

Comparing and indexing metagenomes at a large scale using random projections

Roland Faure¹, Hasin Abrar², Haonan Wu², Stephanie Won², Adrita Hossain Nakshi², Rayan Chikhi¹, David Koslicki², Paul Medvedev²

¹Institut Pasteur

²The Pennsylvania State University

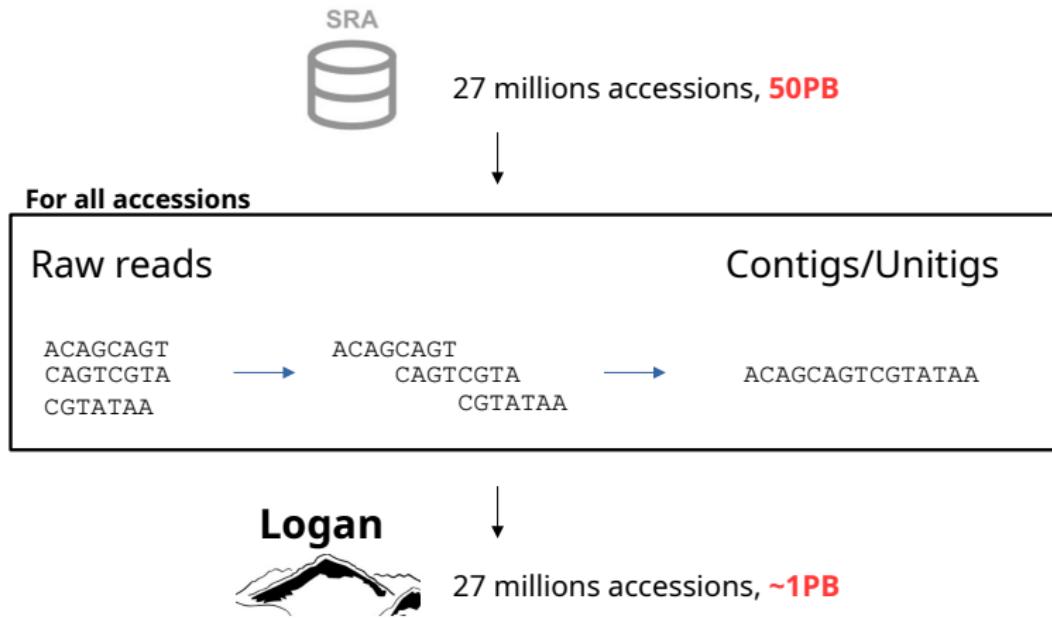
SeqBIM 2025

The SRA database

SRA: All public sequencing reads, **50 PBases** (as of Dec 2023)

Slide Credits: Teo Lemane

The Logan project



Our problem: comparing all metagenomes



David Koslicki
September 2025

I want do an all-vs-all comparison
of all 5M metagenomes of Logan!

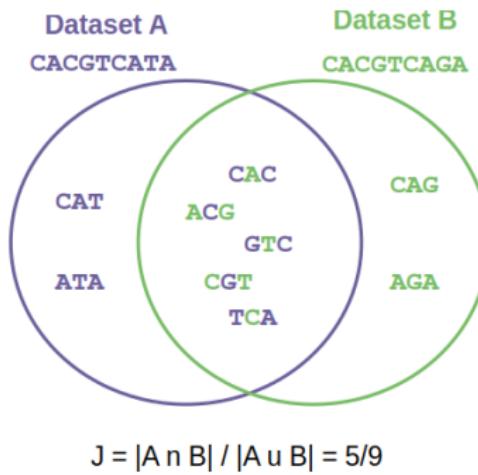
- ▶ Phase 1: Compute all-vs-all distance between metagenomes using Logan's unitigs
- ▶ Phase 2: Find patterns in metagenomes

It's a big problem!

- ▶ 5M metagenomes: 100TB sequences
- ▶ All-vs-all means $5M \times 5M = 25,000$ billion comparisons

It's a big problem!

- ▶ 5M metagenomes: 100TB sequences
- ▶ All-vs-all means $5M \times 5M = 25,000$ billion comparisons
- ▶ Use Jaccard index as a distance



First try: sourmash



David Koslicki
September 2025

Let's try with Sourmash

- ▶ Sourmash: take 1/1000 k-mer, measure Jaccard on this sketch

First try: sourmash



David Koslicki
September 2025

Let's try with Sourmash

- ▶ Sourmash: take 1/1000 k-mer, measure Jaccard on this sketch
- ▶ Estimated total time it will take: 15 CPU.years

First try: sourmash



- ▶ Sourmash: take 1/1000 k-mer, measure Jaccard on this sketch
- ▶ Estimated total time it will take: 15 CPU.years

First try: sourmash



- ▶ Sourmash: take 1/1000 k-mer, measure Jaccard on this sketch
- ▶ Estimated total time it will take: 15 CPU.years
- ▶ The problem: even taking 1/1000 k-mer, sketches get big

First try: sourmash



- ▶ Sourmash: take 1/1000 k-mer, measure Jaccard on this sketch
- ▶ Estimated total time it will take: 15 CPU.years
- ▶ The problem: even taking 1/1000 k-mer, sketches get big

Huge all-vs-all Jaccard computation: existing methods

- ▶ Based on **fixed-size** sketches

Huge all-vs-all Jaccard computation: existing methods

- ▶ Based on **fixed-size** sketches
- ▶ MinHash (e.g. Mash)
- ▶ HyperLogLog (e.g. Dashing2)
- ▶ DotHash / Random projections (e.g. Hypergen)

Comparing and indexing metagenomes at a large scale using random projections

Roland Faure¹, Hasin Abrar², Haonan Wu², Stephanie Won², Adrita Hossain Nakshi², Rayan Chikhi¹, David Koslicki², Paul Medvedev²

