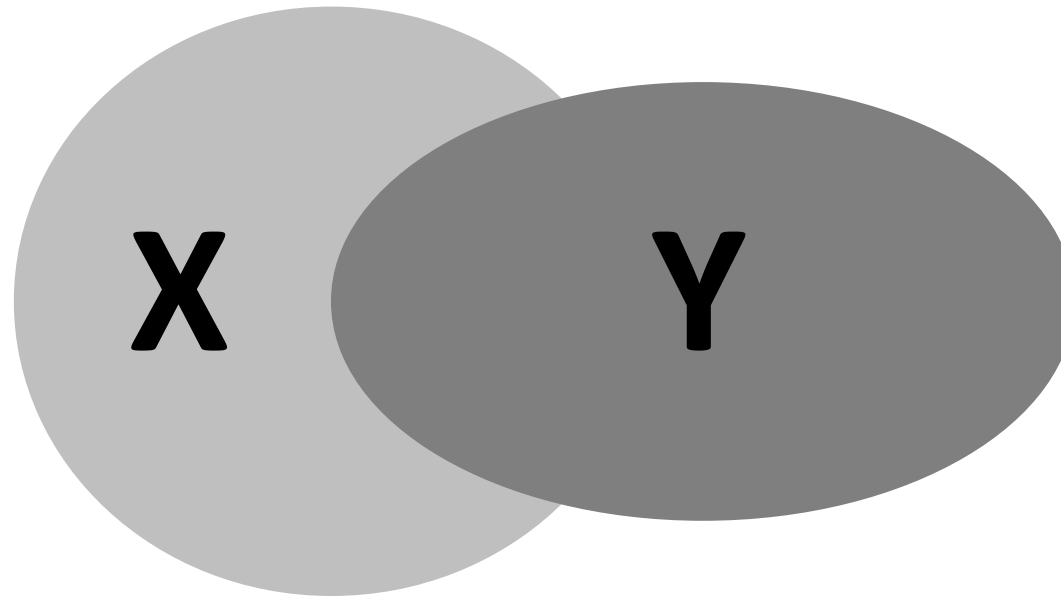


SET SIMILARITY



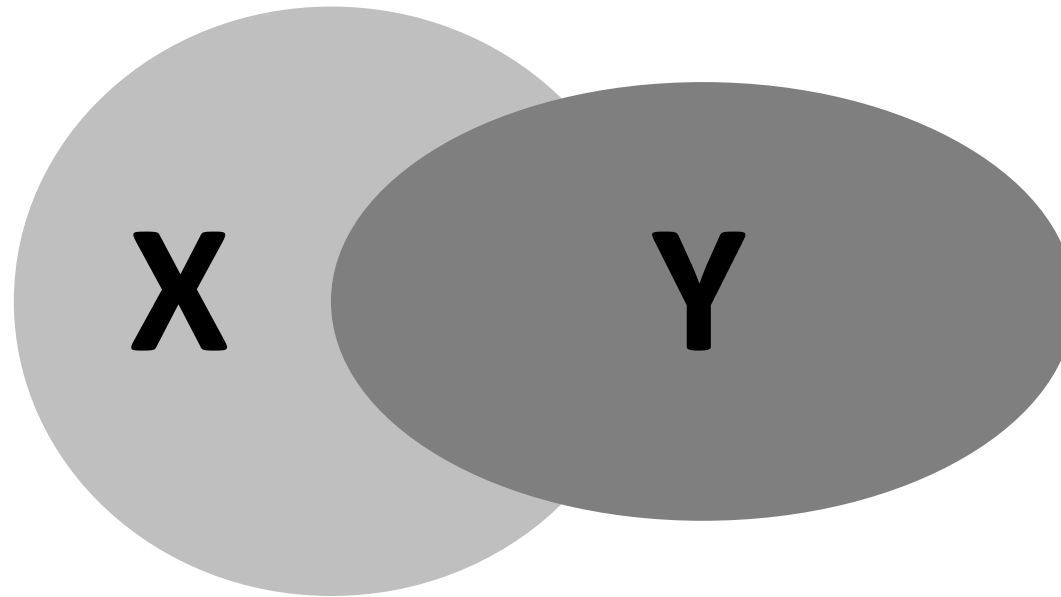
SET SIMILARITY

How can we measure the „similarity” of X and Y?



SET SIMILARITY

? How can we measure the „similarity“ of X and Y?



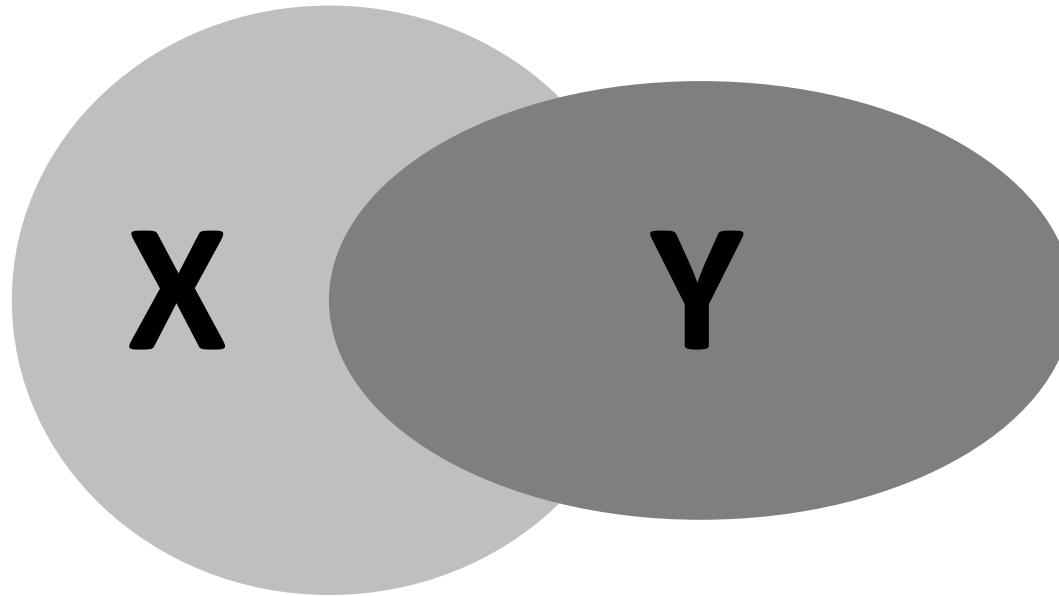
Jaccard Index:

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

SET SIMILARITY



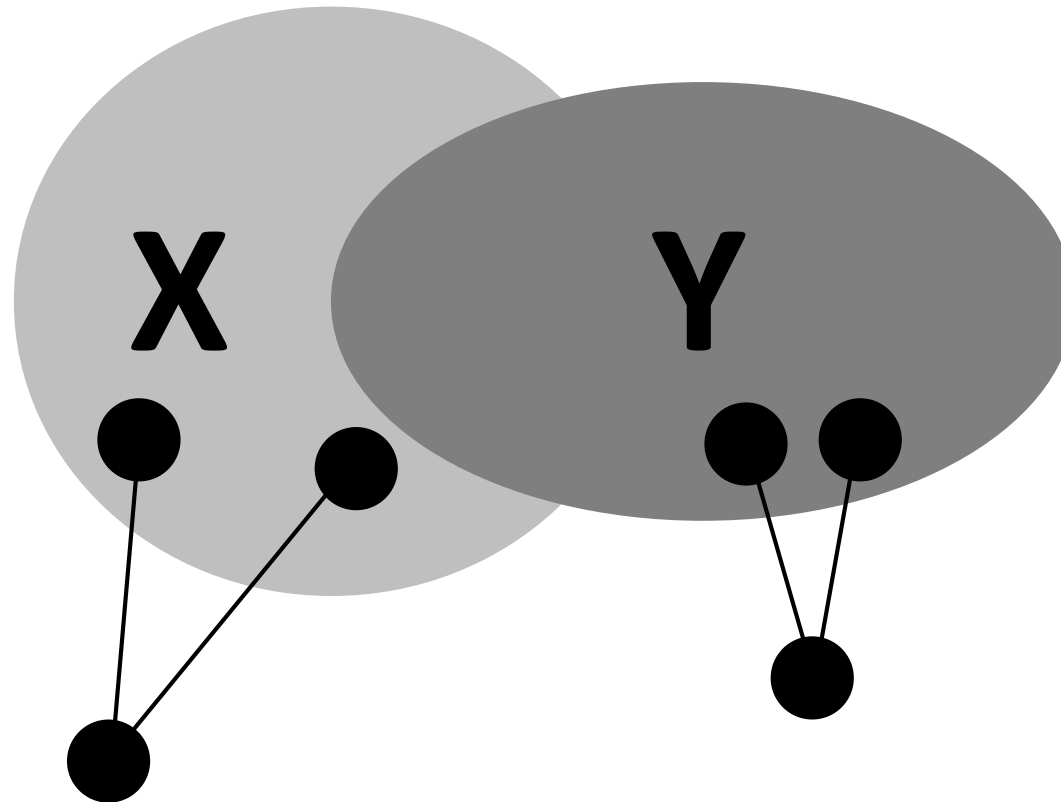
What are X and Y in practice?



SET SIMILARITY

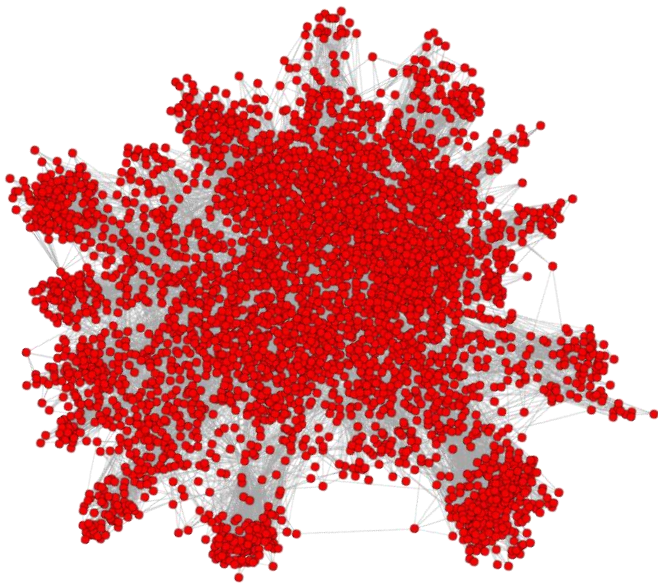


What are X and Y in practice?

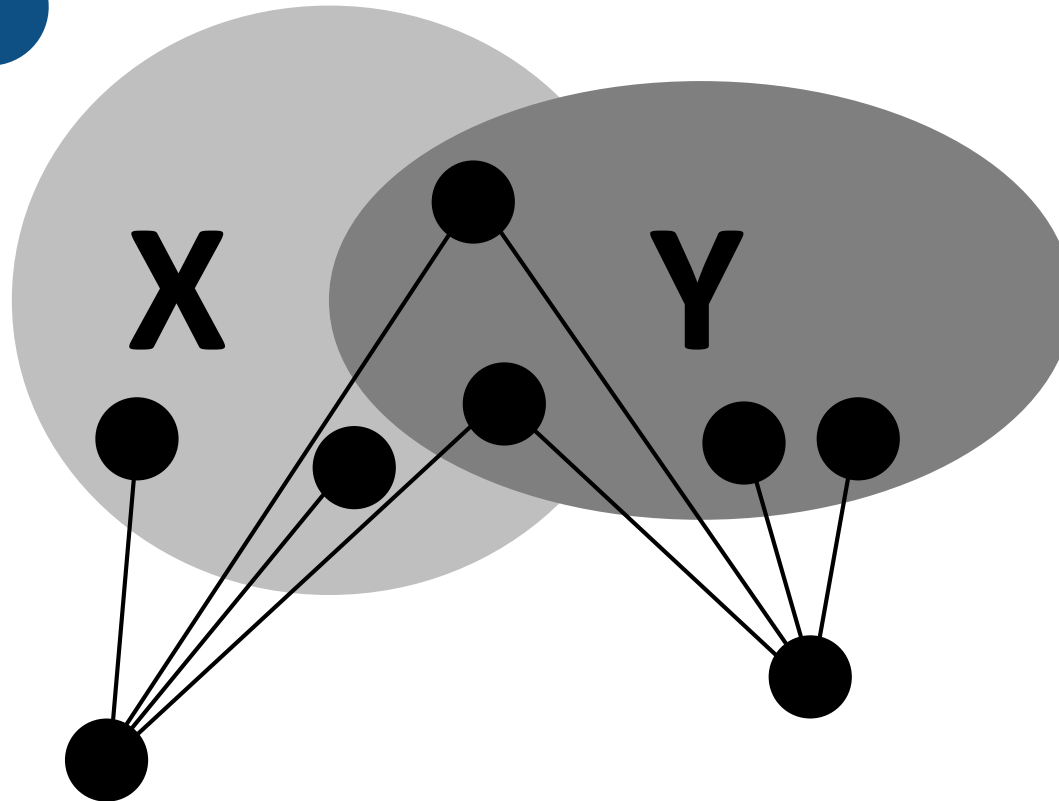


SET SIMILARITY

Neighborhoods of
vertices in a graph

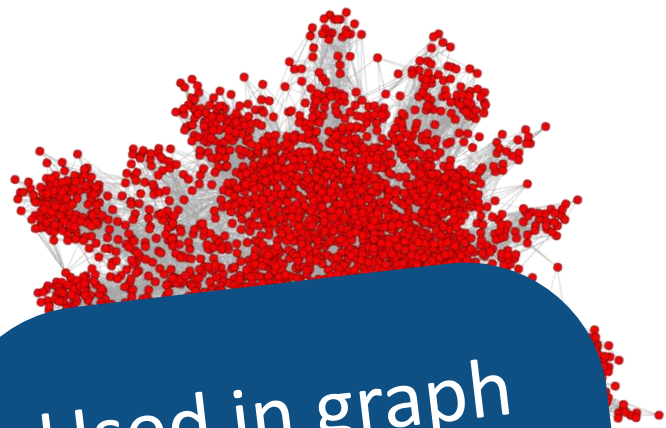


What are X and Y in
practice?



SET SIMILARITY

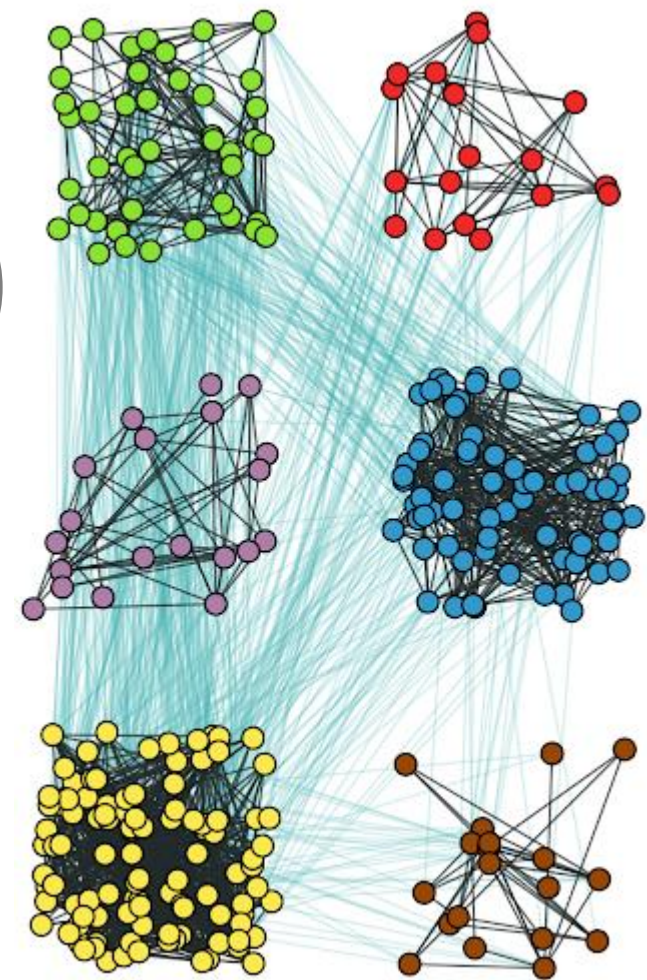
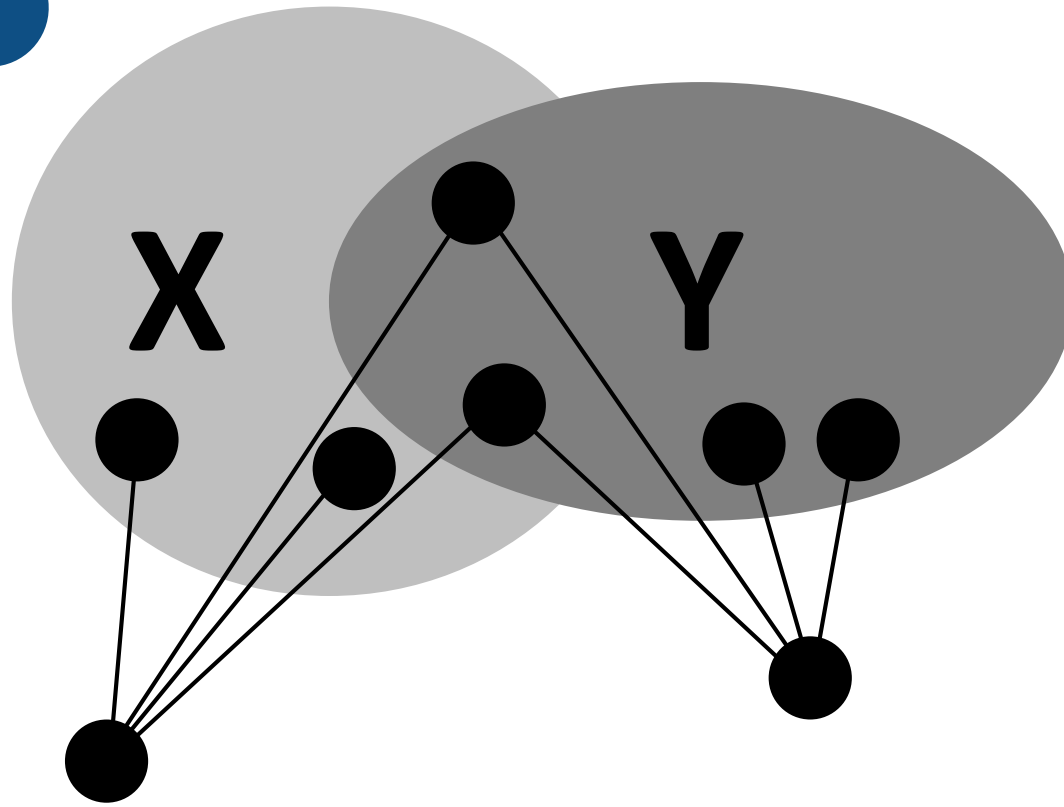
Neighborhoods of vertices in a graph



Used in graph analytics (clustering, link prediction, ...)



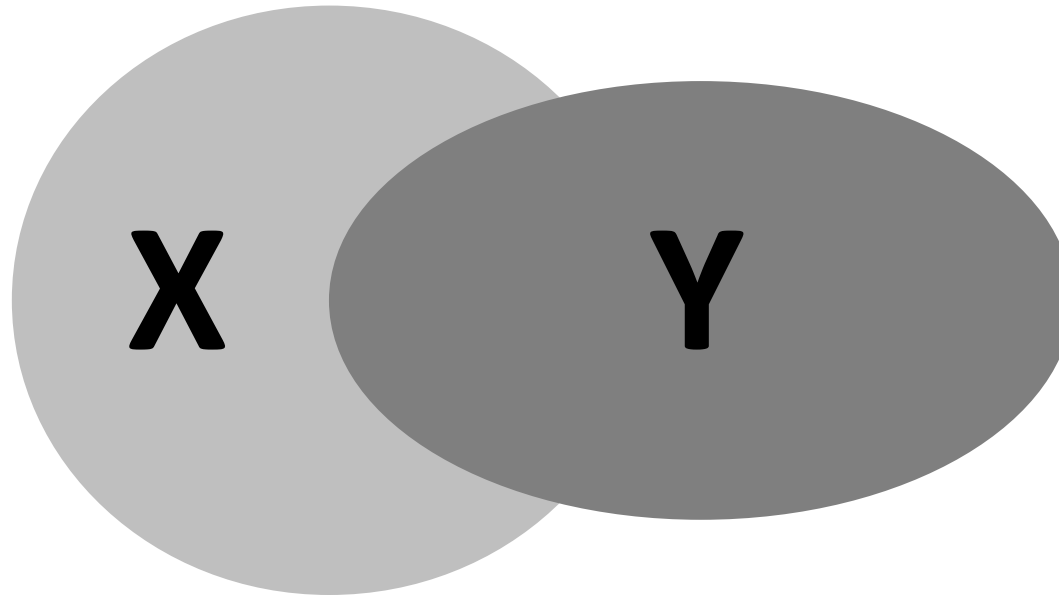
What are X and Y in practice?



