

A new way to handle large data: MSR sketching

Roland Faure^{1,2,3}

¹Université libre de Bruxelles (ULB) - Belgium

²Université de Rennes, IRISA - France

³Institut Pasteur, Paris - France

Penn State, September 2025

About me: postdoc, since Feb. 2025



Institut Pasteur, Paris

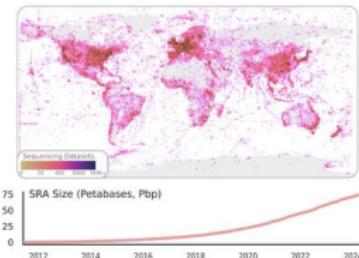


Rayan Chikhi
Institut Pasteur, Paris
Focus: massive genomics

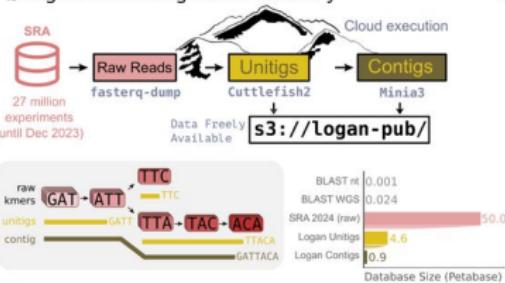
My postdoc : Index & Search the Logan database