¹Institut Pasteur

²The Pennsylvania State University

SeqBIM 2025

DotHash/random projections: idea

Set of k-mers

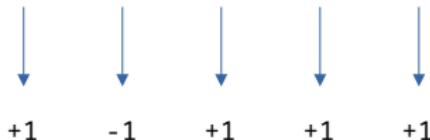
CAC, AGC, TTT, GCA, CAT

DotHash/random projections: idea

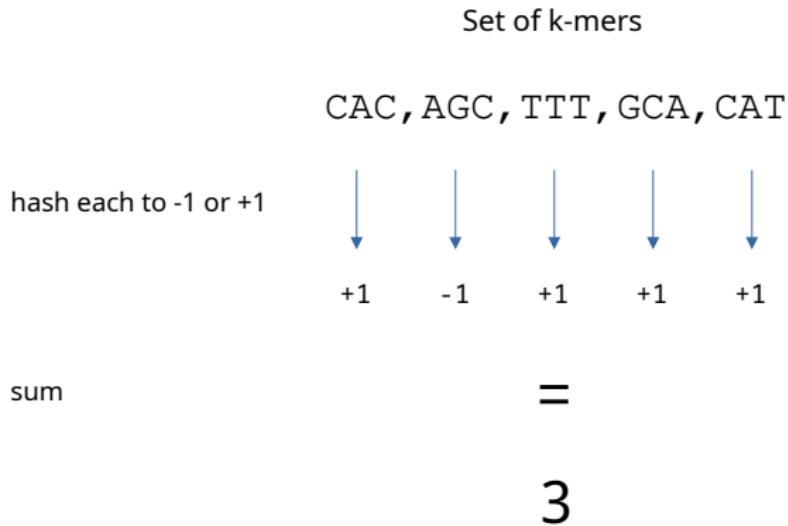
Set of k-mers

CAC, AGC, TTT, GCA, CAT

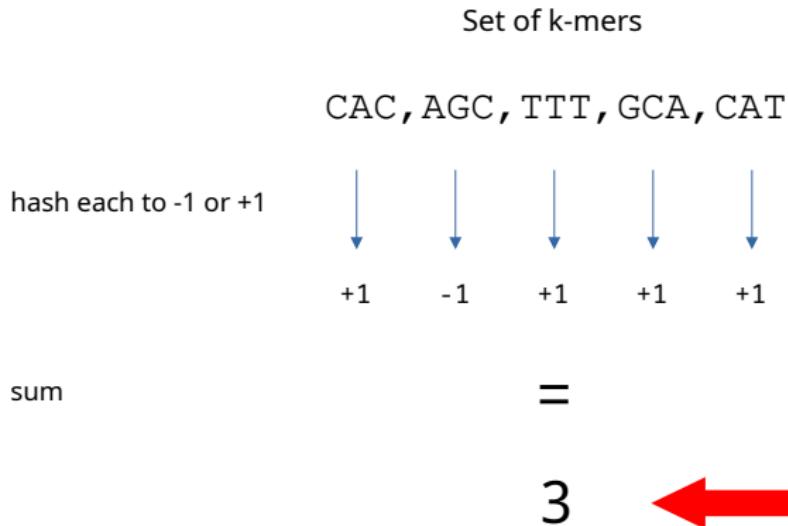
hash each to -1 or +1



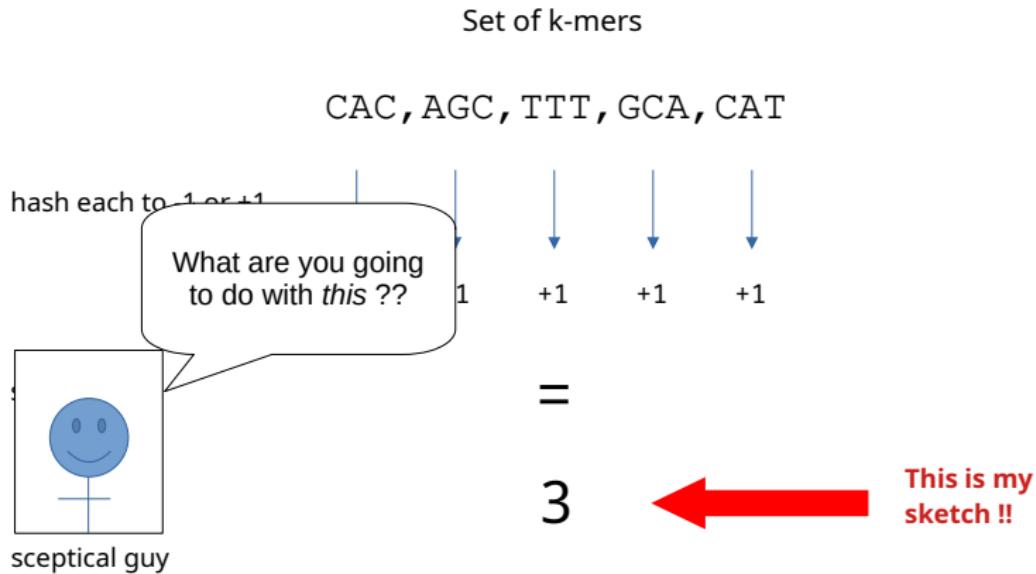
DotHash/random projections: idea



DotHash/random projections: idea



DotHash/random projections: idea



DotHash/random projections: idea

Human gut sequencing

[ERX14244107: PromethION sequencing](#)

1 OXFORD_NANOPORE (PromethION) run: 3.7M spots, 30.2G bases, 24.4Gb downloads

Submitted by: STANFORD UNIVERSITY SCHOOL OF MEDICINE

Study: Long read metagenomics reveals phage dynamics in the human gut microbiome

[PRJEB88320](#) • [ERP171428](#) • All experiments • All runs
[show Abstract](#)