SET SIMILARITY



What are X and Y in practice?



SET SIMILARITY

Detected object & ground-truth object

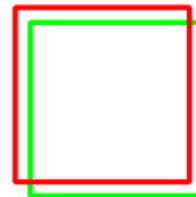
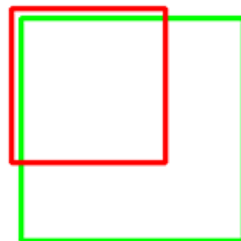
? What are X and Y in practice?



IoU: 0.4034

IoU: 0.7330

IoU: 0.9264



Poor

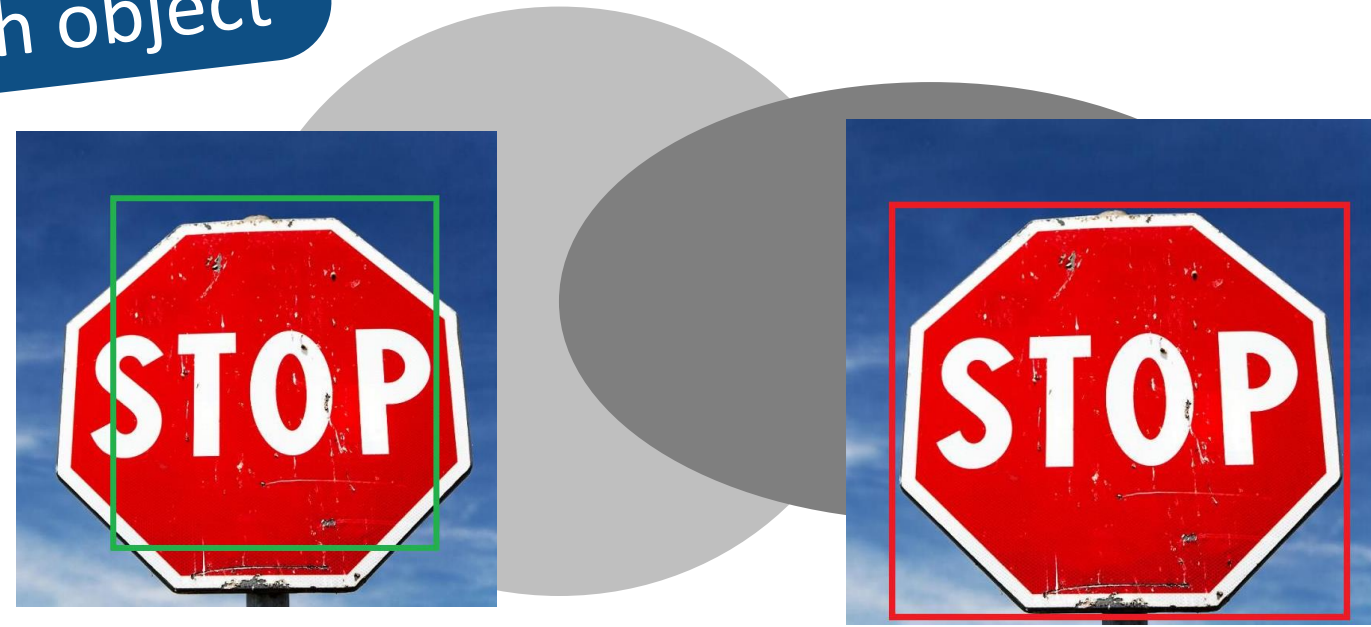
Good

Excellent

SET SIMILARITY

Detected object & ground-truth object

? What are X and Y in practice?

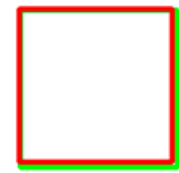
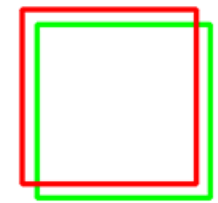
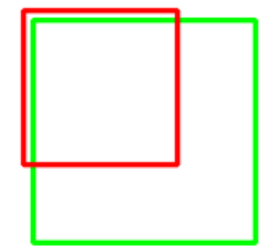


Used in image recognition

IoU: 0.4034

IoU: 0.7330

IoU: 0.9264

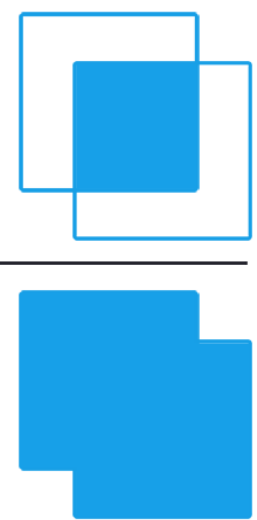


Poor

Good

Excellent

$$IoU = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$

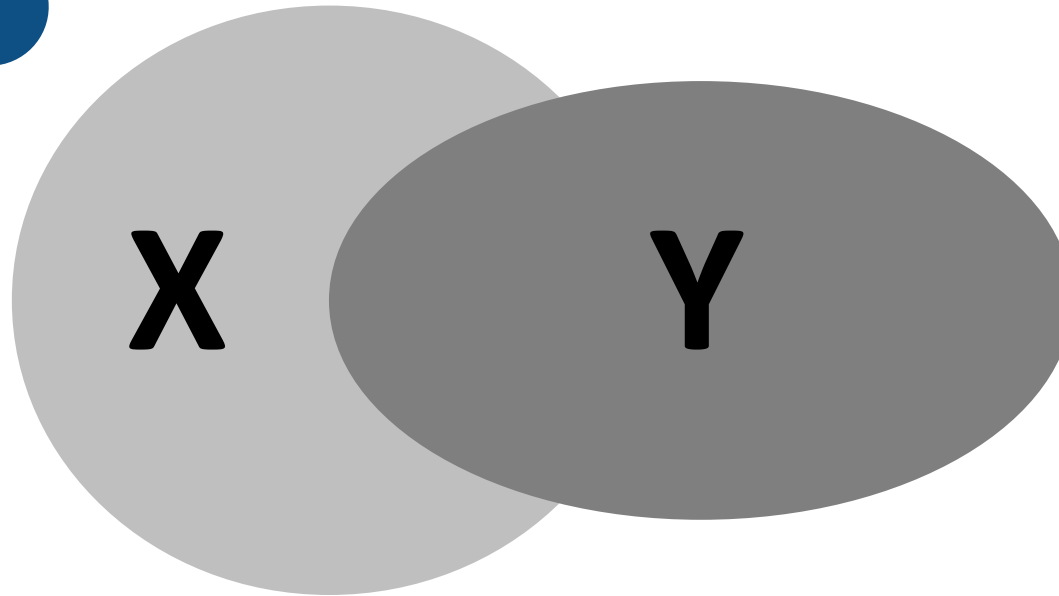


SET SIMILARITY

Sequences of
genomes



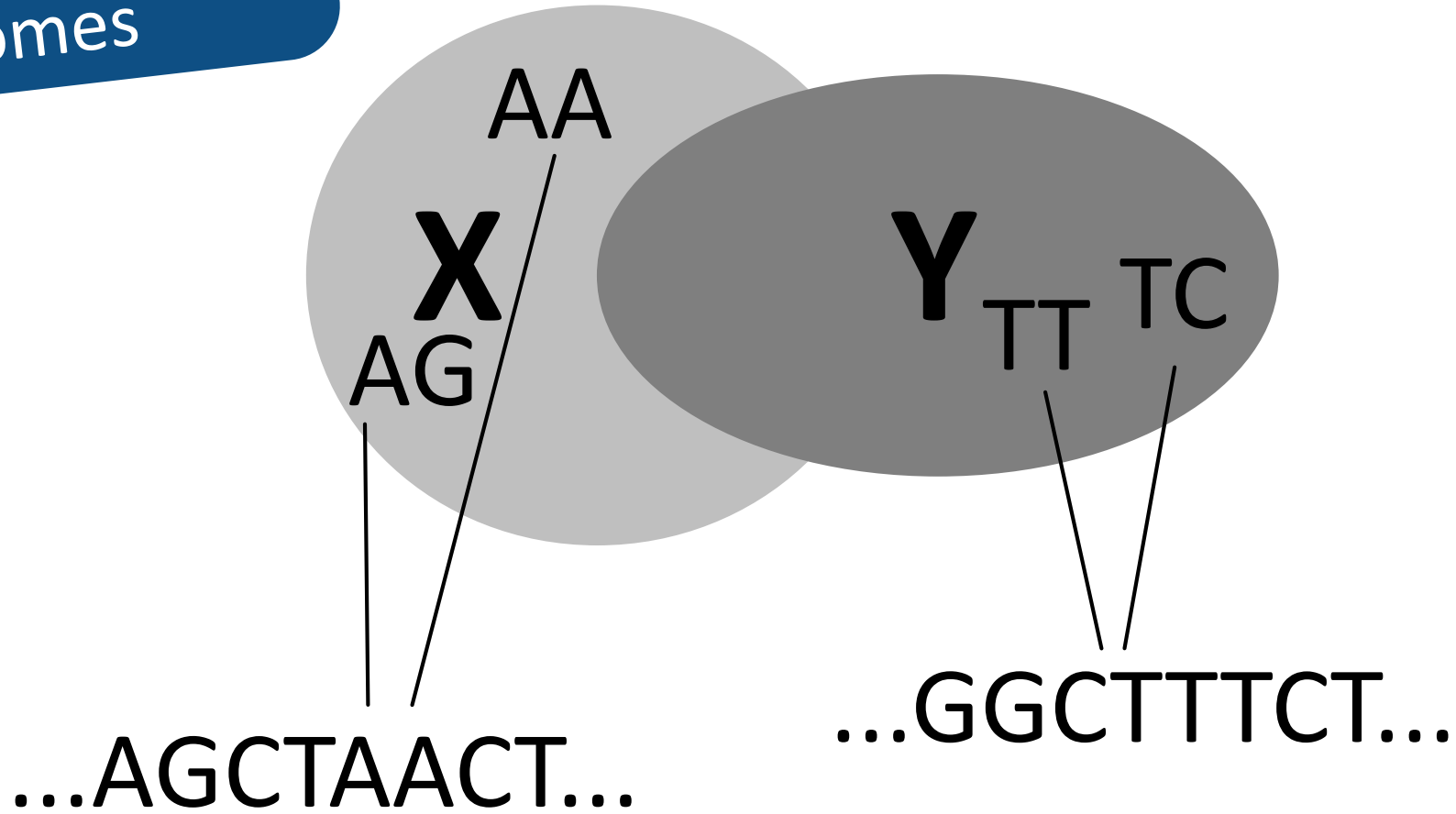
What are X and Y in
practice?



SET SIMILARITY

Sequences of genomes

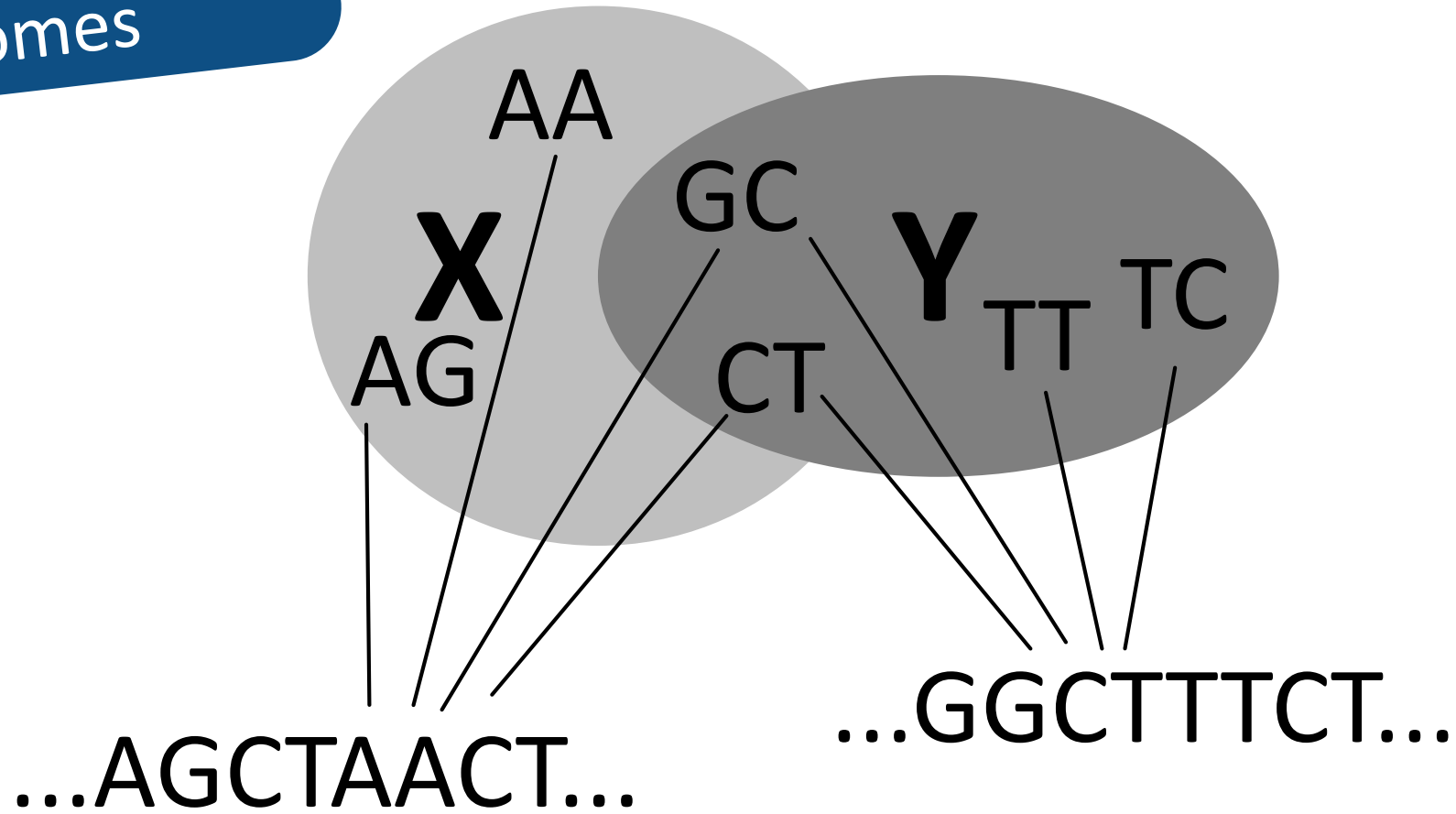
? What are X and Y in practice?



SET SIMILARITY

Sequences of genomes

? What are X and Y in practice?



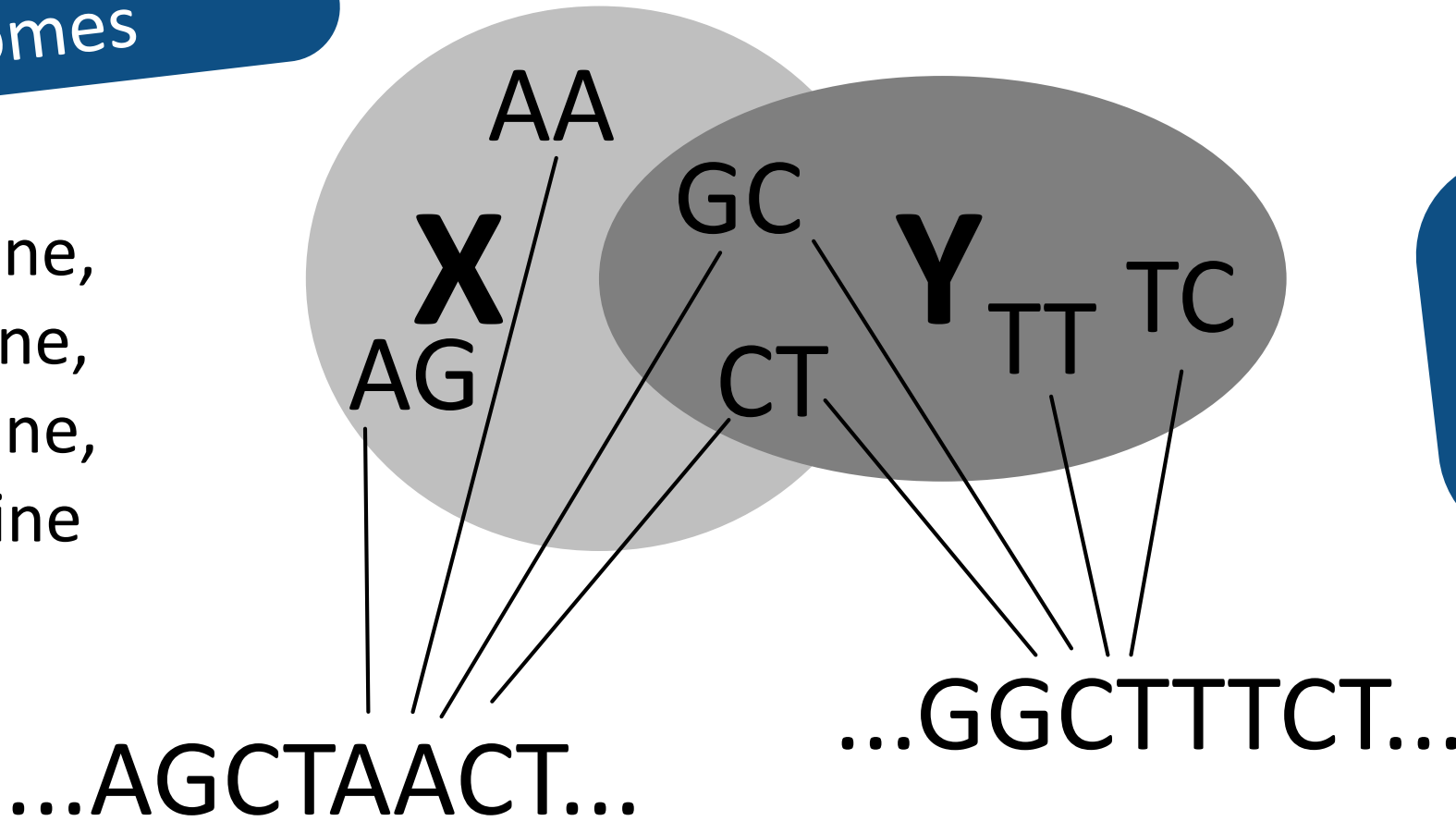
SET SIMILARITY

Sequences of genomes



What are X and Y in practice?

A: adenine,
T: thymine,
G: guanine,
C: cytosine



Used in genome assembly

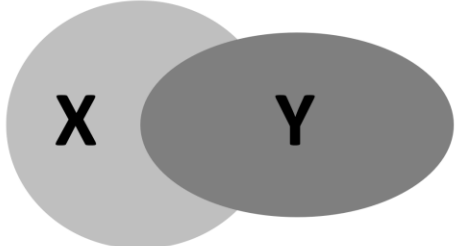
SCOPE OF WORK

SCOPE OF WORK

SPCL spl.inf.ethz.ch @spl_eth ETH zürich

SET SIMILARITY

How can we measure the „similarity“ of X and Y?



Jaccard Index:

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

Navigation icons: back, forward, search, refresh, close

Part 1 „SimilarityAtScale“: the first communication-efficient distributed algorithm to compute the general Jaccard similarity index and distance

SCOPE OF WORK

SET SIMILARITY

How can we measure the „similarity“ of X and Y?

Jaccard Index:

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

Part 1 „SimilarityAtScale“: the first communication-efficient distributed algorithm to compute the general Jaccard similarity index and distance

Part 2 „GenomeAtScale“: the first tool for fast, scalable, accurate, and large-scale derivations of Jaccard index between genome sequences

SET SIMILARITY

Sequences of genomes

What are X and Y in practice?

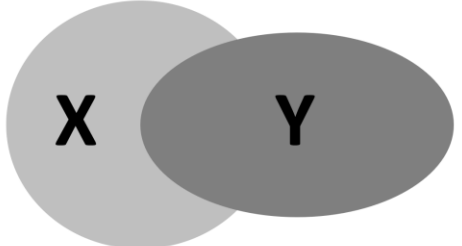
A: adenine,
T: thymine,
G: guanine,
C: cytosine

...AGCTAACT... ...GGCTTTCT...

SPCL spcl.inf.ethz.ch @spcl_eth ETH zürich

SET SIMILARITY

How can we measure the „similarity“ of X and Y?