The Logan project

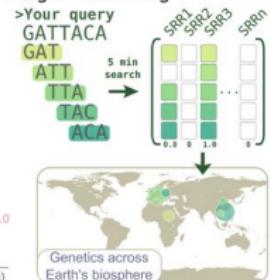
a SRA accessions



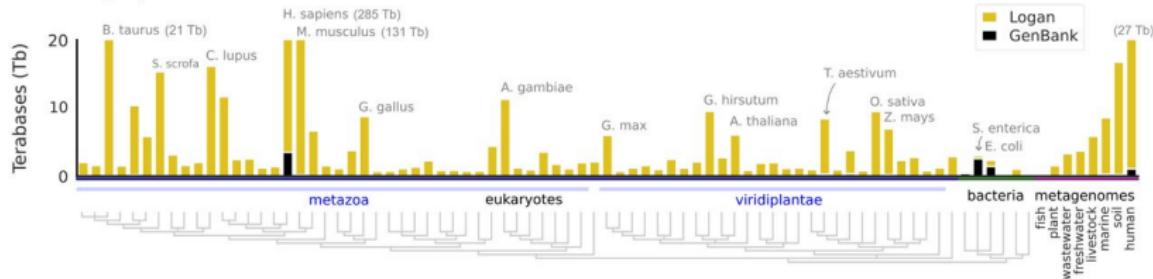
b Logan: SRA-wide genome assembly



c logan-search.org



d Assembly expansion across the tree of life

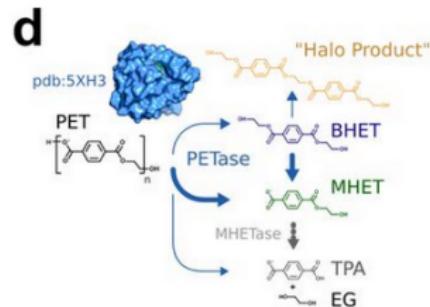


Logan: Planetary-Scale Genome Assembly Surveys Life's Diversity, *biorXiv*, 2025

The Logan project: exciting example

PAST SUCCESS • STARTUP

Plastivores: Plastic-Degrading Super-Microbes and Enzymes

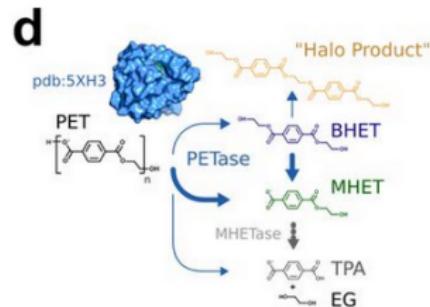


- ▶ 213 known PETases

The Logan project: exciting example

PAST SUCCESS • STARTUP

Plastivores: Plastic-Degrading Super-Microbes and Enzymes

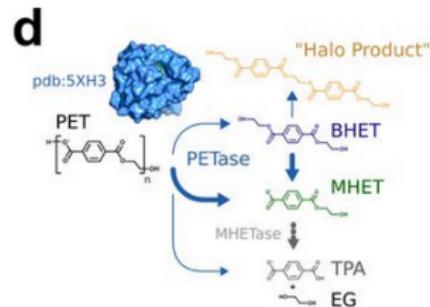


- ▶ 213 known PETases
- ▶ Let's look in Logan for homologs

The Logan project: exciting example

PAST SUCCESS • STARTUP

Plastivores: Plastic-Degrading Super-Microbes and Enzymes

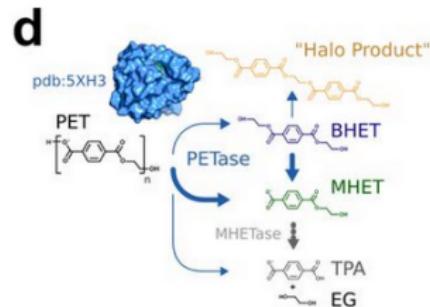


- ▶ 213 known PETases
- ▶ Let's look in Logan for homologs
- ▶ Result: 215M distinct sequences

The Logan project: exciting example

PAST SUCCESS • STARTUP

Plastivores: Plastic-Degrading Super-Microbes and Enzymes

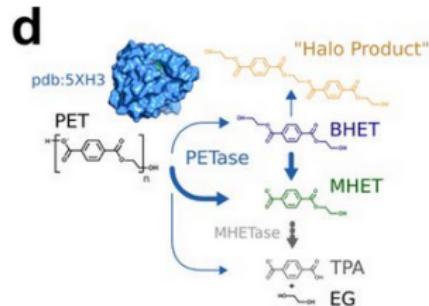


- ▶ 213 known PETases
- ▶ Let's look in Logan for homologs
- ▶ Result: 215M distinct sequences
- ▶ Some of them best than previously known PETases

The Logan project: exciting example

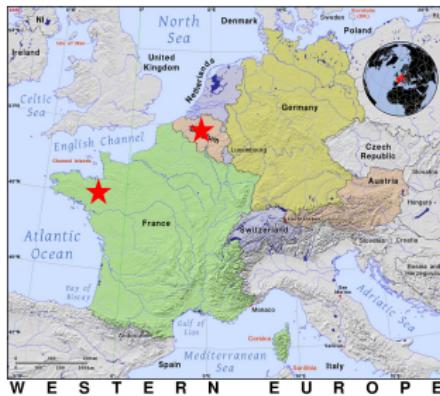
PAST SUCCESS • STARTUP

Plastivores: Plastic-Degrading Super-Microbes and Enzymes



- ▶ 213 known PETases
- ▶ Let's look in Logan for homologs
- ▶ Result: 215M distinct sequences
- ▶ Some of them best than previously known PETases
- ▶ My job: improving speed/cost of the search