Sample: D04_T2

[SAMEA117989194](#) • [ERS23983071](#) • All experiments • All runs
[Organism: human gut metagenome](#)

Soil sequencing

[SRX3026408B: Nanopore shotgun metagenomic sequencing of R2 pool](#)

1 OXFORD_NANOPORE (MinION) run: 1.8M spots, 7.8G bases, 6.3Gb downloads

Design: To simultaneously isolate RNA and DNA from soil samples, the PowerSoil total RNA kit (Qiagen) and PowerSoil DNA kit was used following the manufacturer's recommendations with minor modifications. Briefly, 3 g of soil were used as an initial sample. Mechanical lysis was performed with a mixture of glass beads and bead solution that was processed for three pulses of 30 s at 4000 rpm in a Mini Beadbeater (QIAGEN). After lysis, 100 µL of RLT lysis buffer (Qiagen) was added to each tube and sent to Nanopore Technologies (ONT). Library preparation was performed using genomic DNA that had been previously fragmented to a g-TUBE (Covaris), and subsequently processed with the Ligation Sequencing Kit SQK-LSK109 (Oxford Nanopore Technologies) with the manufacturer's protocol. Sequencing was carried out on a MinION device (version 1B) (Oxford Nanopore Technologies) with Basecalling was performed using Guppy v3.5.7 (Oxford Nanopore Technologies).

Submitted by: Instituto de Nutrición y Tecnología de los Alimentos (INTA)

Study: Metagenome sequencing of soil samples surrounding Pappostipa frigida plants from Paso de Jama, Atacama, Chile
[PRJNA1213082](#) • [SRP561469](#) • All experiments • All runs
[show Abstract](#)



sketch



DotHash/random projections: idea

Human gut sequencing

[ERX14244107: PromethION sequencing](#)

1 OXFORD_NANOPORE (PromethION) run: 3.7M spots, 30.2G bases, 24.4Gb downloads

Submitted by: STANFORD UNIVERSITY SCHOOL OF MEDICINE

Study: Long read metagenomics reveals phage dynamics in the human gut microbiome

[PRJEB88320](#) • [ERP171428](#) • All experiments • All runs

[show Abstract](#)

Sample: D04_T2

[SAMEA117989194](#) • [ERS23983071](#) • All experiments • All runs

Organism: [human gut metagenome](#)

Soil sequencing

[SRX3026408B: Nanopore shotgun metagenomic sequencing of R2 pool](#)

1 OXFORD_NANOPORE (MinION) run: 1.8M spots, 7.8G bases, 6.3Gb downloads

Design: To simultaneously isolate RNA and DNA from soil samples, the PowerSoil total RNA kit (Qiagen) and PowerSoil DNA kit was used following the manufacturer's recommendations with minor modifications. Briefly, 3 g of soil were used as an initial sample. Mechanical lysis was performed with a mixture of glass beads and bead solution that was processed for three pulses of 30 s at 4000 rpm in a MiniMAG (QIAGEN). After lysis, 100 µL of RLT lysis buffer (Qiagen) was added to each tube and sent to Nanopore Technologies (ONT). Library preparation was performed using genomic DNA that had been previously fragmented to with a g-TUBE (Covaris), and subsequently processed with the Ligation Sequencing Kit SQK-LSK109 (Oxford Nanopore Technologies) with the manufacturer's protocol. Sequencing was carried out on a MinION device (version 1B) (Oxford Nanopore Technologies) with Basecalling was performed using Guppy v5.7 (Oxford Nanopore Technologies).

Submitted by: Instituto de Nutrición y Tecnología de los Alimentos (INTA)

Study: Metagenome sequencing of soil samples surrounding Papposipa frigida plants from Paso de Jama, Atacama, Chile

[PRJNA1213082](#) • [SRP561469](#) • All experiments • All runs

[show Abstract](#)



intersection size: $-276 \times -98 = 27048$ k-mers

DotHash/random projections: idea

Human gut sequencing

[ERX14244107: PromethION sequencing](#)

1 OXFORD_NANOPORE (PromethION) run: 3.7M spots, 30.2G bases, 24.4Gb downloads

Submitted by: STANFORD UNIVERSITY SCHOOL OF MEDICINE

Study: Long read metagenomics reveals phage dynamics in the human gut microbiome

[PRJEB88320](#) • [ERP171428](#) • All experiments • All runs

[show Abstract](#)

Sample: D04_T2

[SAMEA117989194](#) • [ERS23983071](#) • All experiments • All runs

Organism: [human gut metagenome](#)

Soil sequencing

[SRX3026408B: Nanopore shotgun metagenomic sequencing of R2 pool](#)

1 OXFORD_NANOPORE (MinION) run: 1.8M spots, 7.6G bases, 6.3Gb downloads

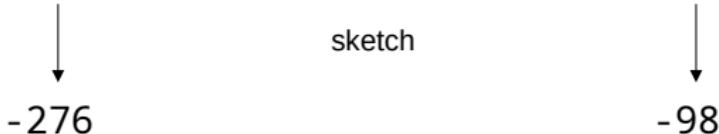
Design: To simultaneously isolate RNA and DNA from soil samples, the PowerSoil total RNA kit (Qiagen) and PowerSoil DNA kit was used following the manufacturer's recommendations with minor modifications. Briefly, 3 g of soil were used as an initial sample. Mechanical lysis was performed with a mixture of glass beads and bead solution that was processed for three pulses of 30 s at 4000 rpm in a MiniMAG (QIAGEN). After lysis, 100 µL of RLT lysis buffer (Qiagen) was added to each tube and sent to Nanopore Technologies (ONT). Library preparation was performed using genomic DNA that had been previously fragmented to with a g-TUBE (Covaris), and subsequently processed with the Ligation Sequencing Kit SQK-LSK109 (Oxford Nanopore Technologies) with the manufacturer's protocol. Sequencing was carried out on a MinION device (version 1B) (Oxford Nanopore Technologies) with Basecalling was performed using Guppy v5.7 (Oxford Nanopore Technologies).