Jaccard Index:

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

Navigation icons: back, forward, search, refresh, close

Part 1 „SimilarityAtScale“: the first communication-efficient distributed algorithm to compute the general Jaccard similarity index and distance

SIMILARITY AT SCALE

SIMILARITY AT SCALE

(1) Sets X_i $i = 1, \dots, n$

SIMILARITY AT SCALE

(1) Sets X_i $i = 1, \dots, n$

(2) Indicator matrix

$$\mathbf{A} \in \mathbb{B}^{m \times n} \quad a_{ij} = \begin{cases} 1 & : i \in X_j \\ 0 & : \text{otherwise} \end{cases}$$

SIMILARITY AT SCALE

(1) Sets X_i $i = 1, \dots, n$

(2) Indicator matrix

$$\mathbf{A} \in \mathbb{B}^{m \times n} \quad a_{ij} = \begin{cases} 1 & : i \in X_j \\ 0 & : \text{otherwise} \end{cases}$$

Describes which element (out of all m possible ones) belongs to which set (out of all n possible ones)

SIMILARITY AT SCALE

(1) Sets X_i $i = 1, \dots, n$

Describes which element (out of all m possible ones) belongs to which set (out of all n possible ones)

(2) Indicator matrix

$$\mathbf{A} \in \mathbb{B}^{m \times n} \quad a_{ij} = \begin{cases} 1 & : i \in X_j \\ 0 & : \text{otherwise} \end{cases}$$

(3) Similarity matrix

$$\mathbf{S} \in \mathbb{R}^{n \times n}$$

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

SIMILARITY AT SCALE

(1) Sets X_i $i = 1, \dots, n$

Describes which element (out of all m possible ones) belongs to which set (out of all n possible ones)

(2) Indicator matrix

$$\mathbf{A} \in \mathbb{B}^{m \times n} \quad a_{ij} = \begin{cases} 1 & : i \in X_j \\ 0 & : \text{otherwise} \end{cases}$$

(3) Similarity matrix

$$\mathbf{S} \in \mathbb{R}^{n \times n}$$

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

Describes Jaccard similarity between each set

SIMILARITYATSCALE: COMPUTE SIMILARITY MATRIX $S \in \mathbb{R}^{n \times n}$

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

SIMILARITYATSCALE: COMPUTE SIMILARITY MATRIX $S \in \mathbb{R}^{n \times n}$

Cardinalities of intersections
of each pair of data samples

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

SIMILARITYATSCALE: COMPUTE SIMILARITY MATRIX $S \in \mathbb{R}^{n \times n}$

Cardinalities of intersections
of each pair of data samples

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

Cardinalities of unions of
each pair of data samples

SIMILARITYATSCALE: COMPUTE SIMILARITY MATRIX $S \in \mathbb{R}^{n \times n}$

Cardinalities of intersections
of each pair of data samples

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|} \rightarrow B \in \mathbb{N}^{n \times n}$$

Cardinalities of unions of
each pair of data samples

SIMILARITYATSCALE: COMPUTE SIMILARITY MATRIX $S \in \mathbb{R}^{n \times n}$

Cardinalities of intersections
of each pair of data samples

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

$B \in \mathbb{N}^{n \times n}$

$C \in \mathbb{N}^{n \times n}$

Cardinalities of unions of
each pair of data samples

SIMILARITYATSCALE: COMPUTE SIMILARITY MATRIX $S \in \mathbb{R}^{n \times n}$

Cardinalities of intersections of each pair of data samples

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

$B \in \mathbb{N}^{n \times n}$

Can be obtained from B with simple set theory

$C \in \mathbb{N}^{n \times n}$

Cardinalities of unions of each pair of data samples

$$|X_i \cup X_j| = |X_i| + |X_j| - |X_i \cap X_j|$$

SIMILARITYATSCALE: COMPUTE SIMILARITY MATRIX $S \in \mathbb{R}^{n \times n}$

Cardinalities of intersections of each pair of data samples

Critical element!

$$s_{ij} = J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

$B \in \mathbb{N}^{n \times n}$

Can be obtained from B with simple set theory

$C \in \mathbb{N}^{n \times n}$

Cardinalities of unions of each pair of data samples

$$|X_i \cup X_j| = |X_i| + |X_j| - |X_i \cap X_j|$$

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B \in \mathbb{N}^{n \times n}$

$$b_{ij} = |X_i \cap X_j| = \sum_k a_{ki} a_{kj}, \text{ so } B = A^T A$$

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B \in \mathbb{N}^{n \times n}$

$$b_{ij} = |X_i \cap X_j| = \sum_k a_{ki} a_{kj}, \text{ so } B = A^T A$$

Sparse-matrix product

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B \in \mathbb{N}^{n \times n}$

$$b_{ij} = |X_i \cap X_j| = \sum_k a_{ki} a_{kj}, \text{ so } B = A^T A$$

Largest computational challenge: for the considered workloads, the indicator matrix A is incredibly sparse

Sparse-matrix product

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B \in \mathbb{N}^{n \times n}$

$$b_{ij} = |X_i \cap X_j| = \sum_k a_{ki} a_{kj}, \text{ so } B = A^T A$$

Largest computational challenge: for the considered workloads, the indicator matrix A is incredibly sparse

Sparse-matrix product

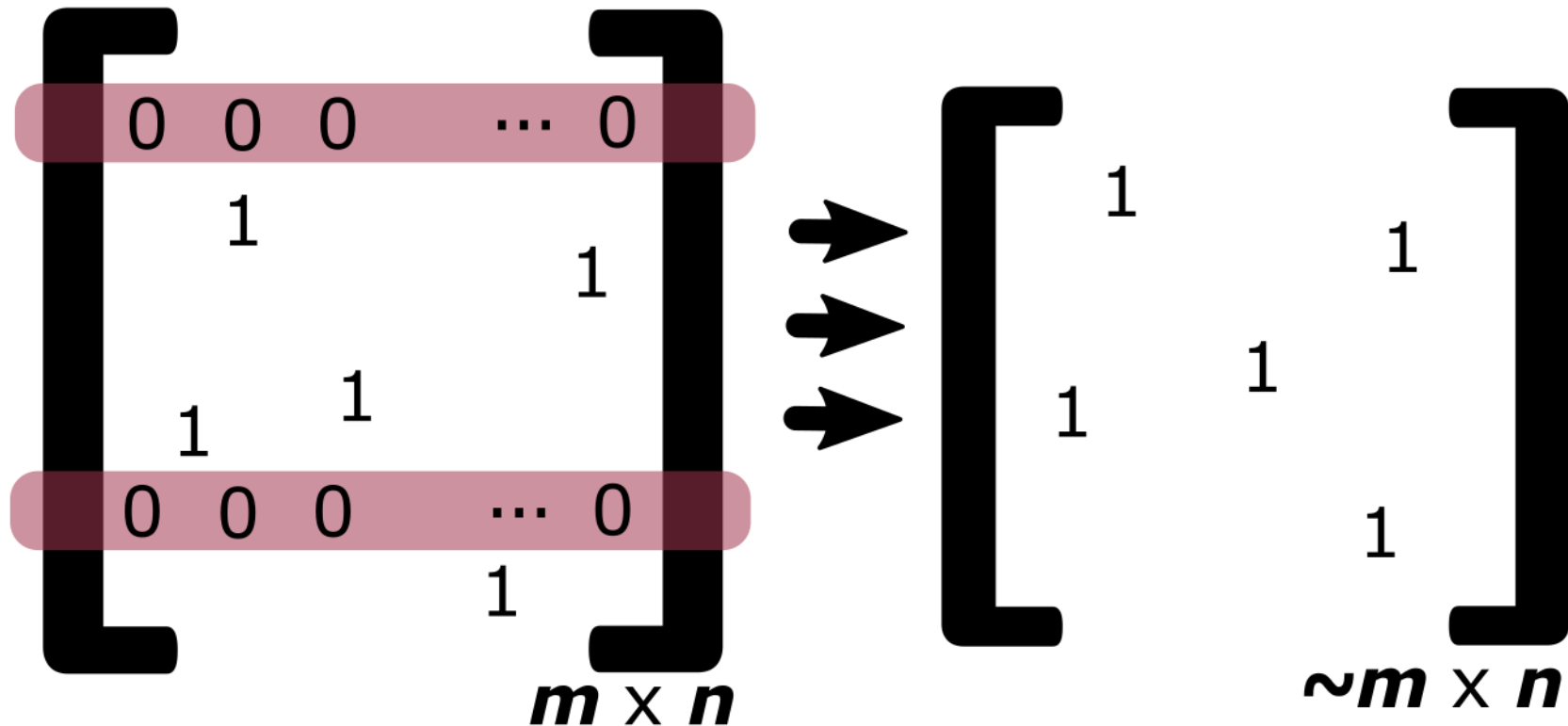
We use three techniques to tackle it

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Technique 1: remove zero rows from A

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Technique 1: remove zero rows from A

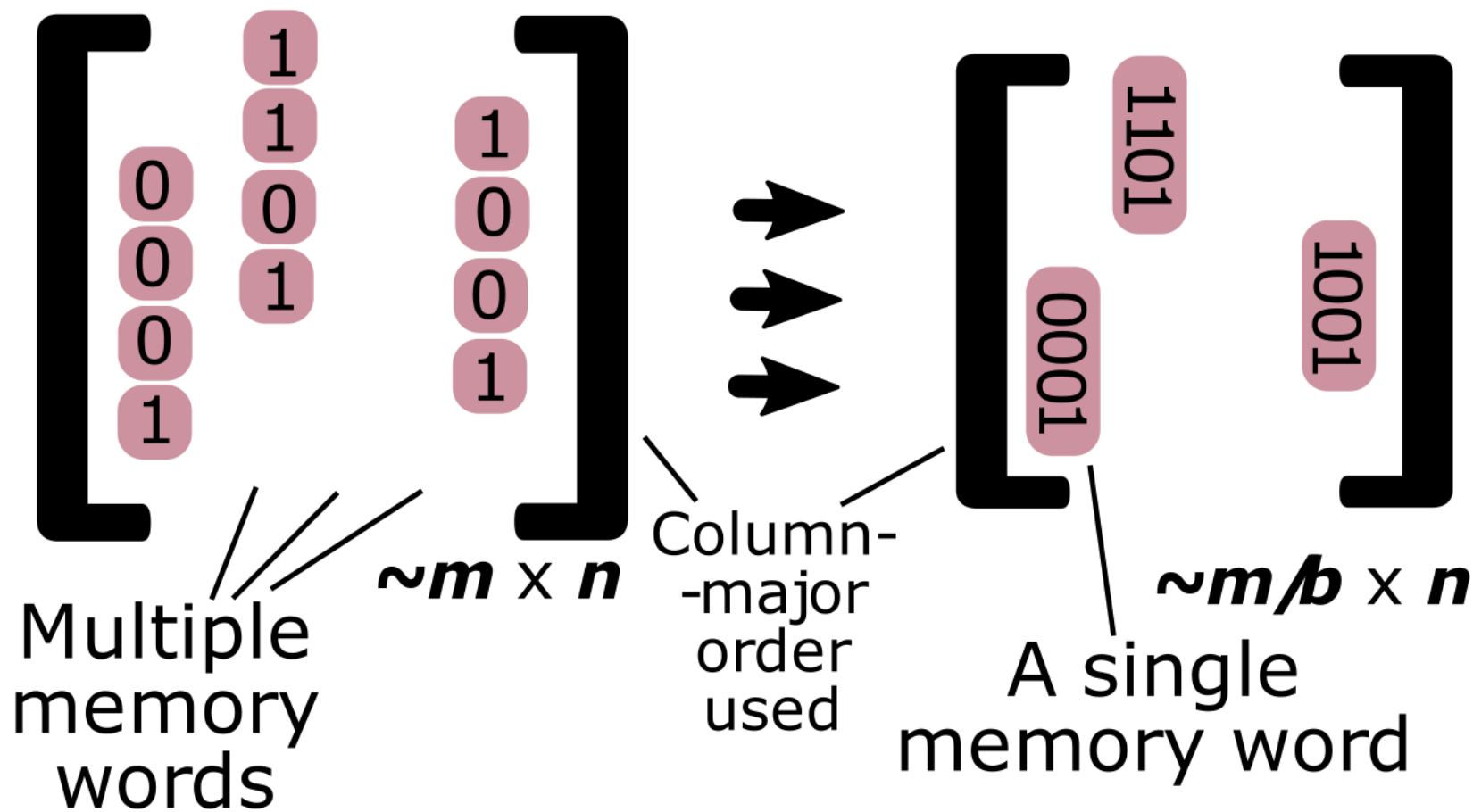


SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Technique 2: Compress (mask row segments into bit vectors)

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Technique 2: Compress (mask row segments into bit vectors)

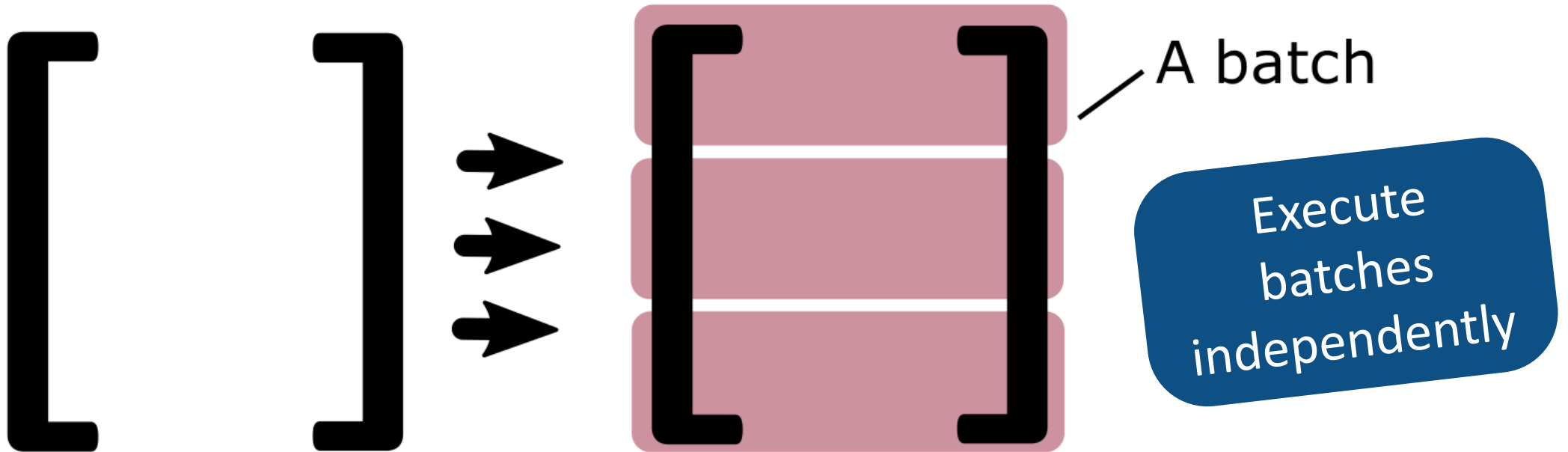


SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Technique 3: Divide into batches to further alleviate the large size of A

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Technique 3: Divide into batches to further alleviate the large size of A



SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Full algebraic formulations
of all the techniques

$$A = \begin{bmatrix} A^{(1)} \\ \vdots \\ A^{(r)} \end{bmatrix}, \text{ where } A^{(l)} \in \mathbb{B}^{\tilde{m} \times n}, \forall l \in \{1, \dots, r\}.$$

$$B = \sum_{l=1}^r A^{(l)T} A^{(l)}, \quad \hat{a} = \sum_{l=1}^r \hat{a}^{(l)}, \quad \text{where } \hat{a}_i^{(l)} = \sum_k a_{ki}^{(l)}$$

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Full algebraic formulations
of all the techniques

$$A = \begin{bmatrix} A^{(1)} \\ \vdots \\ A^{(r)} \end{bmatrix}, \text{ where } A^{(l)} \in \mathbb{B}^{\tilde{m} \times n}, \forall l \in \{1, \dots, r\}.$$

Theoretical analysis of
communication cost

$$B = \sum_{l=1}^r A^{(l)T} A^{(l)}, \quad \hat{a} = \sum_{l=1}^r \hat{a}^{(l)}, \text{ where } \hat{a}_i^{(l)} = \sum_k a_{ki}^{(l)}$$

$$O\left(\left(1 + \frac{z}{M\sqrt{cp}}\right) \cdot \alpha + \left(\frac{z}{\sqrt{cp}} + \frac{cn^2}{p}\right) \cdot \beta\right)$$

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Full algebraic formulations
of all the techniques

$$A = \begin{bmatrix} A^{(1)} \\ \vdots \\ A^{(r)} \end{bmatrix}, \text{ where } A^{(l)} \in \mathbb{B}^{\tilde{m} \times n}, \forall l \in \{1, \dots, r\}.$$

Theoretical analysis of
communication cost

$$B = \sum_{l=1}^r A^{(l)T} A^{(l)}, \quad \hat{a} = \sum_{l=1}^r \hat{a}^{(l)}, \text{ where } \hat{a}_i^{(l)} = \sum_k a_{ki}^{(l)}$$

$$O\left(\left(1 + \frac{z}{M\sqrt{cp}}\right) \cdot \alpha + \left(\frac{z}{\sqrt{cp}} + \frac{cn^2}{p}\right) \cdot \beta\right)$$

Public implementation
based on CTF [1]

SIMILARITYATSCALE: COMPUTE INTERSECTION MATRIX $B = A^T A$

Full algebraic formulations
of all the techniques

$$A = \begin{bmatrix} A^{(1)} \\ \vdots \\ A^{(r)} \end{bmatrix}, \text{ where } A^{(l)} \in \mathbb{B}^{\tilde{m} \times n}, \forall l \in \{1, \dots, r\}.$$

Theoretical analysis of
communication cost

$$B = \sum_{l=1}^r A^{(l)T} A^{(l)} \quad \text{and} \quad \hat{a}_i^{(l)} = \sum_k a_{ki}^{(l)}$$

Check the paper for details 😊

$$O\left(\left(1 + \frac{z}{M\sqrt{cp}}\right) \cdot \alpha + \left(\frac{cn}{\sqrt{cp}} + \frac{cn}{p}\right) \cdot \beta\right)$$

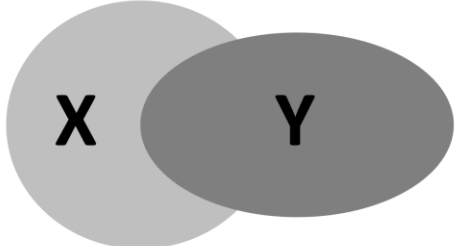
Public implementation
based on CTF [1]

SPCL

spcl.inf.ethz.ch
@spcl_eth ETH zürich

SET SIMILARITY

How can we measure the „similarity“ of X and Y?



Jaccard Index:

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

Navigation icons: back, forward, search, refresh, close

Part 1 „SimilarityAtScale“: the first communication-efficient distributed algorithm to compute the general Jaccard similarity index and distance

SET SIMILARITY

How can we measure the „similarity“ of X and Y?

Jaccard Index:

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

Part 1 „SimilarityAtScale“: the first communication-efficient distributed algorithm to compute the general Jaccard similarity index and distance

Part 2 „GenomeAtScale“: the first tool for fast, scalable, accurate, and large-scale derivations of Jaccard index between genome sequences

SET SIMILARITY

Sequences of genomes

What are X and Y in practice?

A: adenine, T: thymine, G: guanine, C: cytosine

...AGCTAACT... ...GGCTTTCT...

GENOME SEQUENCE COMPARISON: K-MERS

GENOME SEQUENCE COMPARISON: K-MERS

Example sequence

AGCTAACT

A: adenine,
T: thymine,
G: guanine,
C: cytosine

GENOME SEQUENCE COMPARISON: K-MERS

Example sequence

AGCTAACT

2-mers:

AG, GC, CT, TA, AA, AC, CT

A: adenine,
T: thymine,
G: guanine,
C: cytosine

GENOME SEQUENCE COMPARISON: K-MERS

Example sequence

AGCTAACT

2-mers:

AG, GC, CT, TA, AA, AC, CT

3-mers:

AGC, GCT, CTA, TAA, AAC, ACT

A: adenine,
T: thymine,
G: guanine,
C: cytosine

GENOME SEQUENCE COMPARISON: K-MERS

Example sequence

AGCTAACT

2-mers:

AG, GC, CT, TA, AA, AC, CT

3-mers:

AGC, GCT, CTA, TAA, AAC, ACT

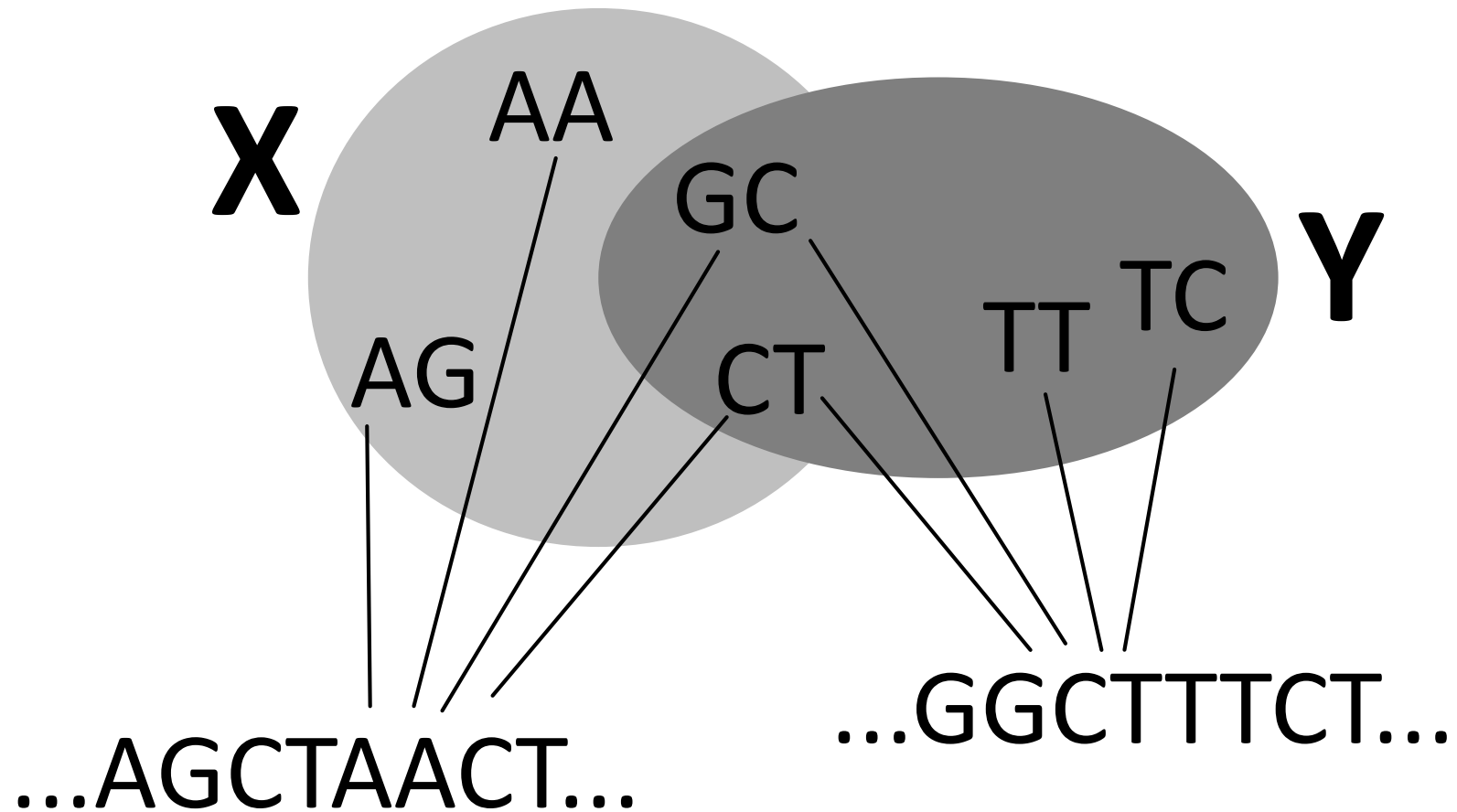
⋮

⋮

A: adenine,
T: thymine,
G: guanine,
C: cytosine

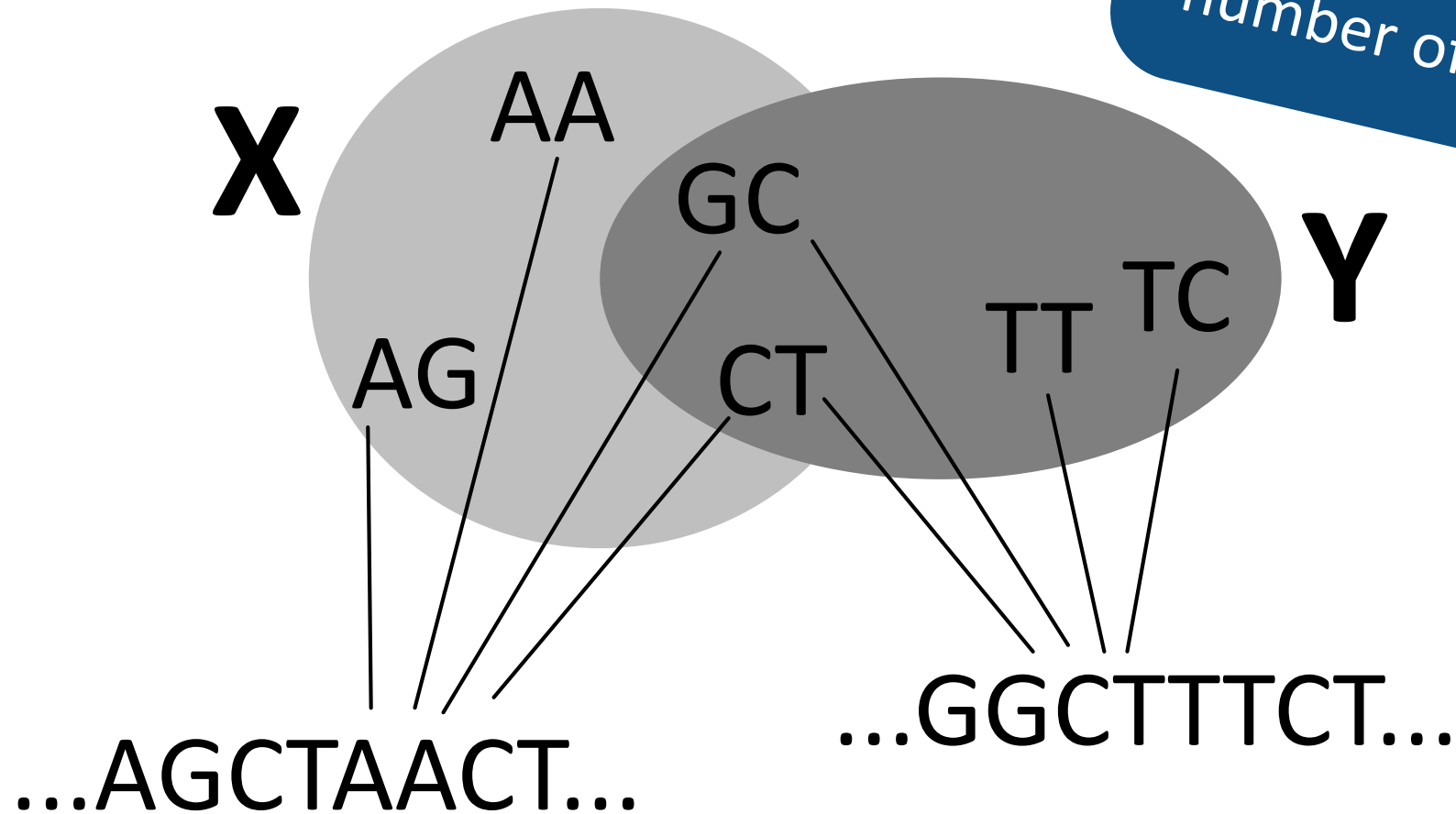
GENOME SEQUENCE COMPARISON BASED ON K-MERS

2-mers:



GENOME SEQUENCE COMPARISON BASED ON K-MERS

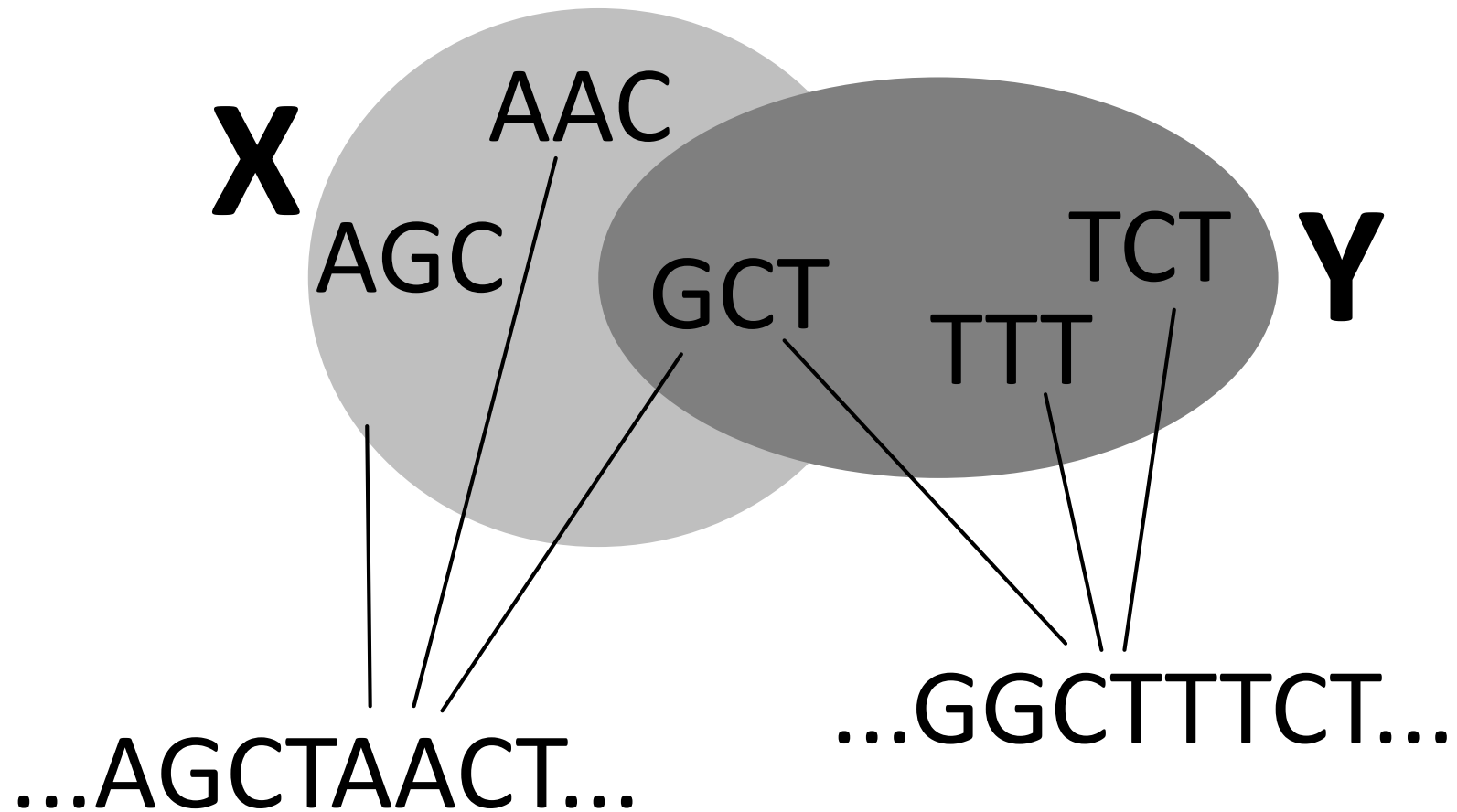
2-mers:



Similarity based on the number of common 2-mers

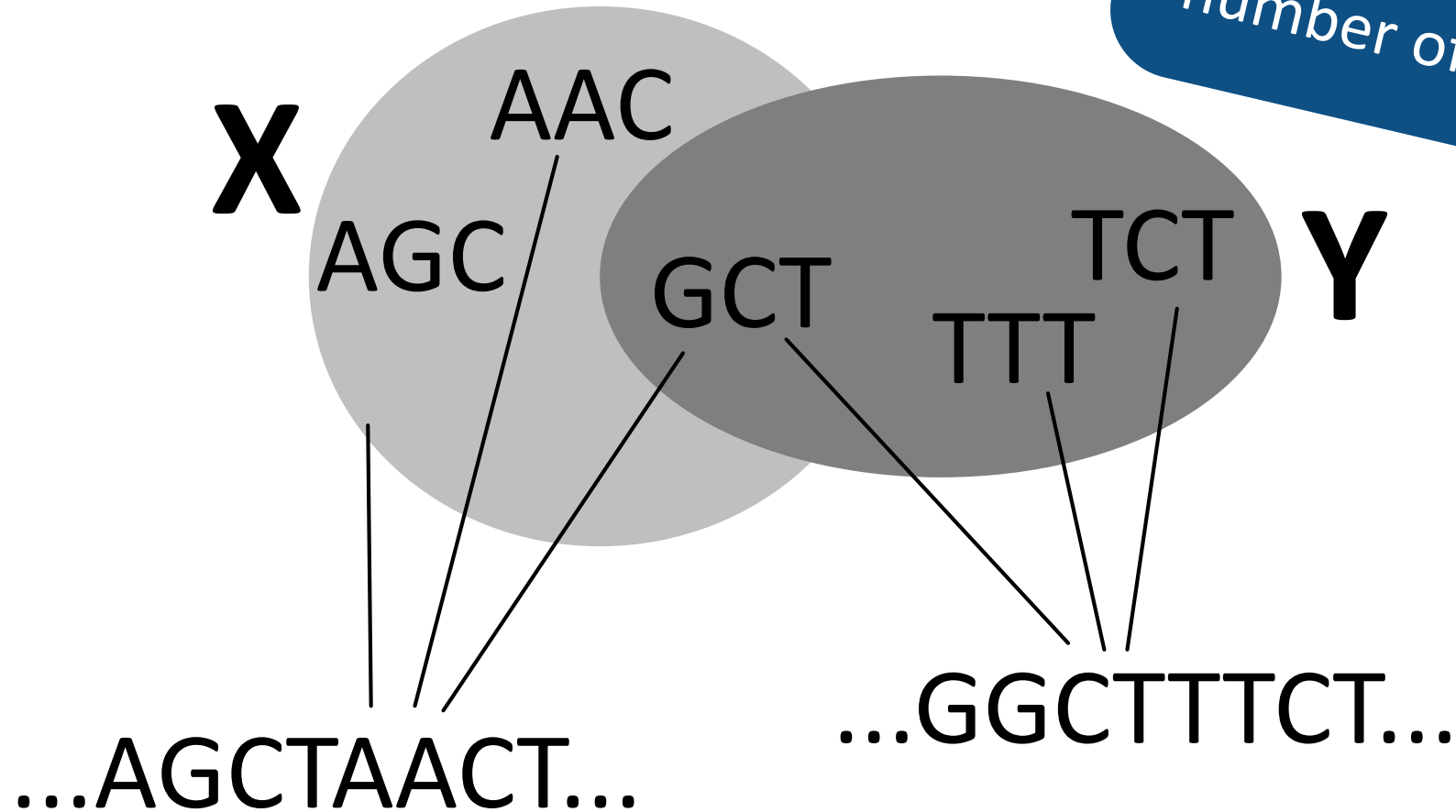
GENOME SEQUENCE COMPARISON BASED ON K-MERS

3-mers:



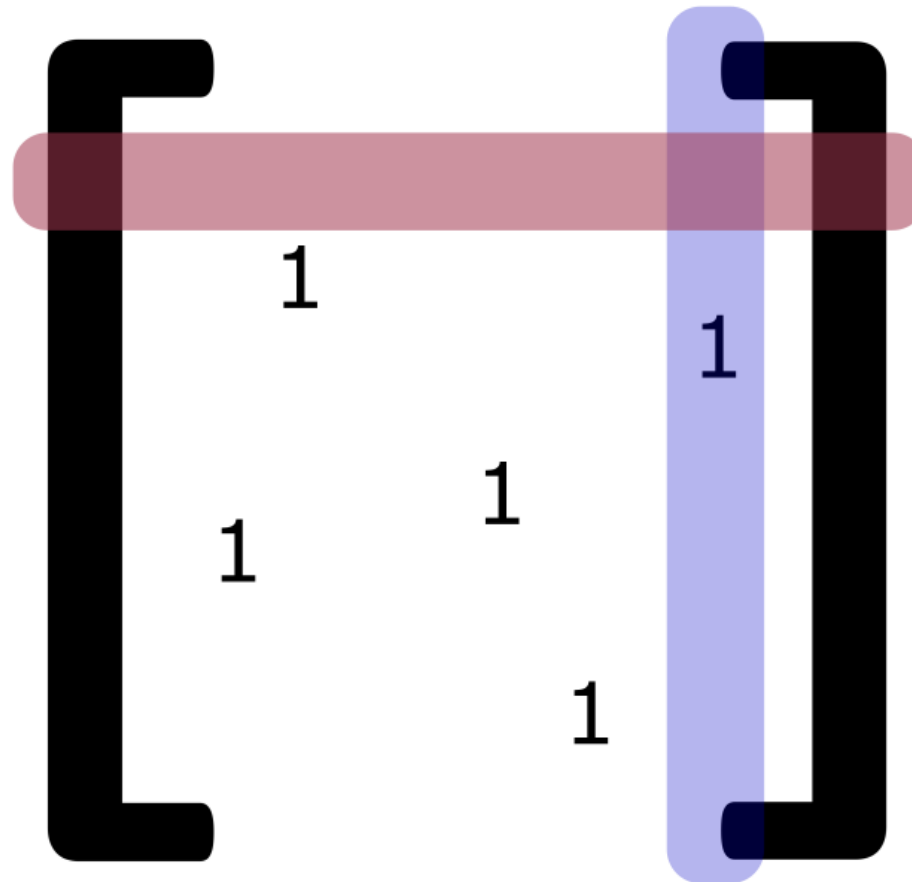
GENOME SEQUENCE COMPARISON BASED ON K-MERS

3-mers:



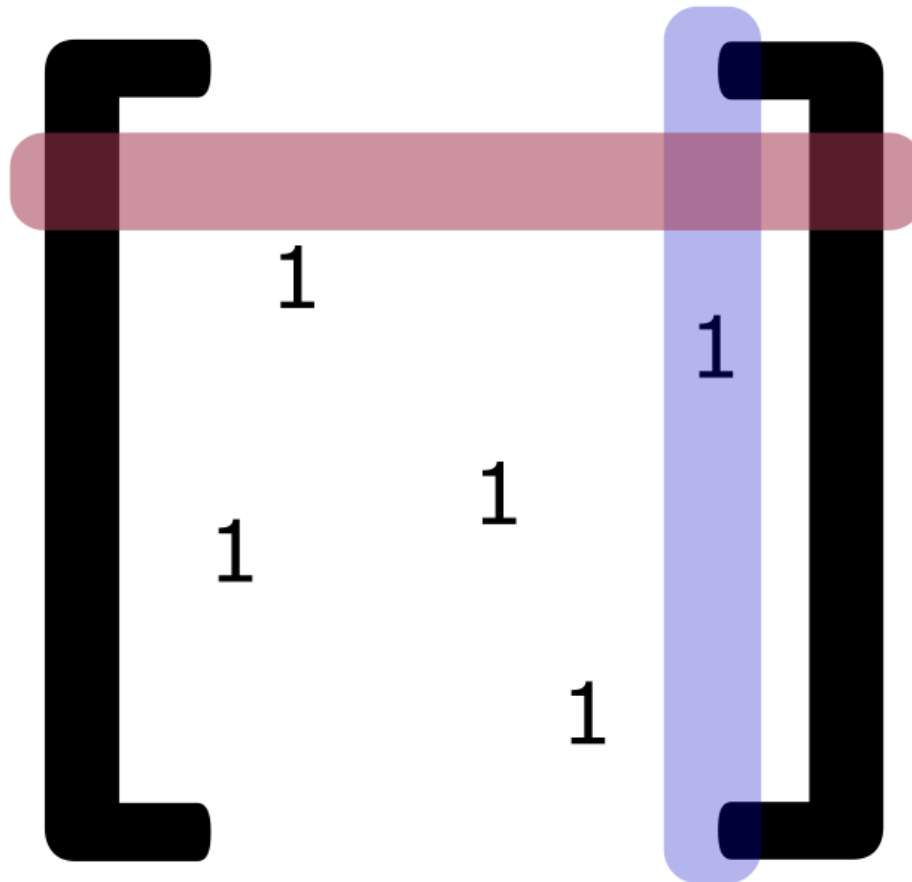
Similarity based on the number of common 3-mers