Submitted by: Instituto de Nutrición y Tecnología de los Alimentos (INTA)

Study: Metagenome sequencing of soil samples surrounding Papposipa frigida plants from Paso de Jama, Atacama, Chile

[PRJNA123082](#) • [SRP561449](#) • All experiments • All runs

[show Abstract](#)



intersection size: $-276 \times -98 = 27048$ k-mers

unbiased estimator

DotHash/random projections: idea

Human gut sequencing

[ERX14244107: PromethION sequencing](#)

1 OXFORD_NANOPORE (PromethION) run: 3.7M spots, 30.2G bases, 24.4Gb downloads

Submitted by: STANFORD UNIVERSITY SCHOOL OF MEDICINE

Study: Long read metagenomics reveals phage dynamics in the human gut microbiome

[PRJEB88320](#) • [ERP171428](#) • All experiments • All runs
[show Abstract](#)

Sample: D04_T2

[SAMEA117989194](#) • [ERS23983071](#) • All experiments • All runs
[Organism: human gut metagenome](#)

↓
- 276

sketch

intersection size: $-276 \times -98 = 27048$ k-mers

unbiased estimator

No way!



Rayan Chikhi,
November 2025

DotHash/random projections: idea explanation

AAA, CCC, TTT**+1 -1 -1****=****-1****AAA, CCC, ATA, CGA****+1 -1 -1 -1****=****-2**

$$-1 \times -2 = (+1 - 1 - 1) \times (+1 - 1 - 1 - 1)$$

DotHash/random projections: idea explanation

AAA, CCC, TTT**+1 -1 -1****=****-1****AAA, CCC, ATA, CGA****+1 -1 -1 -1****=****-2**

$$-1 \times -2 = (+1 - 1 - 1) \times (+1 - 1 - 1 - 1)$$

$$= 1 \times 1 + -1 \times -1 + 1 \times -1 + 1 \times -1 + \dots + -1 \times -1$$

DotHash/random projections: idea explanation

AAA, CCC, TTT

=

-1

AAA, CCC, ATA, CGA

=

-2

$$-1 \times -2 = (+1-1-1) \times (+1-1-1-1)$$

$$= \underbrace{1 \times 1}_{\text{ }} + \underbrace{-1 \times -1}_{\text{ }} + \underbrace{1 \times -1}_{\text{ }} + \underbrace{1 \times -1}_{\text{ }} + \underbrace{1 \times -1}_{\text{ }} + \dots + \underbrace{-1 \times -1}_{\text{ }}$$

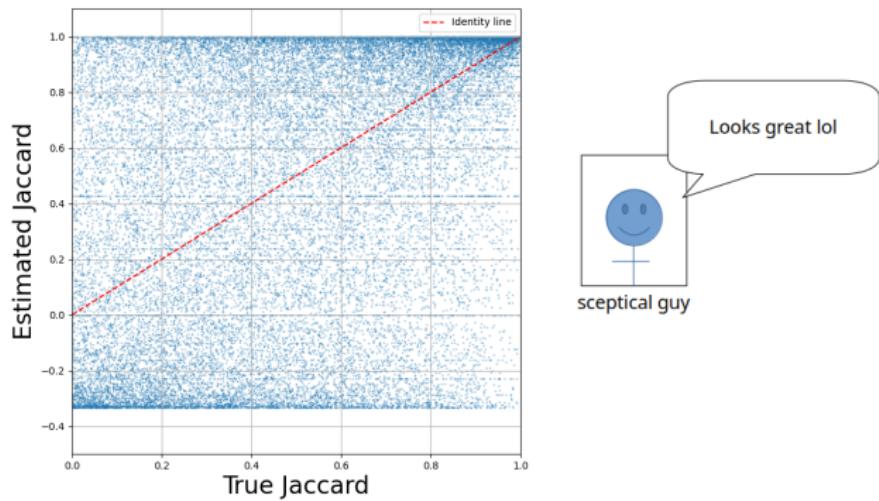
$$= |A \cap B| + \text{sum of random } -1 \text{ and } +1, \text{ on average } 0$$

DotHash: in practice

- ▶ Let's try to compute the Jaccard of 30k pairs of sets using this method

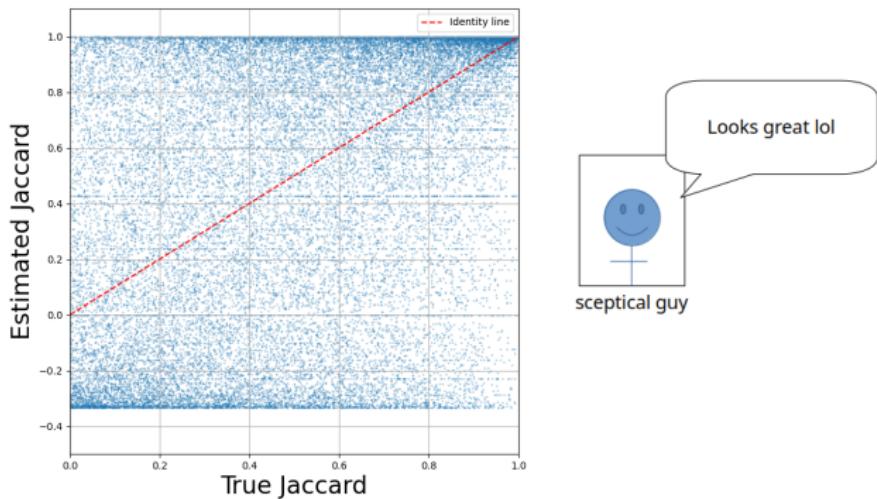
DotHash: in practice

- ▶ Let's try to compute the Jaccard of 30k pairs of sets using this method



DotHash: in practice

- ▶ Let's try to compute the Jaccard of 30k pairs of sets using this method



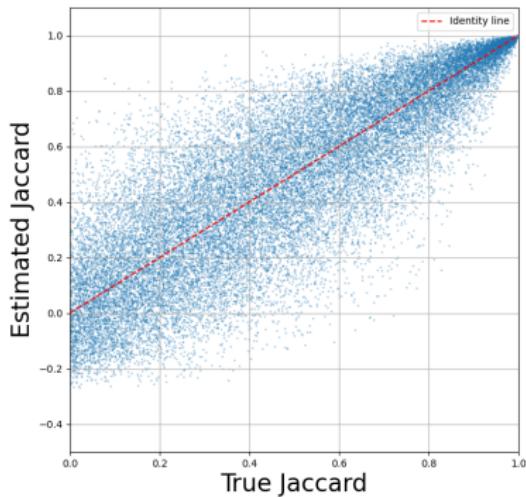
- ▶ It works on average, right?

DotHash: in practice

- ▶ For each pair of datasets, using 10 different hash functions and taking the mean

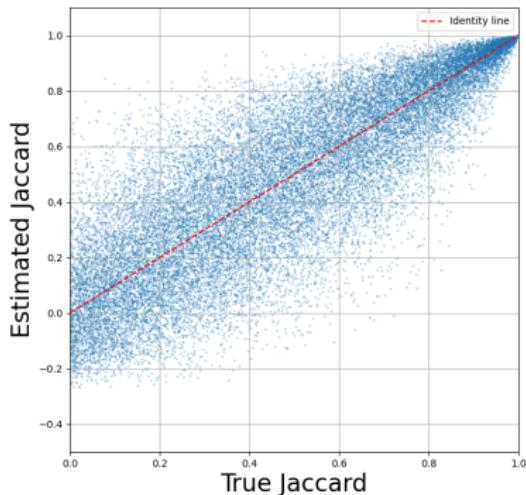
DotHash: in practice

- ▶ For each pair of datasets, using 10 different hash functions and taking the mean



DotHash: in practice

- ▶ For each pair of datasets, using 10 different hash functions and taking the mean



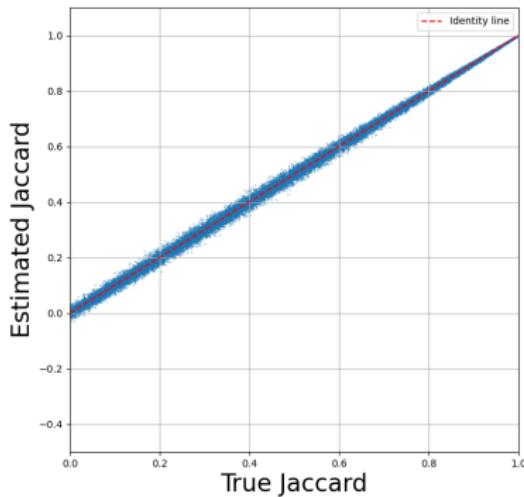
- ▶ It works better!

DotHash: in practice

- ▶ For each pair of datasets, using 2048 different hash functions and taking the mean

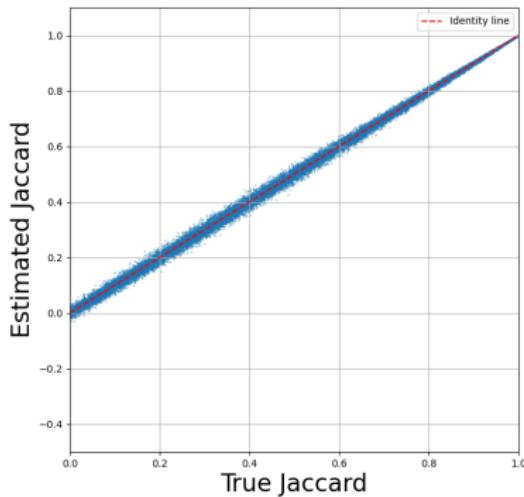
DotHash: in practice

- ▶ For each pair of datasets, using 2048 different hash functions and taking the mean



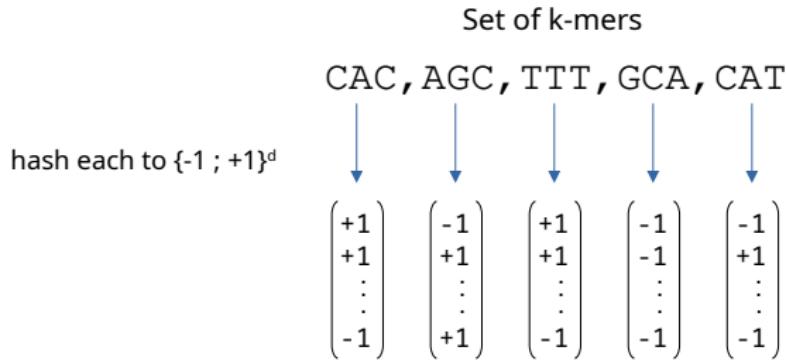
DotHash: in practice

- ▶ For each pair of datasets, using 2048 different hash functions and taking the mean



- ▶ It works!