GENOMEATSCALE: ALGEBRAIC REPRESENTATION OF MATRIX A



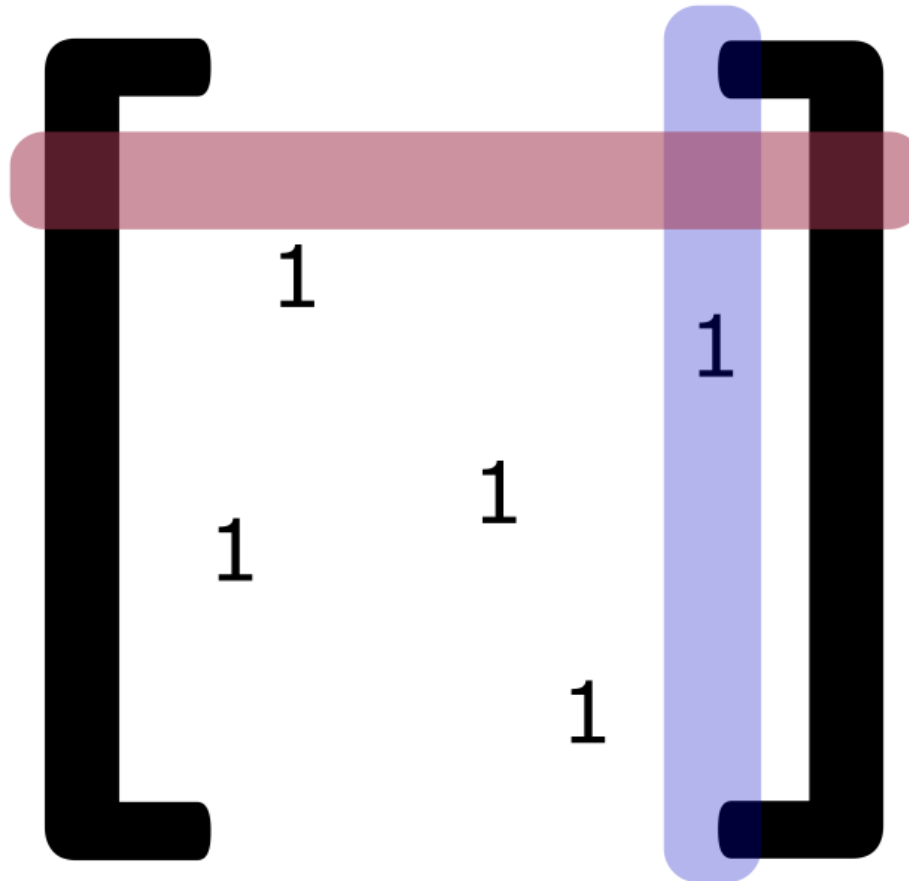
GENOMEATSCALE: ALGEBRAIC REPRESENTATION OF MATRIX A

One **row**
describes
one sequence
***k*-mer**



GENOMEATSCALE: ALGEBRAIC REPRESENTATION OF MATRIX A

One **row**
describes
one sequence
***k*-mer**

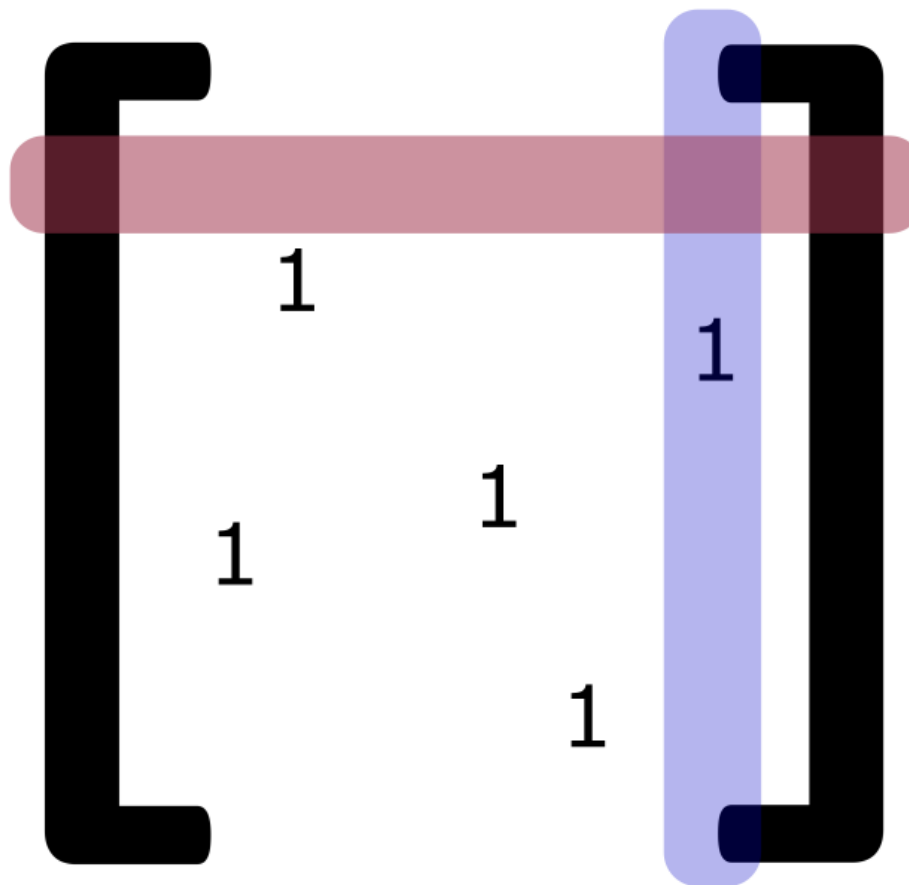


One **column**
describes
one **data**
sample

GENOMEATSCALE: ALGEBRAIC REPRESENTATION OF MATRIX A

One **row**
describes
one sequence
***k*-mer**

"**0**": a given
data sample
does not
contain a
given *k*-mer

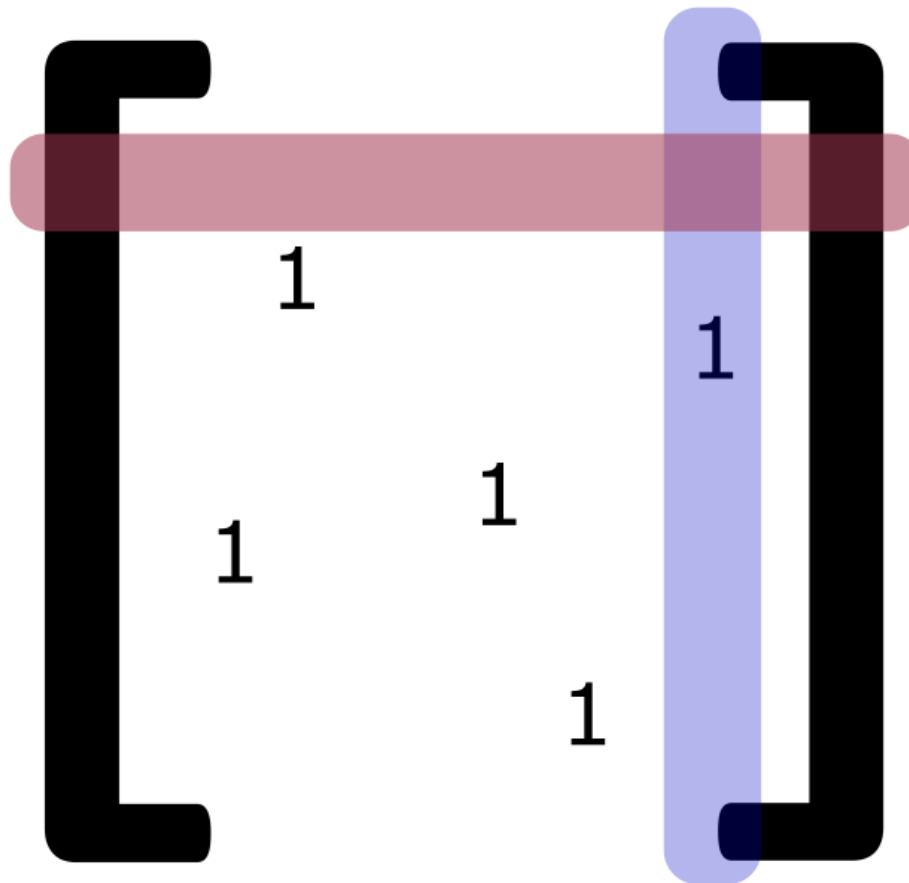


One **column**
describes
one **data**
sample

GENOMEATSCALE: ALGEBRAIC REPRESENTATION OF MATRIX A

One **row**
describes
one sequence
***k*-mer**

"0": a given
data sample
does not
contain a
given *k*-mer



One **column**
describes
one **data**
sample

"1": a given
data sample
contains a
given *k*-mer

GENOMEATSCALE: ALGEBRAIC REPRESENTATION OF MATRIX A

One **row**
describes
one sequence
k-mer

One **column**
describes
one **data**
sample

Seamless integration with
metagenomics projects:
check the paper 😊

"0": a given
data sample
does not
contain a
given *k*-mer

1

"1": a given
data sample
contains a
given *k*-mer

GENOMEATSCALE: SCALE COMPARISON





GENOMEATSCALE: SCALE COMPARISON

Tool	# compute nodes	# samples	Raw input data size	Preprocessed data size	Similarity
DSM [71]	1	435	3.3TB	N/A [‡]	Jaccard
Mash [63]	1	54, 118	N/A [†]	674GB	Jaccard (MinHash)
Libra [29]	10	40	372GB	N/A [‡]	Cosine
GenomeAtScale	1024	446, 506	170TB	1.8TB	Jaccard










GENOMEATSCALE: SCALE COMPARISON

Tool	# compute nodes	# samples	Raw input data size	Preprocessed data size	Similarity
DSM [71]	1	435	3.3TB	N/A [‡]	Jaccard
Mash [63]	1	54, 118	N/A [†]	674GB	Jaccard (MinHash)
Libra [29]	10	40	372GB	N/A [‡]	Cosine
GenomeAtScale	1024	446, 506	170TB	1.8TB	Jaccard





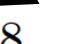







GENOMEATSCALE: SCALE COMPARISON

Tool	# compute nodes	# samples	Raw input data size	Preprocessed data size	Similarity
DSM [71]	1 	435	3.3TB	N/A [‡]	Jaccard
Mash [63]	1 	54, 118	N/A [†]	674GB	Jaccard (MinHash)
Libra [29]	10 	40	372GB	N/A [‡]	Cosine
GenomeAtScale	1024 	446, 506	170TB	1.8TB	Jaccard

GENOMEATSCALE: SCALE COMPARISON

Tool	# compute nodes	# samples	Raw input data size	Preprocessed data size	Similarity
DSM [71]	1 	 435 	3.3TB	N/A [‡]	Jaccard
Mash [63]	1 	 54, 118	N/A [†]	674GB	Jaccard (MinHash)
Libra [29]	10 	 40	372GB	N/A [‡]	Cosine
GenomeAtScale	1024 	446, 506 	170TB	1.8TB	Jaccard

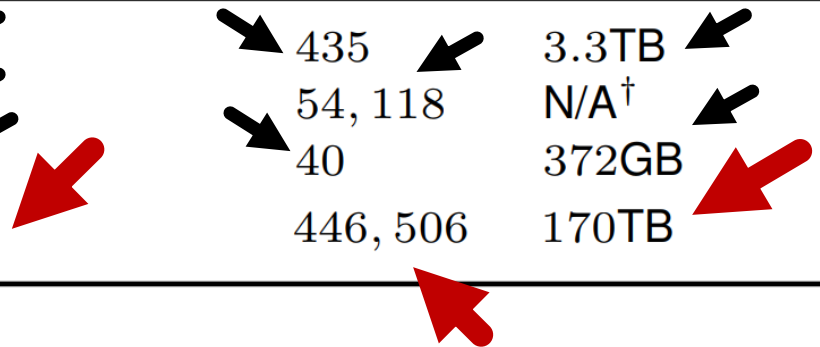
GENOMEATSCALE: SCALE COMPARISON

Tool	# compute nodes	# samples	Raw input data size	Preprocessed data size	Similarity
DSM [71]	1 	435 	3.3TB 	N/A [‡]	Jaccard
Mash [63]	1 	54, 118 	N/A [†] 	674GB	Jaccard (MinHash)
Libra [29]	10 	40 	372GB 	N/A [‡]	Cosine
GenomeAtScale	1024 	446, 506 	170TB 	1.8TB	Jaccard

GENOMEATSCALE: SCALE COMPARISON

„GenomeAtScale” achieves larger problem size and parallelism scales than past approaches

Tool	# compute nodes	# samples	Raw input data size	Preprocessed data size	Similarity
DSM [71]	1	435	3.3TB	N/A [‡]	Jaccard
Mash [63]	1	54, 118	N/A [†]	674GB	Jaccard (MinHash)
Libra [29]	10	40	372GB	N/A [‡]	Cosine
GenomeAtScale	1024	446, 506	170TB	1.8TB	Jaccard



PERFORMANCE ANALYSIS: STAMPEDE2 SUPERCOMPUTER



PERFORMANCE ANALYSIS: STAMPEDE2 SUPERCOMPUTER

Each node has an Intel Xeon Phi 7250 CPU ("Knights Landing") with 68 cores, 96GB RAM, and 16GB of high-speed on-chip MCDRAM



PERFORMANCE ANALYSIS: STAMPEDE2 SUPERCOMPUTER

Each node has an Intel Xeon Phi 7250 CPU (“Knights Landing”) with 68 cores, 96GB RAM, and 16GB of high-speed on-chip MCDRAM

The network: a fat tree topology (six core switches) and the 100 Gb/sec Intel Omni-Path architecture

PERFORMANCE ANALYSIS: STATISTICAL ANALYSIS



PERFORMANCE ANALYSIS: STATISTICAL ANALYSIS

We calculate 95% confidence intervals for the reported mean values by assuming the batch times are normally distributed samples.

PERFORMANCE ANALYSIS: STATISTICAL ANALYSIS

We calculate 95% confidence intervals for the reported mean values by assuming the batch times are normally distributed samples.

The derived confidence intervals are very tight around the means, and we exclude them from the plots to ensure clarity

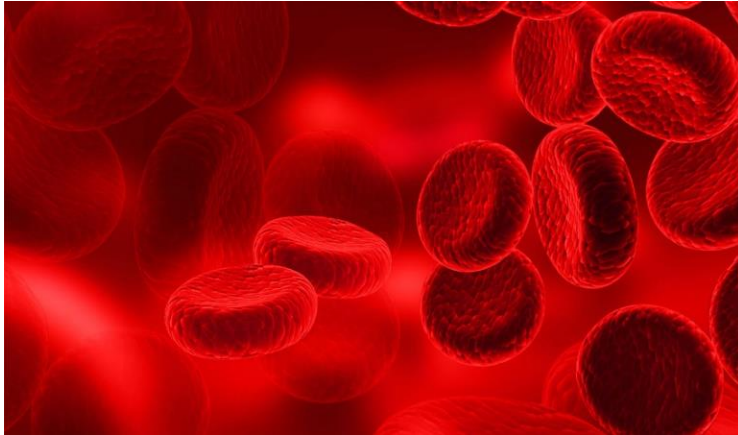
PERFORMANCE ANALYSIS: USED DATASETS

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



[1] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. Nature biotechnology, 2016.

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set



PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

2,580 RNASeq
experiments

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

2,580 RNASeq
experiments

k -mer size
of 19

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

2,580 RNASeq experiments

k -mer size of 19

Indicator matrix (A)
sparsity: $1.5 \cdot 10^{-4}$

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

2,580 RNASeq experiments

k -mer size of 19

Indicator matrix (A)
sparsity: $1.5 \cdot 10^{-4}$

BIGSI [2]

[1] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. Nature biotechnology, 2016.

[2] P. Bradley, H. C. den Bakker, E. P. Rocha, G. McVean, and Z. Iqbal. Ultrafast search of all deposited bacterial and viral genomic data. Nature biotechnology, 2019.

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]

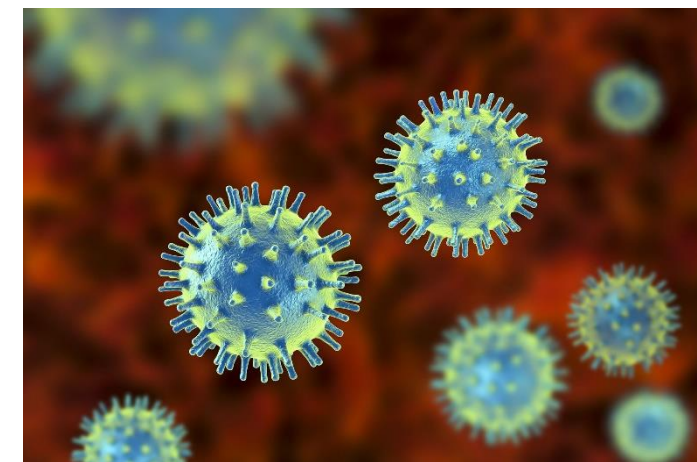


Low-variability set

2,580 RNASeq experiments

k -mer size of 19

Indicator matrix (A)
sparsity: $1.5 \cdot 10^{-4}$



BIGSI [2]

[1] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. Nature biotechnology, 2016.

[2] P. Bradley, H. C. den Bakker, E. P. Rocha, G. McVean, and Z. Iqbal. Ultrafast search of all deposited bacterial and viral genomic data. Nature biotechnology, 2019.

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

2,580 RNASeq experiments

k -mer size of 19

Indicator matrix (A)
sparsity: $1.5 \cdot 10^{-4}$



High-variability set

BIGSI [2]

[1] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. Nature biotechnology, 2016.

[2] P. Bradley, H. C. den Bakker, E. P. Rocha, G. McVean, and Z. Iqbal. Ultrafast search of all deposited bacterial and viral genomic data. Nature biotechnology, 2019.

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

2,580 RNASeq experiments

k -mer size of 19

Indicator matrix (A) sparsity: $1.5 \cdot 10^{-4}$



446,506 samples

High-variability set

BIGSI [2]

[1] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. Nature biotechnology, 2016.

[2] P. Bradley, H. C. den Bakker, E. P. Rocha, G. McVean, and Z. Iqbal. Ultrafast search of all deposited bacterial and viral genomic data. Nature biotechnology, 2019.

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

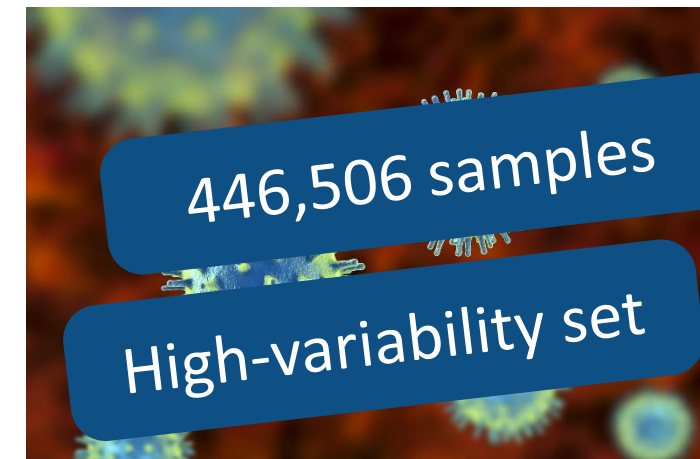
2,580 RNASeq experiments

k-mer size of 19

Indicator matrix (A)
sparsity: $1.5 \cdot 10^{-4}$



k-mer size of 31



446,506 samples

High-variability set

BIGSI [2]

[1] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. Nature biotechnology, 2016.

[2] P. Bradley, H. C. den Bakker, E. P. Rocha, G. McVean, and Z. Iqbal. Ultrafast search of all deposited bacterial and viral genomic data. Nature biotechnology, 2019.

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

2,580 RNASeq experiments

k -mer size of 19

Indicator matrix (A)
sparsity: $1.5 \cdot 10^{-4}$

k -mer size of 31

Indicator matrix (A)
sparsity: $4 \cdot 10^{-12}$

446,506 samples

High-variability set

BIGSI [2]

[1] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. Nature biotechnology, 2016.

[2] P. Bradley, H. C. den Bakker, E. P. Rocha, G. McVean, and Z. Iqbal. Ultrafast search of all deposited bacterial and viral genomic data. Nature biotechnology, 2019.