DotHash: the full method



DotHash: the full method

Set of k-mers

CAC, AGC, TTT, GCA, CAT

hash each to $\{-1; +1\}^d$

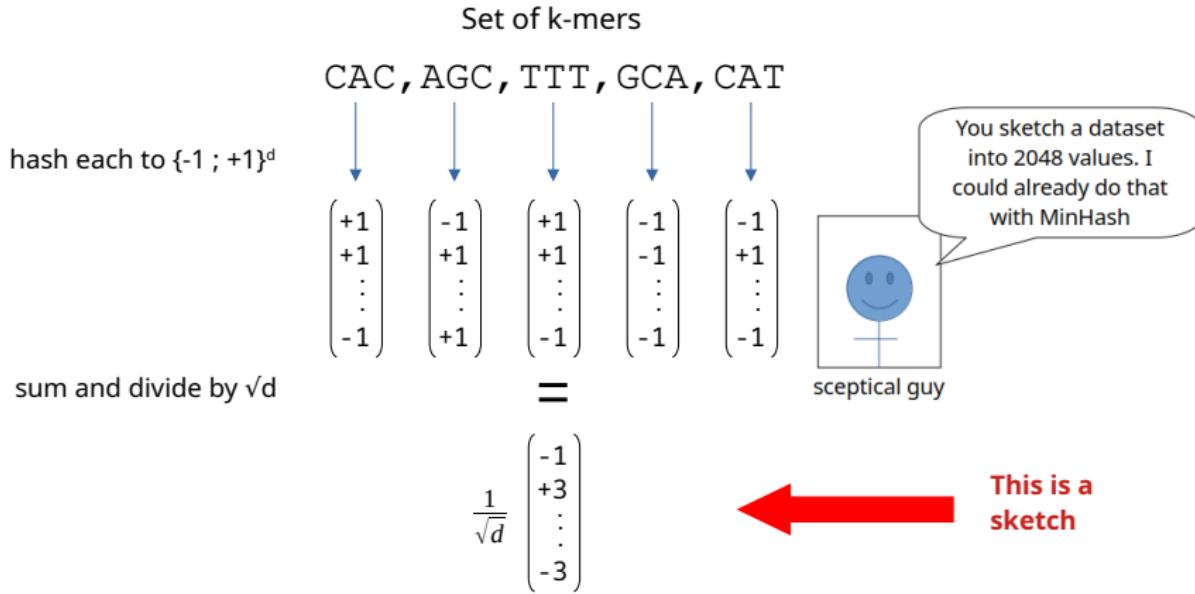
$$\begin{array}{c} \downarrow \\ (+1) \\ +1 \\ \vdots \\ -1 \end{array} \quad \begin{array}{c} \downarrow \\ (-1) \\ +1 \\ \vdots \\ +1 \end{array} \quad \begin{array}{c} \downarrow \\ (+1) \\ +1 \\ \vdots \\ -1 \end{array} \quad \begin{array}{c} \downarrow \\ (-1) \\ -1 \\ \vdots \\ -1 \end{array} \quad \begin{array}{c} \downarrow \\ (-1) \\ +1 \\ \vdots \\ -1 \end{array}$$

sum and divide by \sqrt{d}

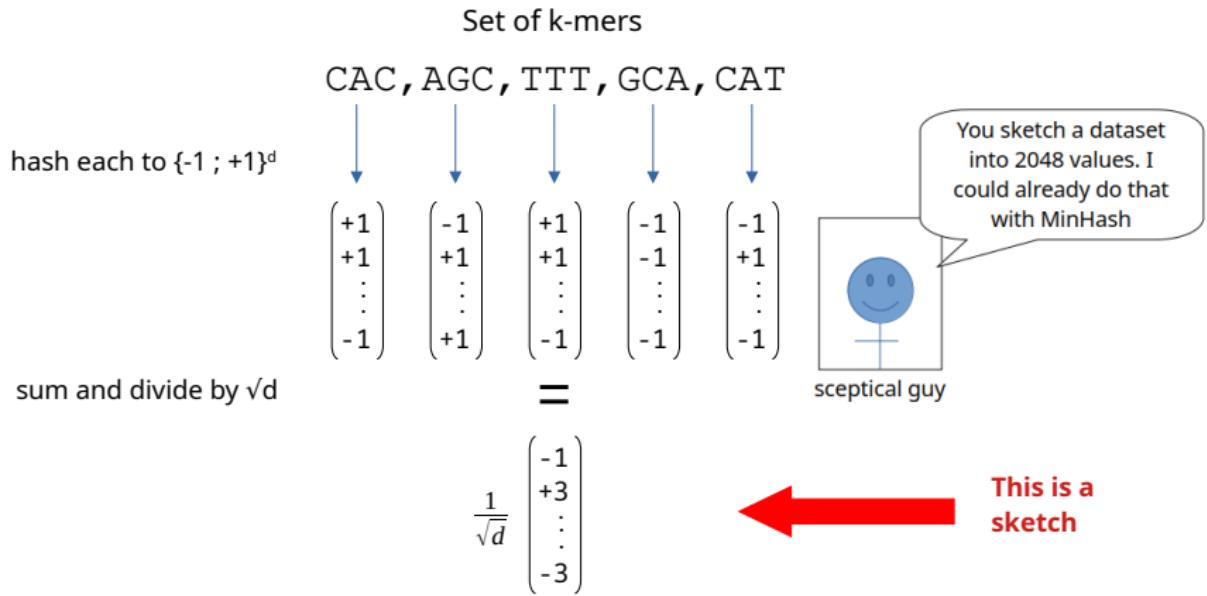
$$\equiv$$
$$\frac{1}{\sqrt{d}} \begin{pmatrix} -1 \\ +3 \\ \vdots \\ -3 \end{pmatrix}$$