PERFORMANCE ANALYSIS: USED DATASETS

BBB/Kingsford [1]



Low-variability set

2,580 RNASeq experiments

k -mer size of 19

Indicator matrix (A) sparsity: $1.5 \cdot 10^{-4}$

Indicator matrix (A) with a pre-determined sparsity p

Synthetic data

k -mer size of 31

Indicator matrix (A) sparsity: $4 \cdot 10^{-12}$

446,506 samples

High-variability set

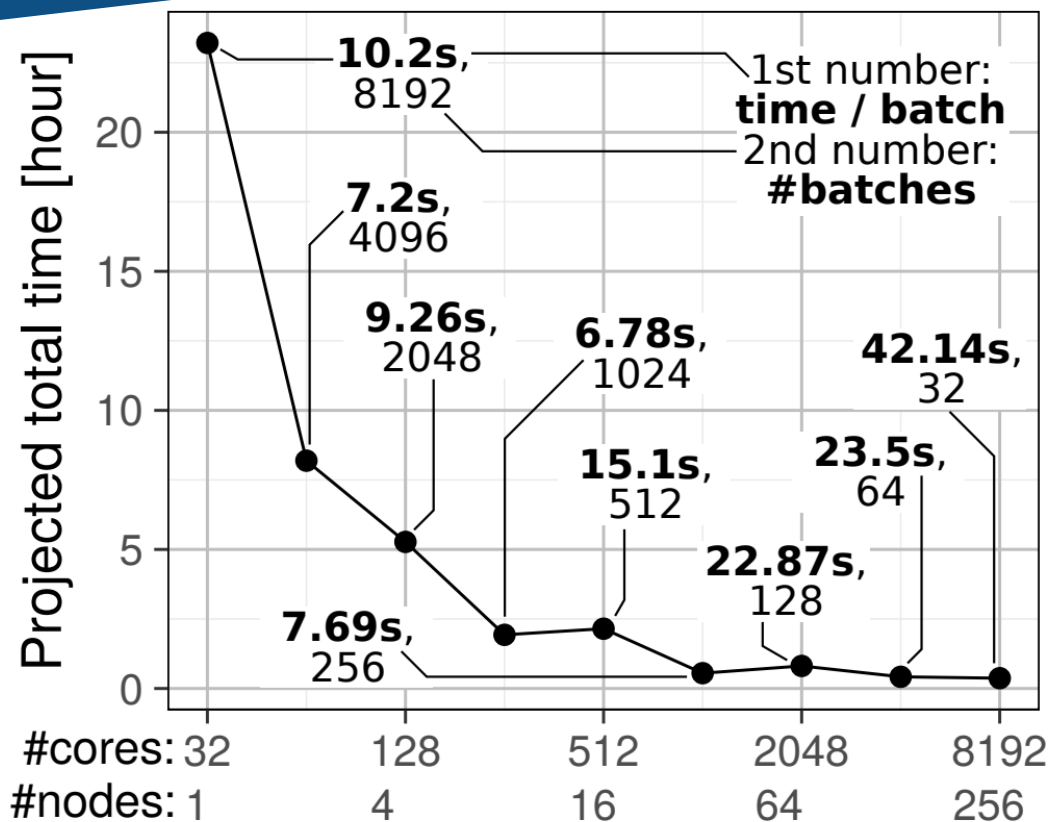
BIGSI [2]

[1] B. Solomon and C. Kingsford. Fast search of thousands of short-read sequencing experiments. Nature biotechnology, 2016.

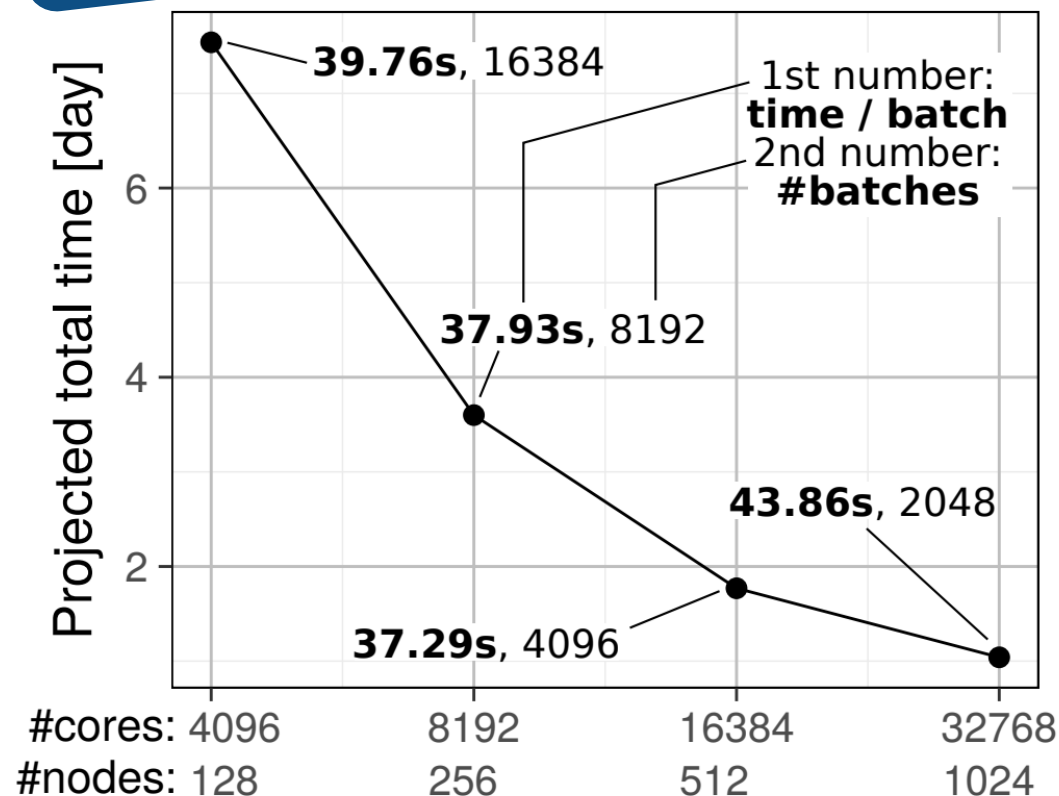
[2] P. Bradley, H. C. den Bakker, E. P. Rocha, G. McVean, and Z. Iqbal. Ultrafast search of all deposited bacterial and viral genomic data. Nature biotechnology, 2019.

PERFORMANCE ANALYSIS: REAL DATA, STRONG SCALING

BBB/Kingsford



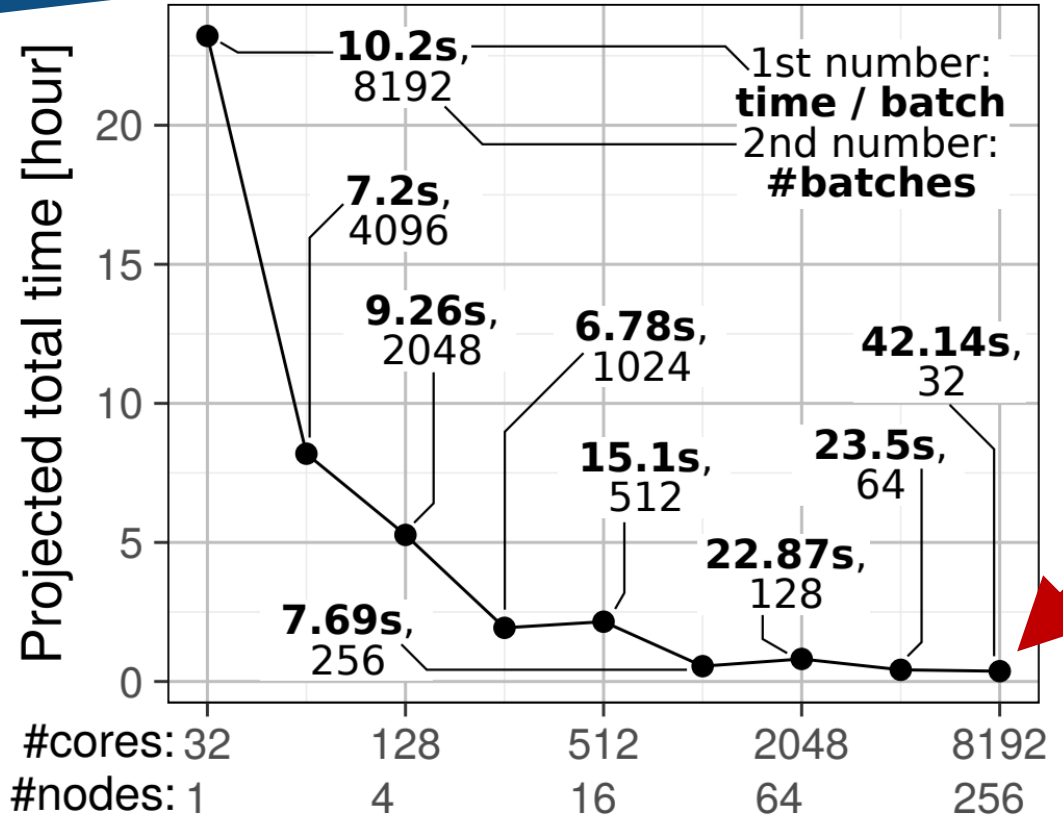
BIGSI



PERFORMANCE ANALYSIS: REAL DATA, STRONG SCALING

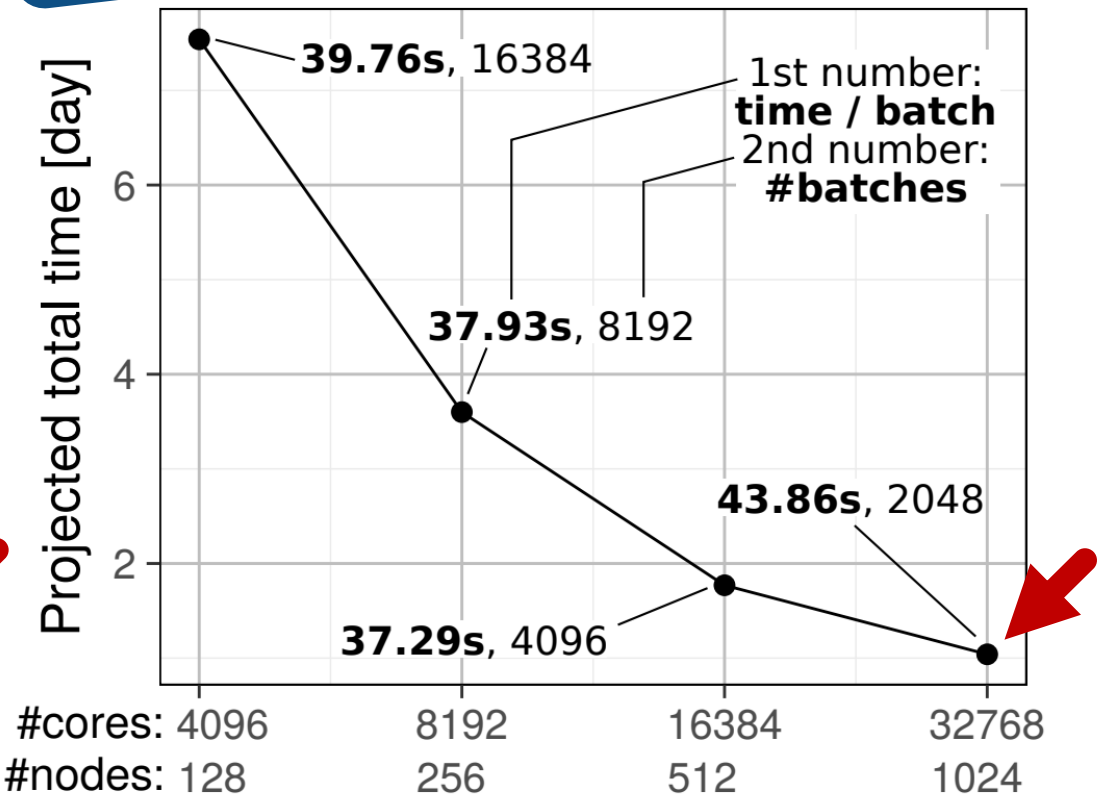
BBB/Kingsford

< 1 hour for whole dataset



BIGSI

< 1 day for whole dataset



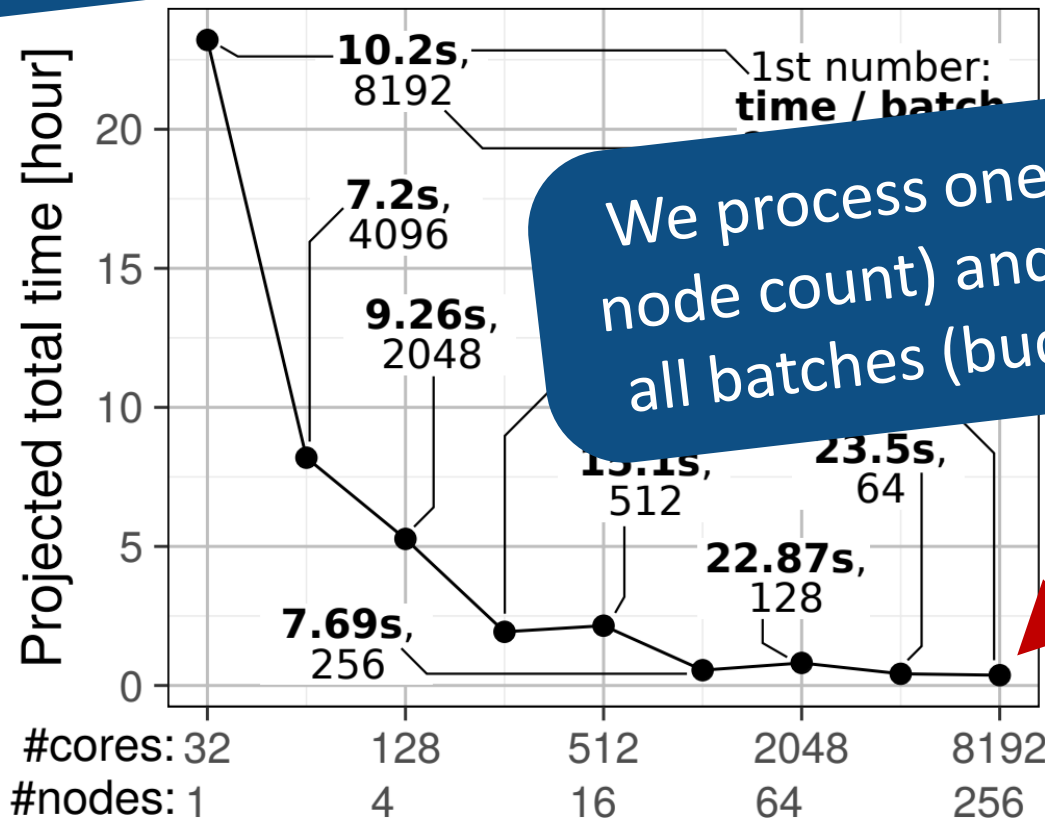
PERFORMANCE ANALYSIS: REAL DATA, STRONG SCALING

BBB/Kingsford

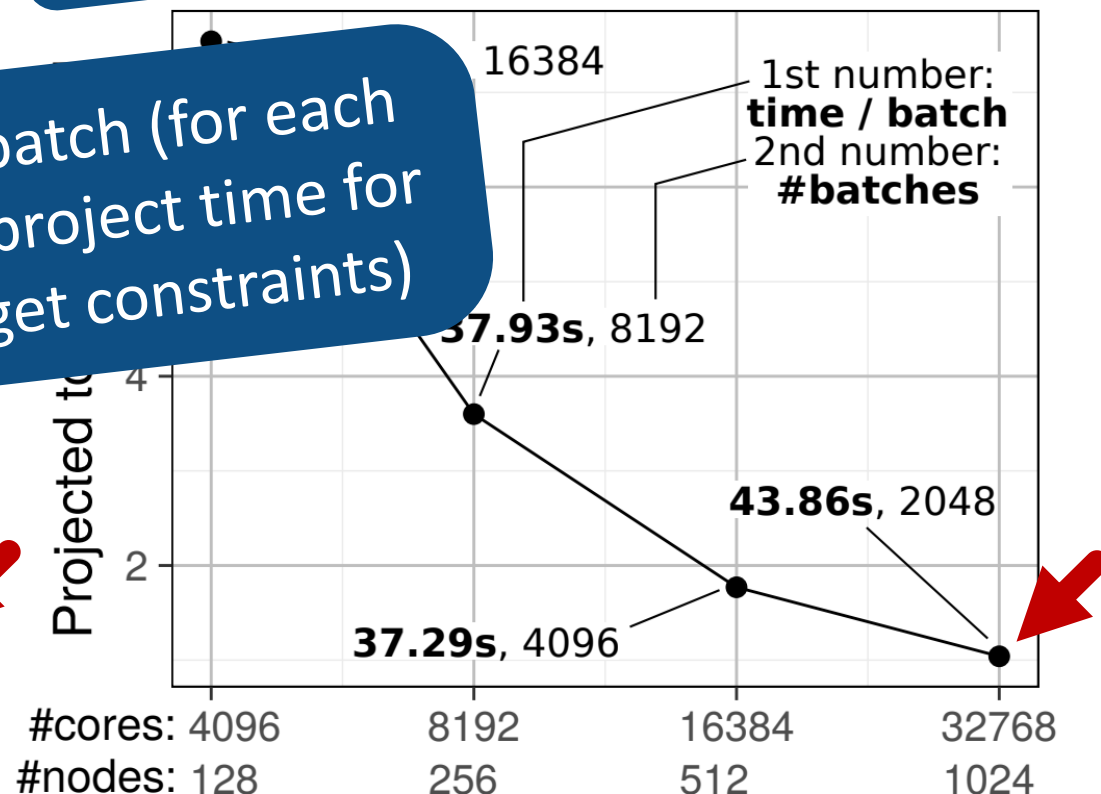
< 1 hour for whole dataset

BIGSI

< 1 day for whole dataset



We process one batch (for each node count) and project time for all batches (budget constraints)



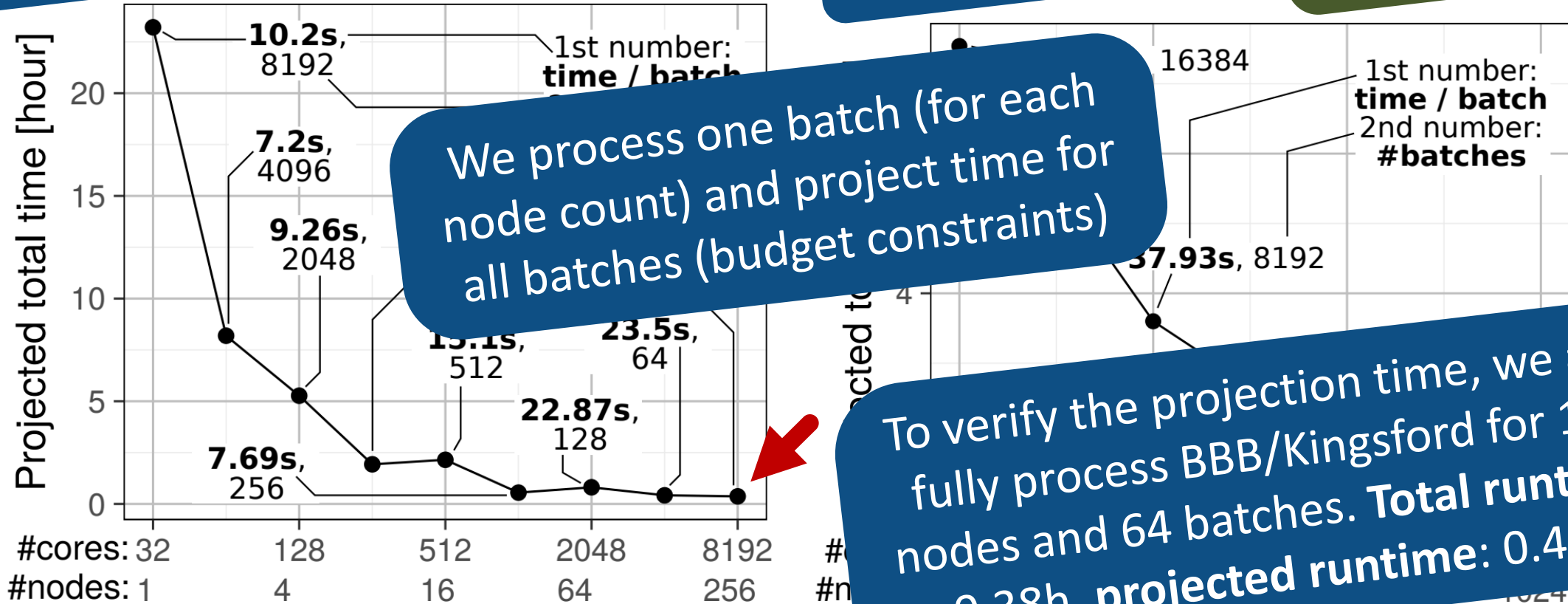
PERFORMANCE ANALYSIS: REAL DATA, STRONG SCALING

BBB/Kingsford

< 1 hour for whole dataset

BIGSI

< 1 day for whole dataset



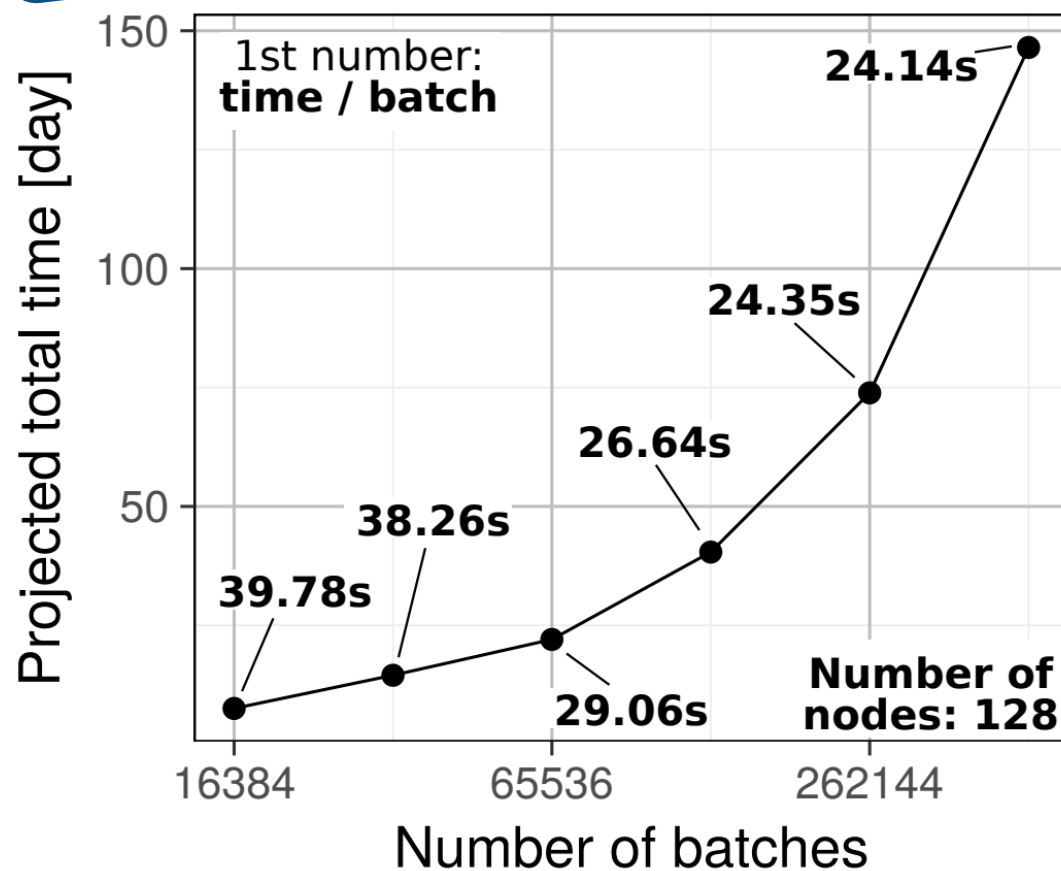
We process one batch (for each node count) and project time for all batches (budget constraints)

To verify the projection time, we also fully process BBB/Kingsford for 128 nodes and 64 batches. Total runtime: 0.38h, projected runtime: 0.42h

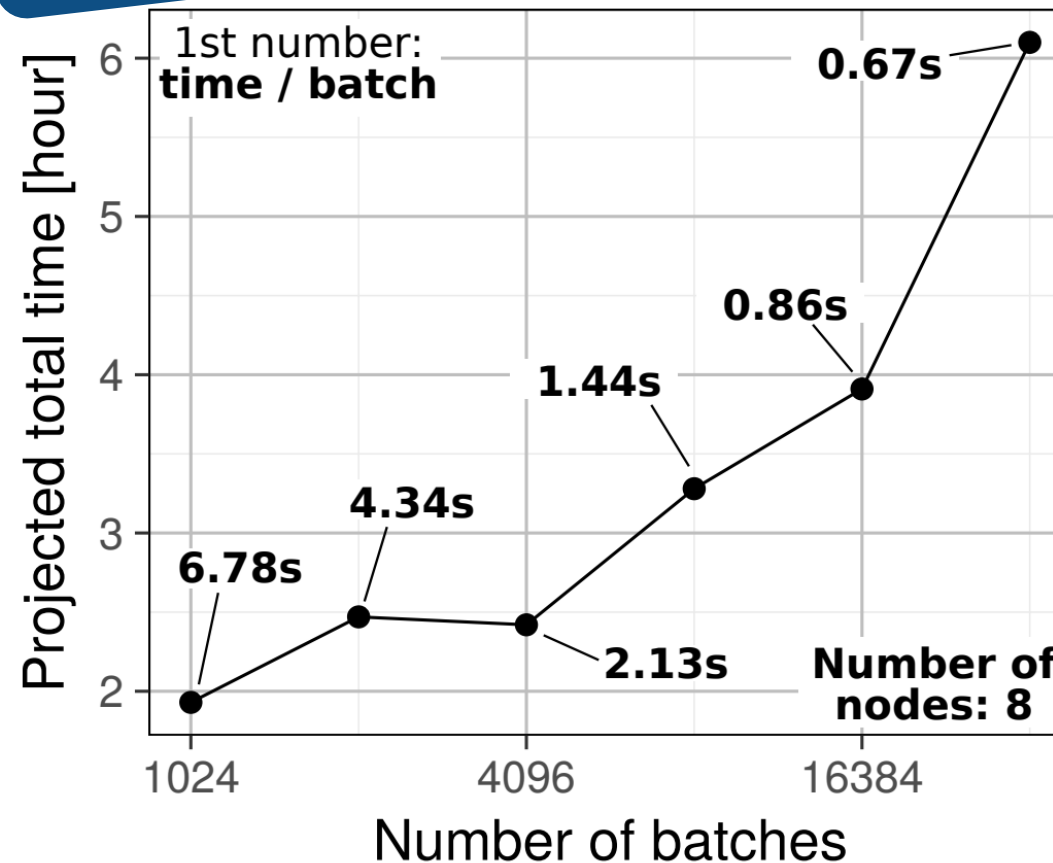


PERFORMANCE ANALYSIS: REAL DATA, SENSITIVITY ANALYSIS (#BATCHES)

BBB/Kingsford

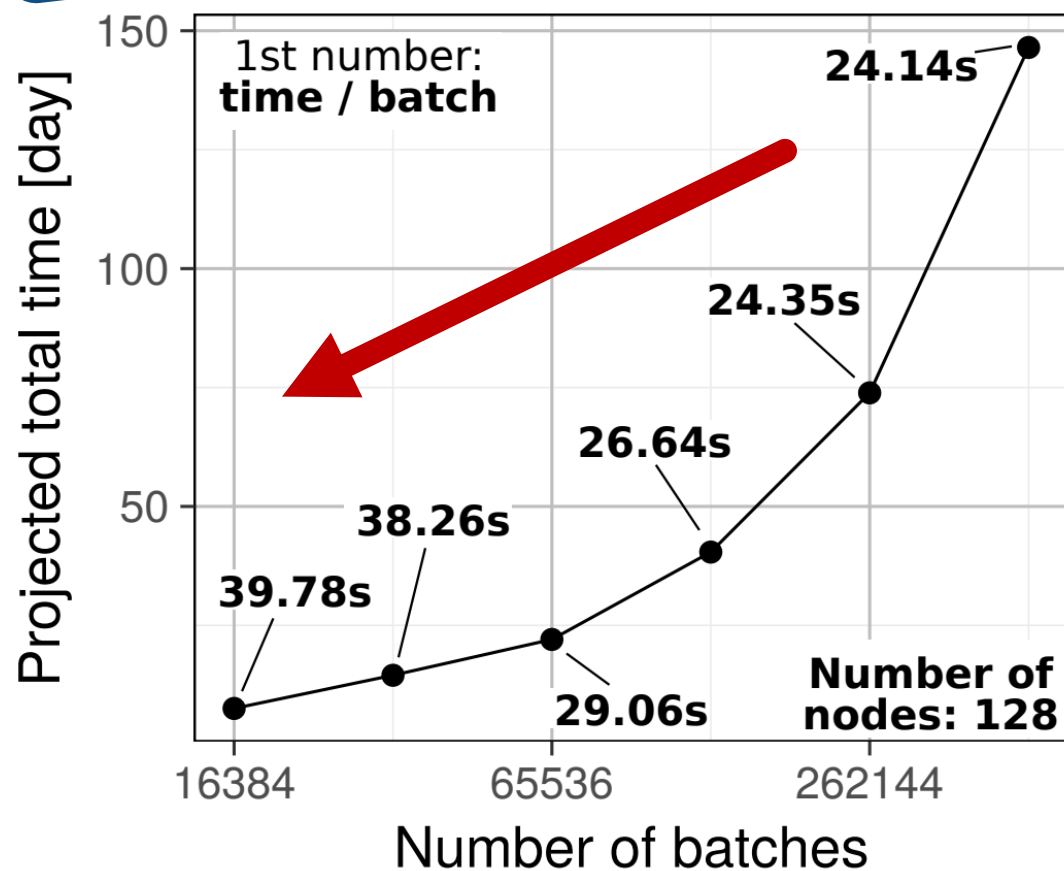


BIGSI

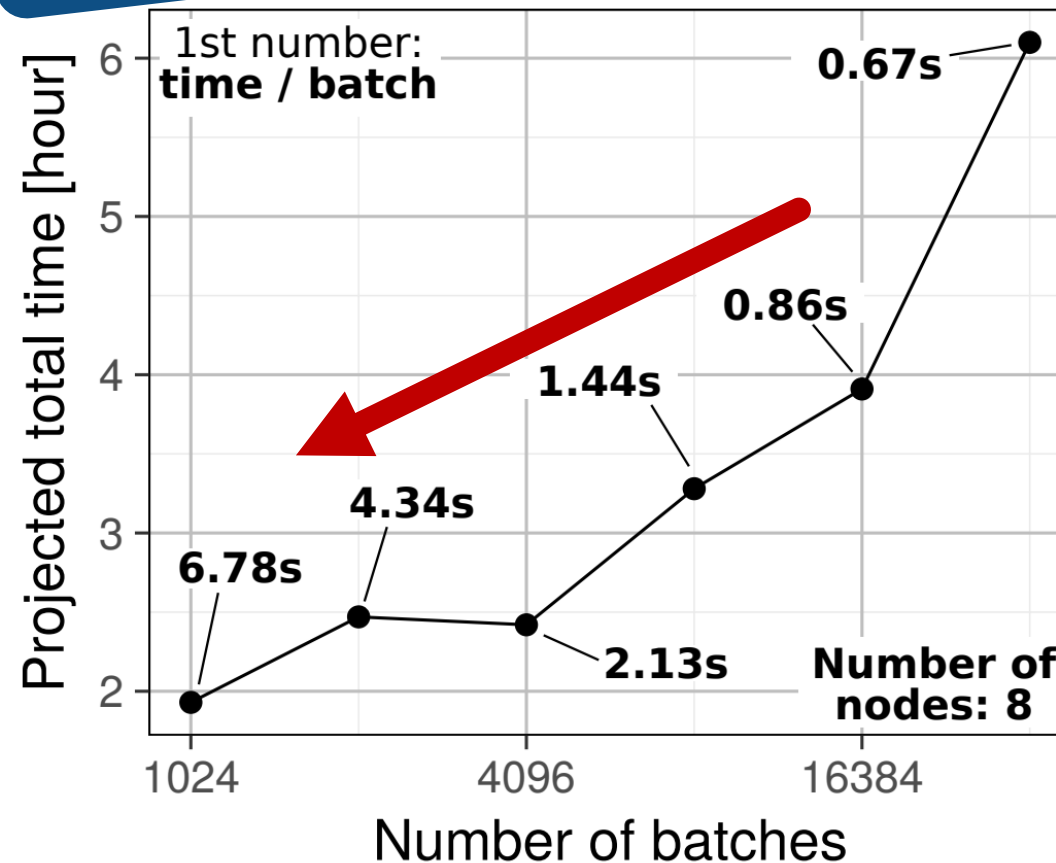


PERFORMANCE ANALYSIS: REAL DATA, SENSITIVITY ANALYSIS (#BATCHES)

BBB/Kingsford

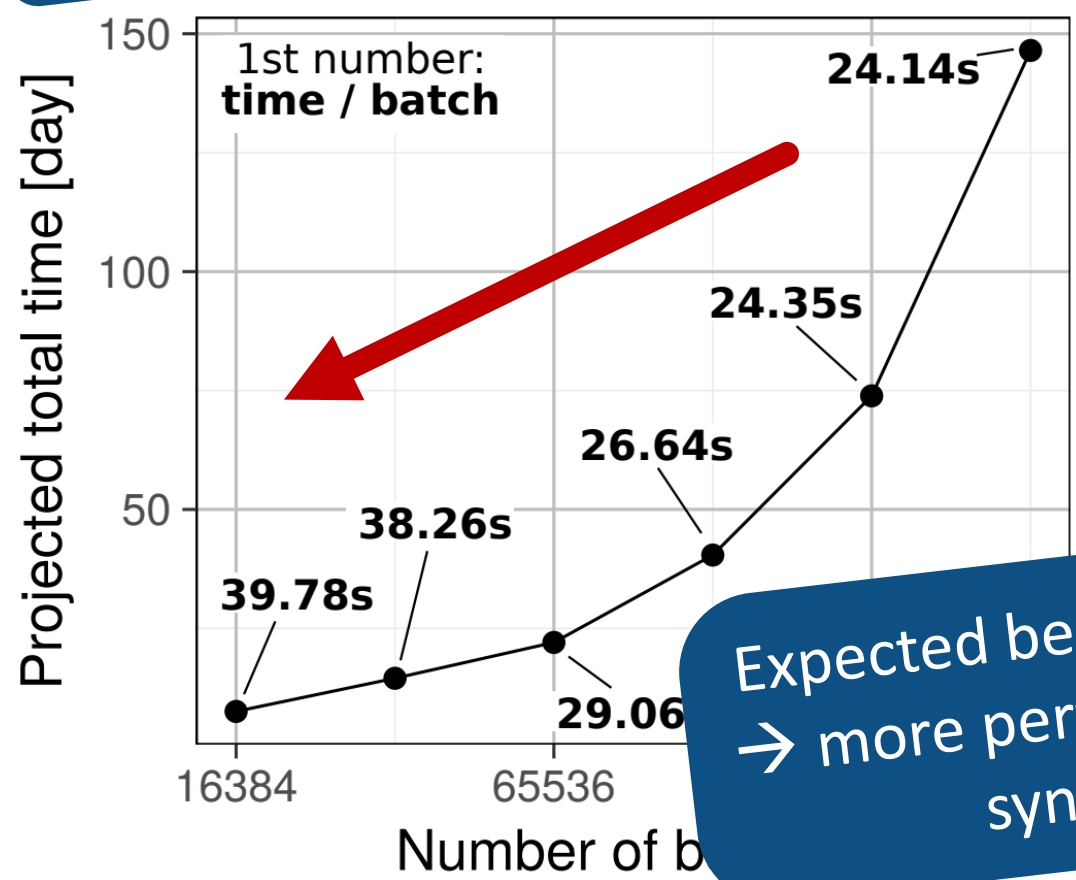


BIGSI

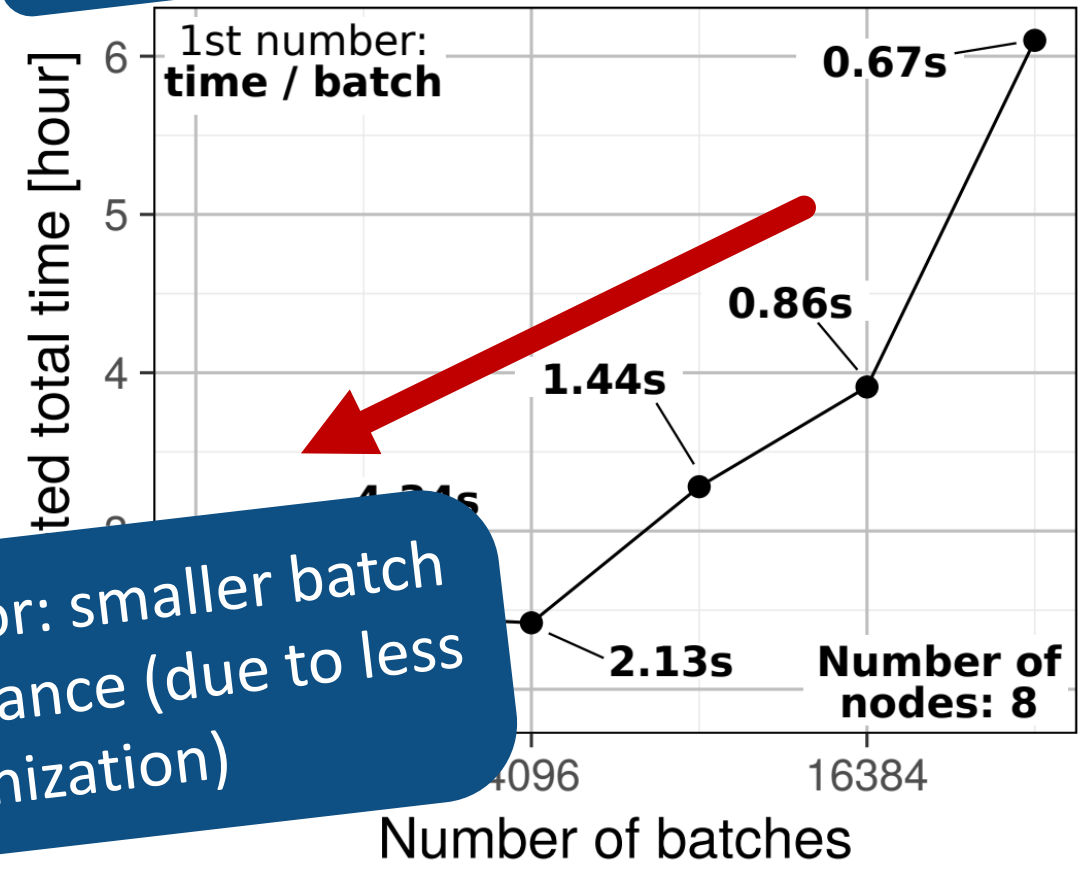


PERFORMANCE ANALYSIS: REAL DATA, SENSITIVITY ANALYSIS (#BATCHES)