This is a sketch

DotHash: the full method



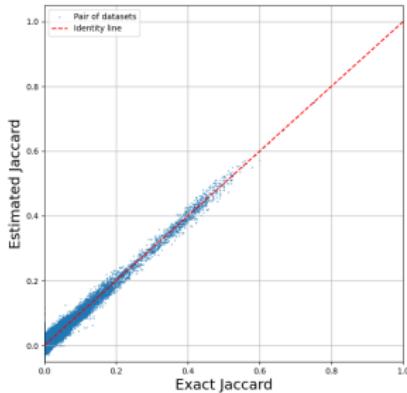
DotHash: the full method



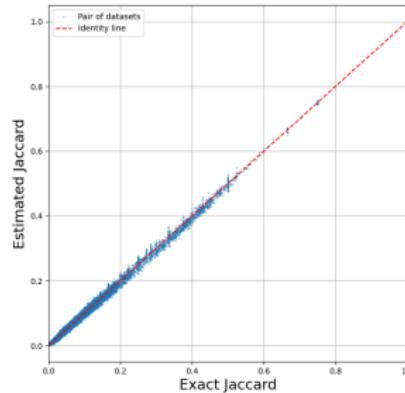
- ▶ Let's benchmark the time for all-vs-all jaccard computation

Comparison with other sketching techniques

- ▶ True vs estimated Jaccard on 35k real datasets



(a) DotHash ($d=1800$)

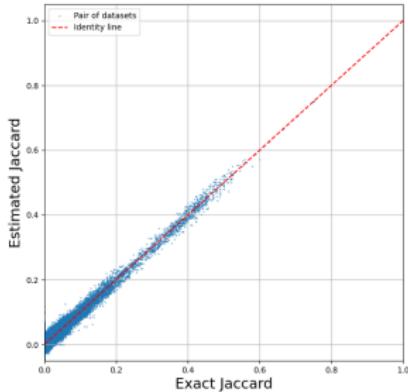


(b) HyperLogLog ($S=2048$)

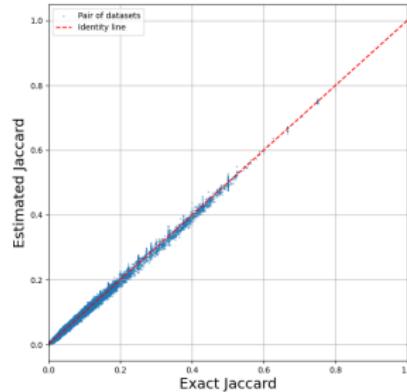
- ▶ Different sketching techniques have different error patterns

Comparison with other sketching techniques

- ▶ True vs estimated Jaccard on 35k real datasets



(a) DotHash ($d=1800$)

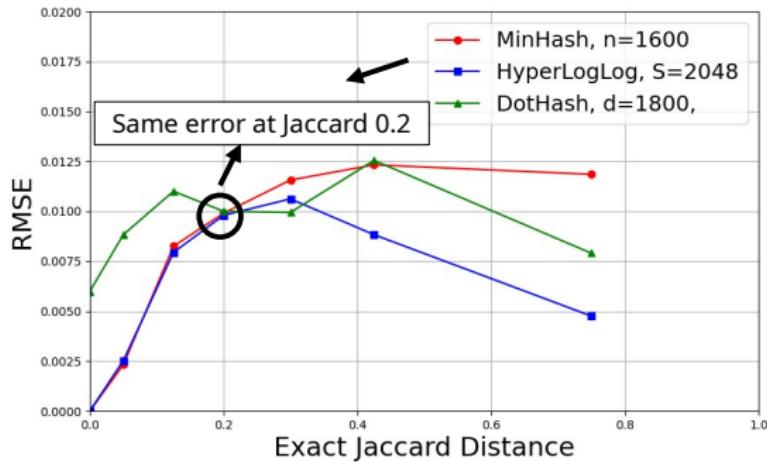


(b) HyperLogLog ($S=2048$)

- ▶ Different sketching techniques have different error patterns
- ▶ Let's calibrate the benchmark

Calibrating the benchmark

- We measured the error on 35k real datasets



- Let's benchmark times & memory

Time benchmark of the sketching techniques

- ▶ Benchmark on 35k real metagenomes
- ▶ Extrapolate the time it would take for the 5M metagenomes and for all Logan

Time benchmark of the sketching techniques

- ▶ Benchmark on 35k real metagenomes
- ▶ Extrapolate the time it would take for the 5M metagenomes and for all Logan

Method	35k genomes	5M genomes (extrapolated)	Logan (extrapolated)
Sourmash	8 CPU.hours	-	-

Time benchmark of the sketching techniques

- ▶ Benchmark on 35k real metagenomes
- ▶ Extrapolate the time it would take for the 5M metagenomes and for all Logan

Method	35k genomes	5M genomes (extrapolated)	Logan (extrapolated)
Sourmash	8 CPU.hours	-	-
Mash	4.7 CPU.hours	10 CPU.years	322 CPU.years

Time benchmark of the sketching techniques

- ▶ Benchmark on 35k real metagenomes
- ▶ Extrapolate the time it would take for the 5M metagenomes and for all Logan

Method	35k genomes	5M genomes (extrapolated)	Logan (extrapolated)
Sourmash	8 CPU.hours	-	-
Mash	4.7 CPU.hours	10 CPU.years	322 CPU.years
Dashing2	12 CPU.minutes	139 CPU.days	13 CPU.years

Time benchmark of the sketching techniques

- ▶ Benchmark on 35k real metagenomes
- ▶ Extrapolate the time it would take for the 5M metagenomes and for all Logan

Method	35k genomes	5M genomes (extrapolated)	Logan (extrapolated)
Sourmash	8 CPU.hours	-	-
Mash	4.7 CPU.hours	10 CPU.years	322 CPU.years
Dashing2	12 CPU.minutes	139 CPU.days	13 CPU.years
HyperGen	13 CPU.minutes	154 CPU.days	15 CPU.years

Time benchmark of the sketching techniques

- ▶ Benchmark on 35k real metagenomes
- ▶ Extrapolate the time it would take for the 5M metagenomes and for all Logan

Method	35k genomes	5M genomes (extrapolated)	Logan (extrapolated)
Sourmash	8 CPU.hours	-	-
Mash	4.7 CPU.hours	10 CPU.years	322 CPU.years
Dashing2	12 CPU.minutes	139 CPU.days	13 CPU.years
HyperGen	13 CPU.minutes	154 CPU.days	15 CPU.years
Custom DotHash	2 CPU.minutes	24 CPU.days	2.3 CPU.years

The strong point of DotHash: practical implementation

- ▶ Comparing datasets is a matrix multiplication

Sketch of 3 datasets

$$\frac{1}{\sqrt{d}} \begin{pmatrix} 23 & 45 & -23 & \dots & -5 \\ -81 & 68 & 111 & \dots & -260 \\ 9 & -4 & 1 & \dots & -10 \end{pmatrix} \times \frac{1}{\sqrt{d}} \begin{pmatrix} 23 & -81 & 9 \\ 45 & 68 & -4 \\ -23 & 111 & 1 \\ \vdots & \ddots & \ddots \\ -5 & -260 & -10 \end{pmatrix} = \begin{pmatrix} 78710.2 & 3245.1 & 292.2 \\ 3245.1 & 97550.3 & -21.2 \\ 292.2 & -21.2 & 6345.2 \end{pmatrix}$$

Size of intersection of datasets 1 and 2

The diagram illustrates the DotHash sketching process. It shows three datasets represented as vectors of length d . The first vector has entries 23, 45, -23, ..., -5. The second vector has entries -81, 68, 111, ..., -260. The third vector has entries 9, -4, 1, ..., -10. These vectors are multiplied by a scalar $\frac{1}{\sqrt{d}}$ to produce a sketch. The sketch is a 3x3 matrix where each entry is the dot product of two vectors. The diagonal entries are 78710.2, 3245.1, and 292.2. The (1,2) entry is circled in red and highlighted with an arrow pointing to it, labeled "Size of intersection of datasets 1 and 2".

- ▶ Highly hardware-optimizable, e.g. SIMD, GPU

Back to our original problem

- ▶ Comparing all-vs-all jaccard of 5M metagenomes: done in one night!



Let's move on to analyse how these datasets are organized!

David Koslicki
November 2025

Back to our original problem

- ▶ Comparing all-vs-all jaccard of 5M metagenomes: done in one night!



Let's move on to analyse how these datasets are organized!

David Koslicki
November 2025

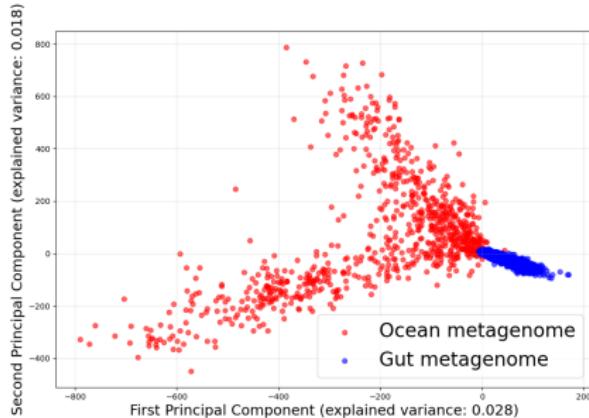
- ▶ DotHash vectors can still be interesting!

Manipulating vectors: analysis

- We can run a PCA directly on the sketches!

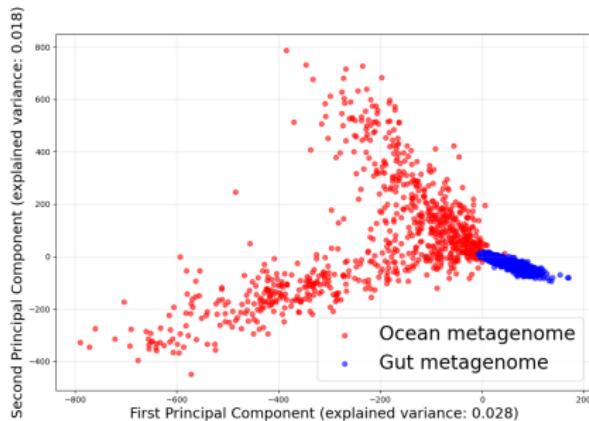
Manipulating vectors: analysis

- ▶ We can run a PCA directly on the sketches!
- ▶ Example: 7101 metagenomes either from Tara Ocean or American Gut Project



Manipulating vectors: analysis

- ▶ We can run a PCA directly on the sketches!
- ▶ Example: 7101 metagenomes either from Tara Ocean or American Gut Project



- ▶ Also methods for clustering, machine learning, indexing, compressing...

Take-home messages

- ▶ DotHash hashes sets into vectors

Take-home messages

- ▶ DotHash hashes sets into vectors
- ▶ Very computationally efficient

Take-home messages

- ▶ DotHash hashes sets into vectors
- ▶ Very computationally efficient
- ▶ Imprecise for low jaccard distances

Take-home messages

- ▶ DotHash hashes sets into vectors
- ▶ Very computationally efficient
- ▶ Imprecise for low jaccard distances
- ▶ Many powerful methods/implementation exist to manipulate vectors

Acknowledgments



Paul
Medvedev



David
Koslicki



Stephanie
Won



Hasin Abrar



Haonan Wu



Rayan Chikhi

Space taken by the sketches

Method	35k genomes	5M genomes (extrapolated)
Mash	408M	52G
Dashing2	548M	245G
Custom DotHash	120M	15G