BBB/Kingsford



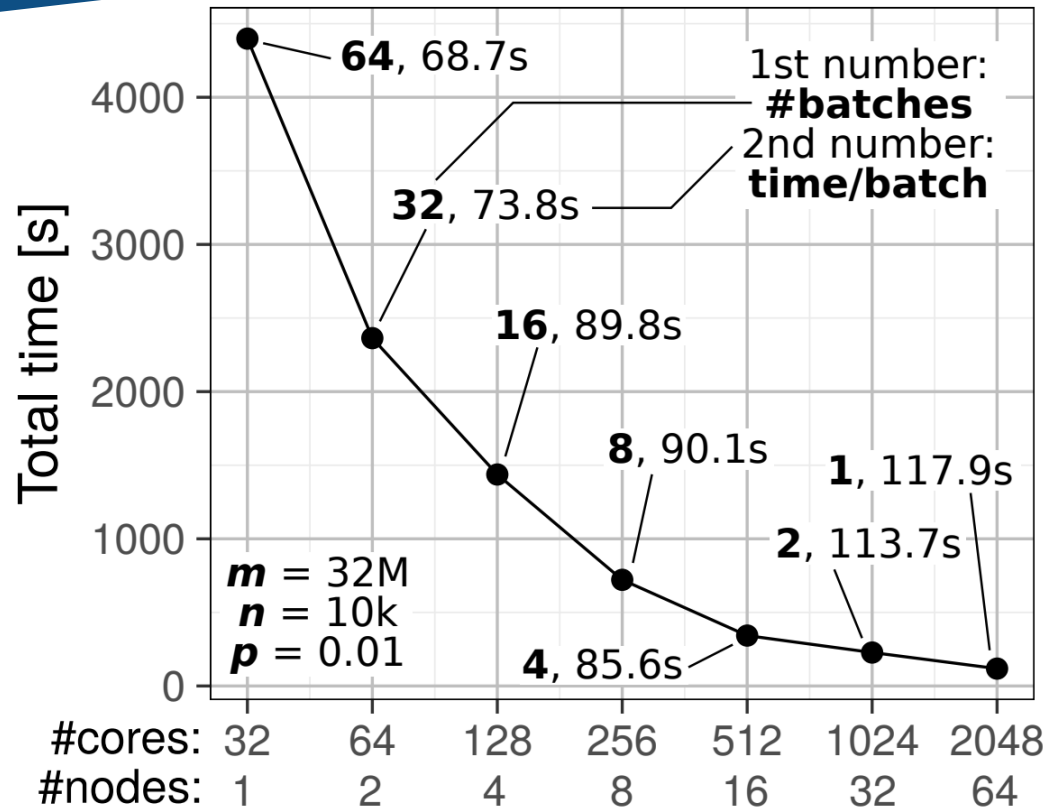
BIGSI



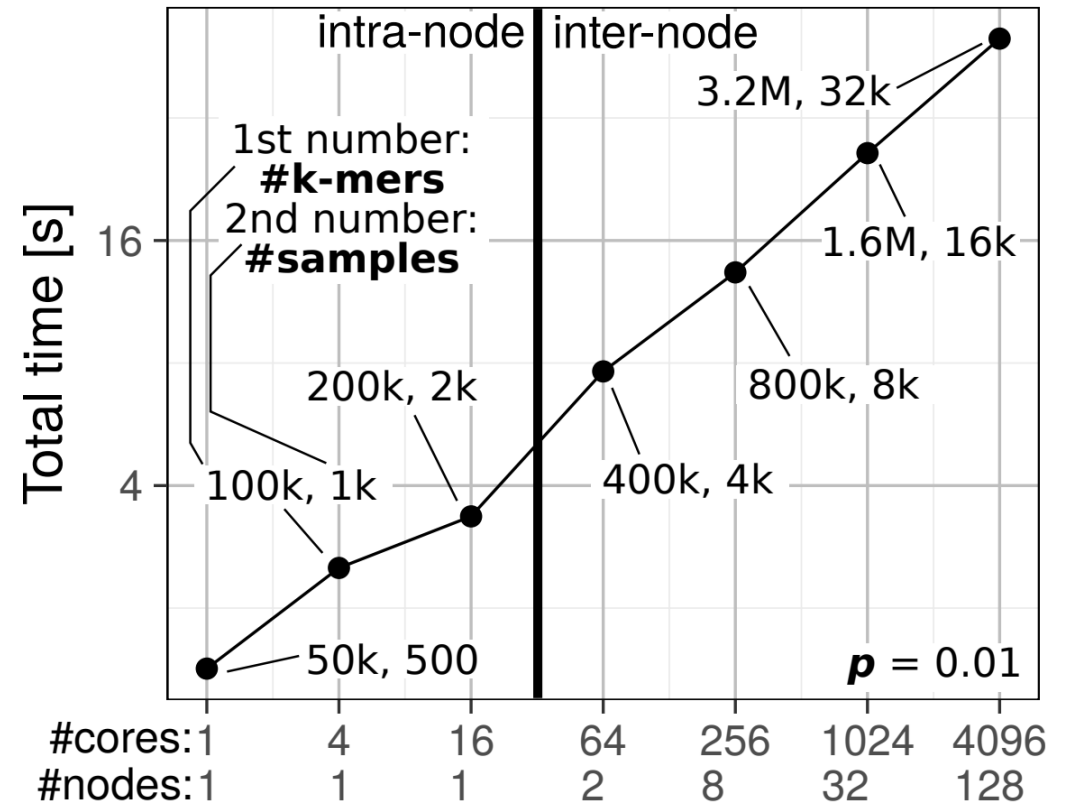
Expected behavior: smaller batch
 → more performance (due to less synchronization)

PERFORMANCE ANALYSIS: SYNTHETIC DATA

Strong scaling



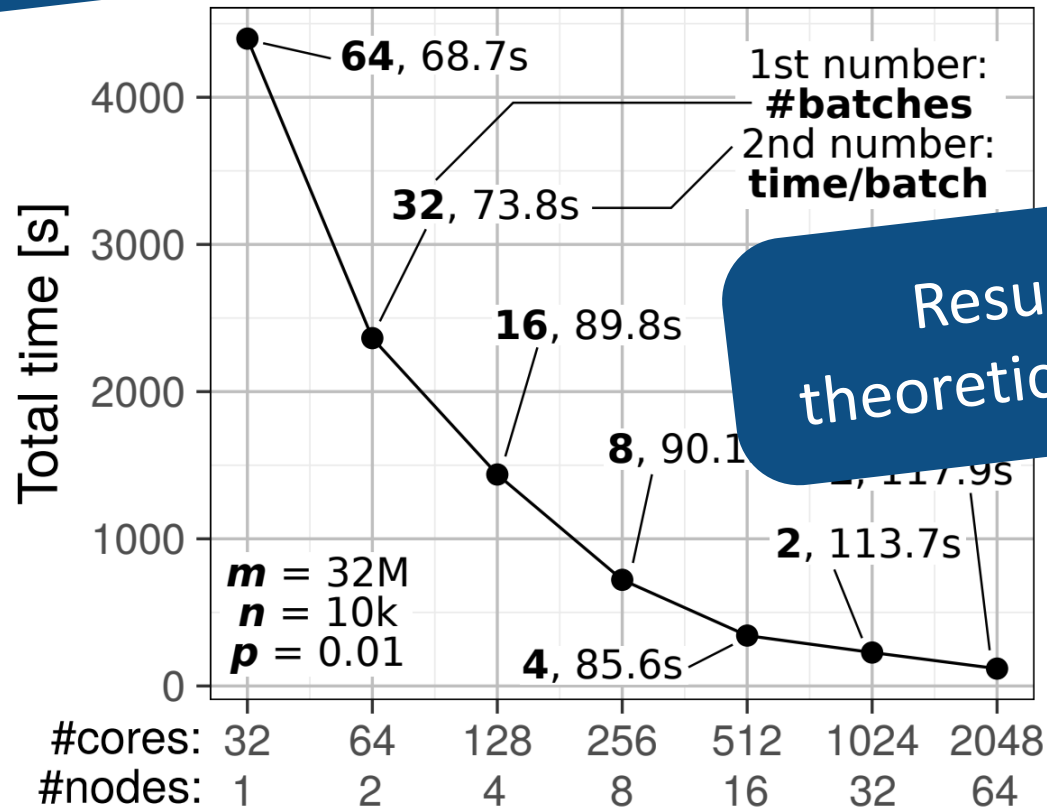
Weak scaling



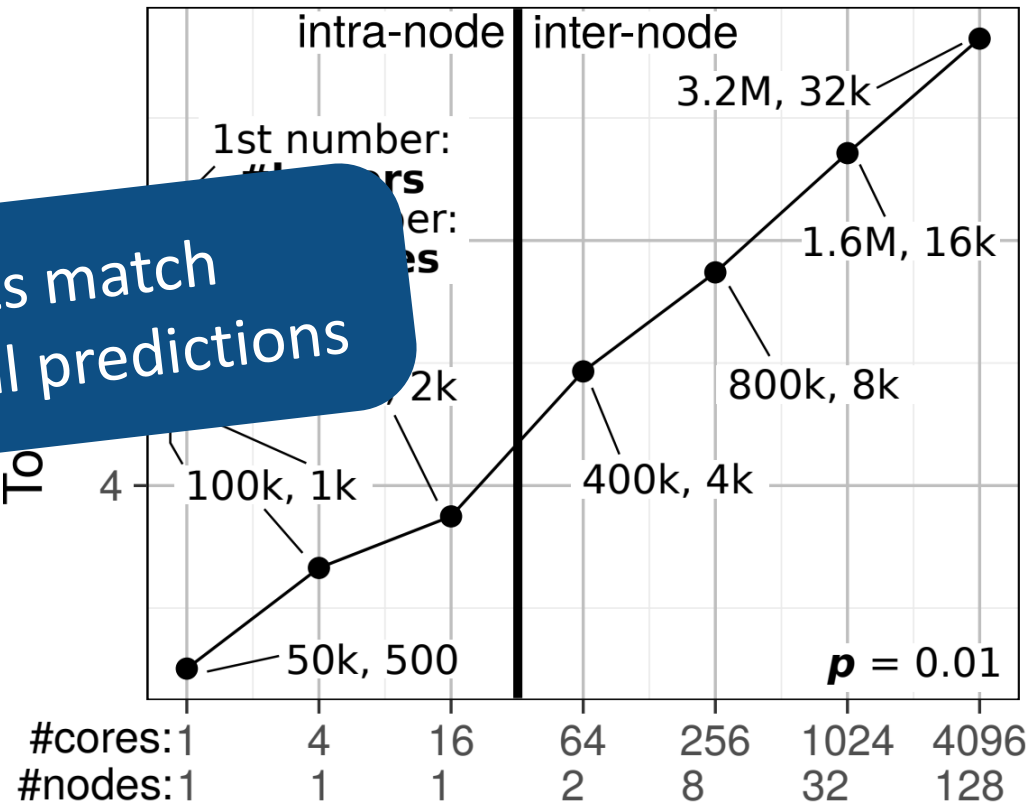
PERFORMANCE ANALYSIS: SYNTHETIC DATA

Strong scaling

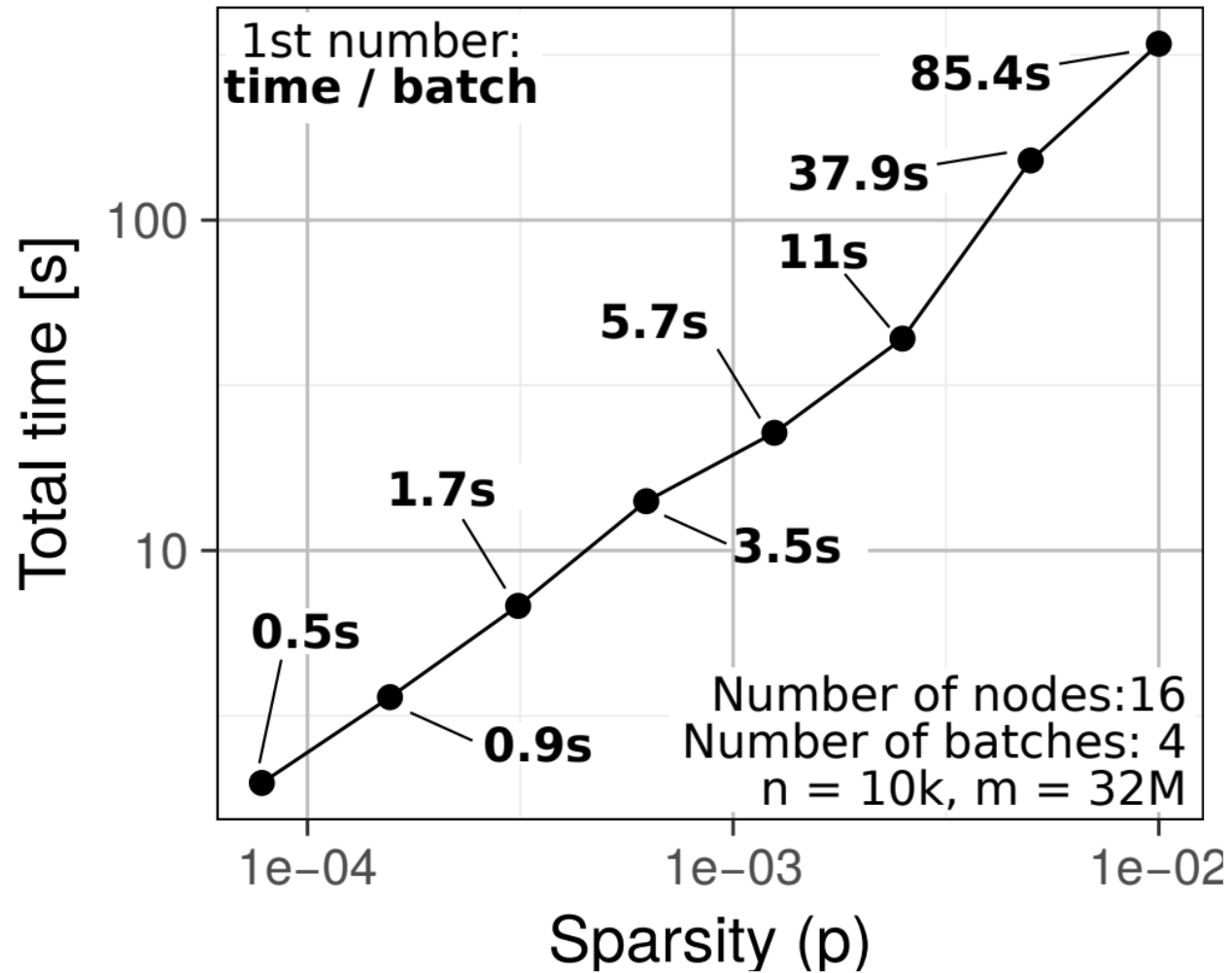
Weak scaling



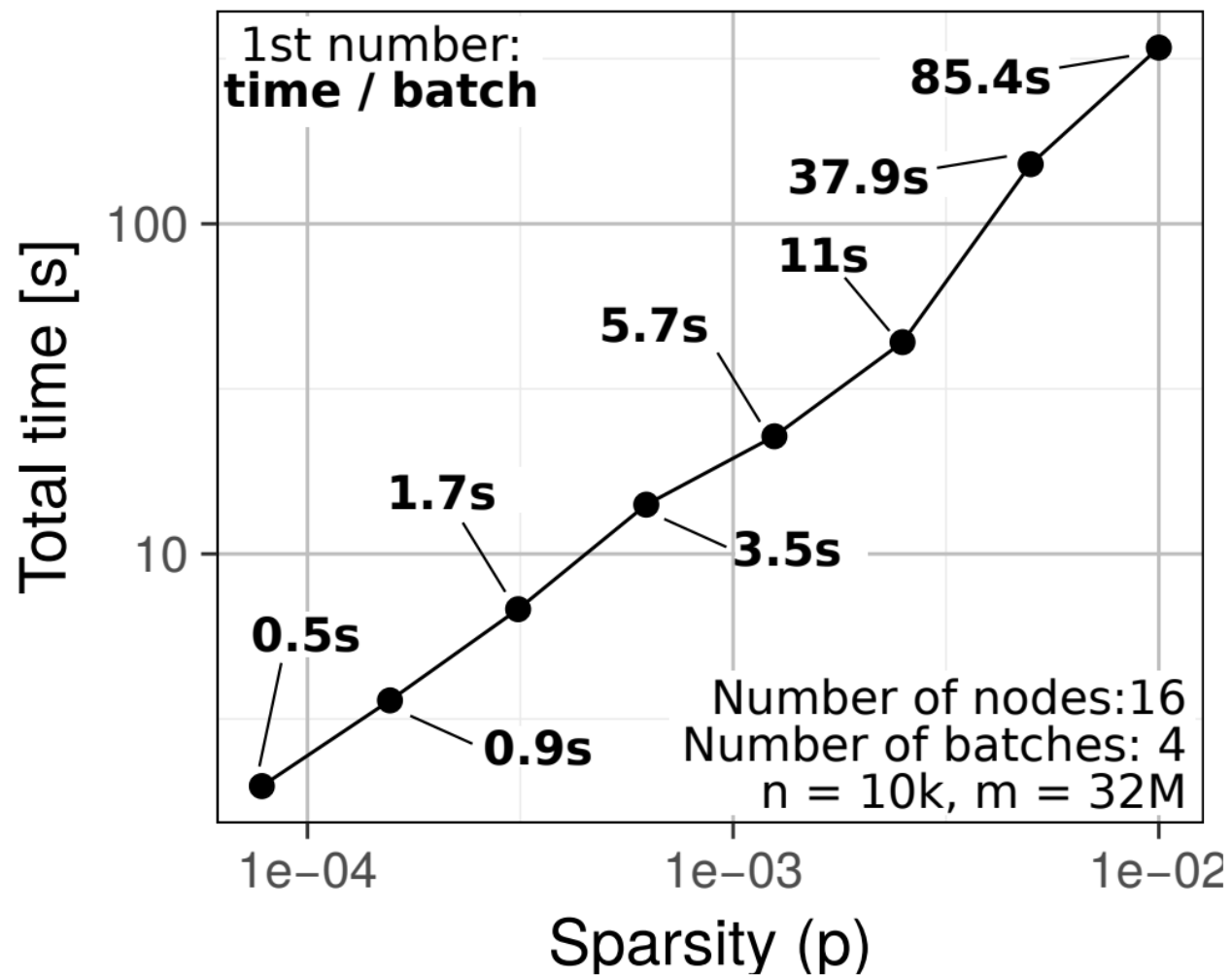
Results match theoretical predictions



PERFORMANCE ANALYSIS: SYNTHETIC DATA, SPARSITY ANALYSIS

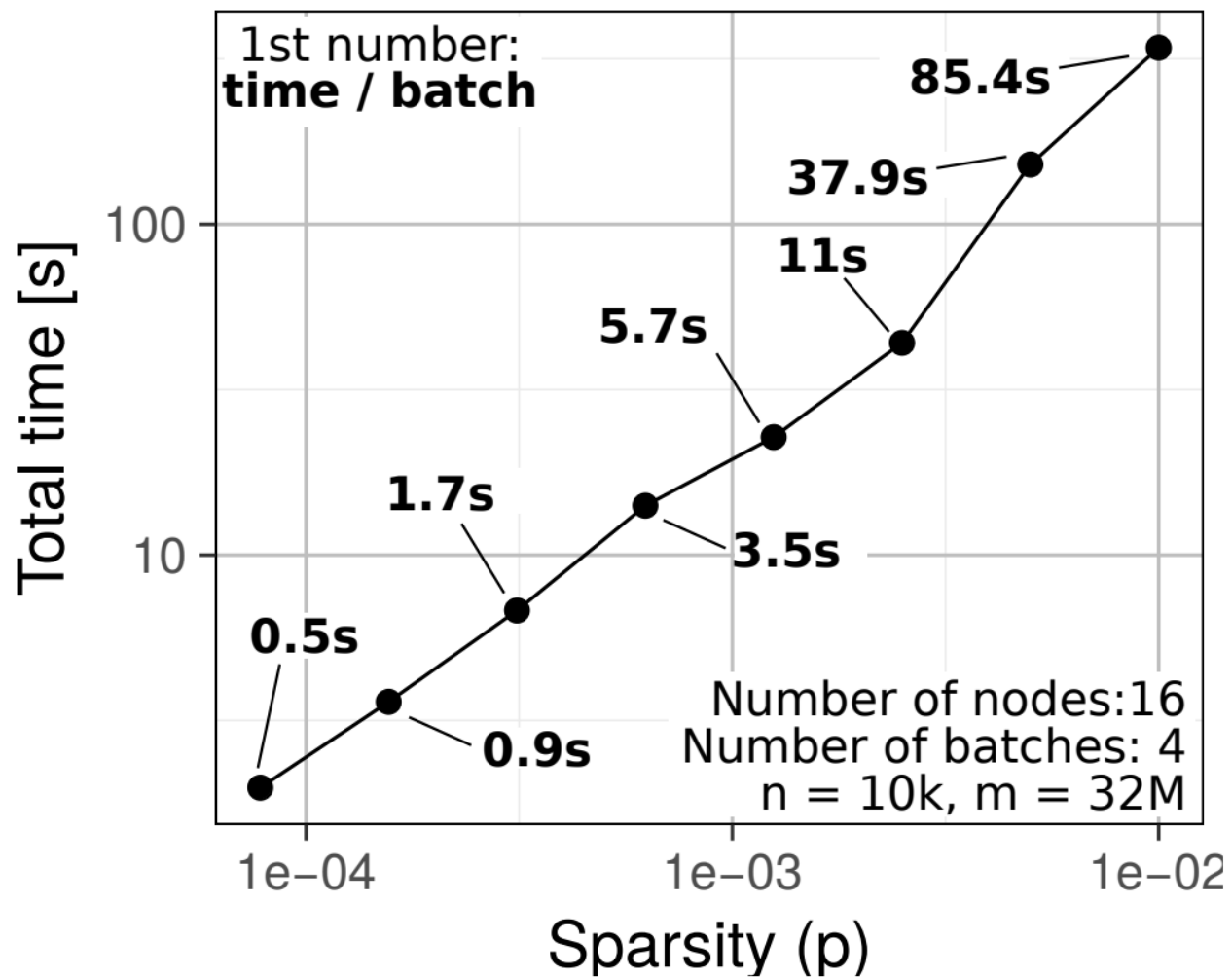


PERFORMANCE ANALYSIS: SYNTHETIC DATA, SPARSITY ANALYSIS



Sparsity (p) corresponds to the probability of the occurrence of a particular k-mer

PERFORMANCE ANALYSIS: SYNTHETIC DATA, SPARSITY ANALYSIS



Sparsity (p) corresponds to the probability of the occurrence of a particular k -mer

Nearly ideal scaling of the total runtime with the decreasing data sparsity