About me: Ph.D., 2021-2024



Jean-François Flot

Université Libre de Bruxelles

Focus: assembling wild genomes



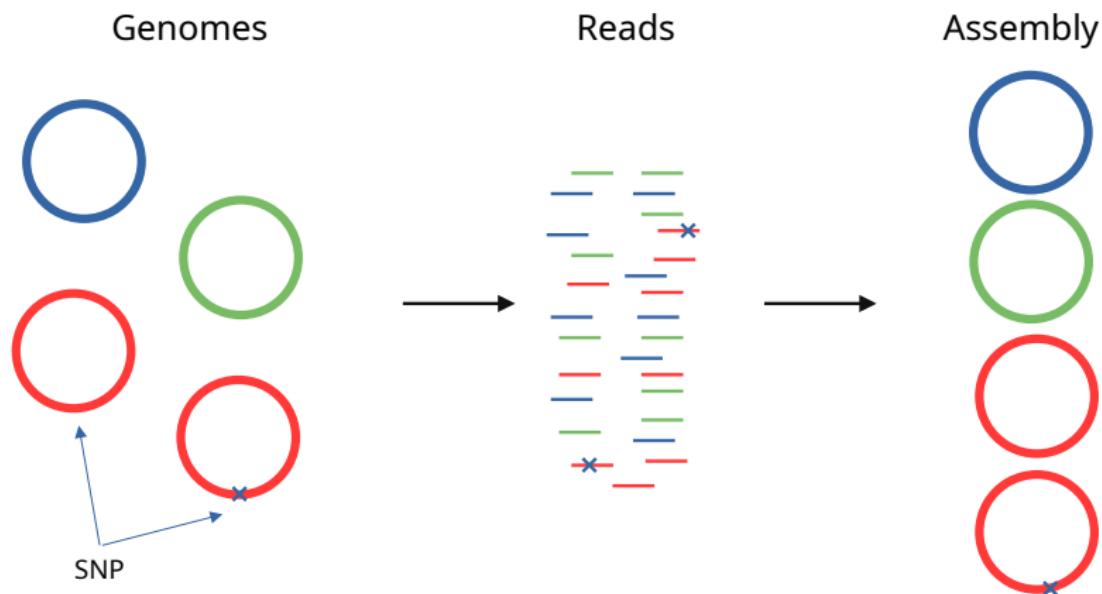
Dominique Laveneir

Université de Rennes

Focus: computational methods

My Ph.D. : Haplotype assembly from long reads

Focus of my Ph.D.: Metagenome assembly



(Meta)genome assembly is a big computation

(Meta)genome assembly is a big computation

- ▶ Assembling a human gut metagenome (HiFi, 250Gpb)

(Meta)genome assembly is a big computation

- ▶ Assembling a human gut metagenome (HiFi, 250Gpb)

metaFlye

4 days, 256GB RAM



(Meta)genome assembly is a big computation

- ▶ Assembling a human gut metagenome (HiFi, 250Gpb)

metaFlye

4 days, 256GB RAM



hifiasm_meta

11 days, 454GB RAM



(Meta)genome assembly is a big computation

- ▶ Assembling a human gut metagenome (HiFi, 250Gpb)

metaFlye

4 days, 256GB RAM



hifiasm_meta

11 days, 454GB RAM



metaMDBG

19h, 10G RAM



metaMDBG: the trick is sketching input reads

CAGAC**TACG**ATATTT**TGCT**GACTCATGCGCG**TTTG**G



k-mer subsampling

TACG

TGCT **TGCT**

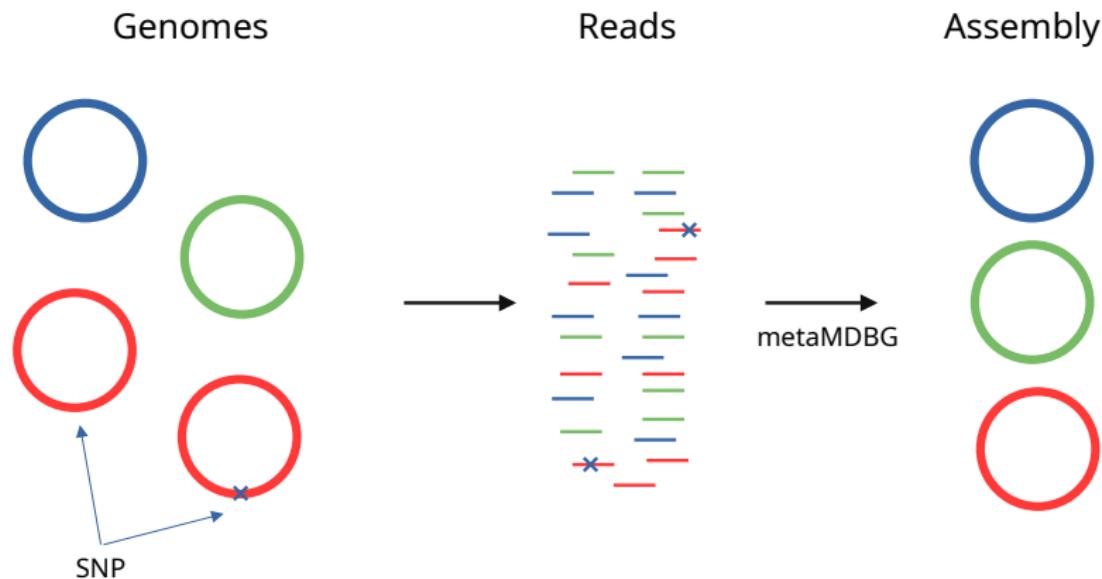


expensive computation

...

- ▶ minimizers, FracMinHash, seed-chain, strobemers...
- ▶ minimap2, Mash, BLAST, **metaMDBG**...

metaMDBG loses strain diversity



- ▶ metaMDBG is very fast, but some variants are lost!

k-mer sketching loses SNPs

SNP



CAGACTA A GATATTTTGCTGACTCAT	→	AGAC	ATTT	CTCA
CAGACTACGAT ATTTTGCTGACTCAT	→	AGAC	ATTT	CTCA

k-mer sketching loses SNPs

SNP

→

CAGACTA A GATATTTTGCTGACTCAT	→	AGAC	ATTT	CTCA
CAGACTACGAT ATTTTGCTGACTCAT	→	AGAC	ATTT	CTCA

- ▶ Is k-mer subsampling really the only way to sketch sequences ?

k-mer sketching loses SNPs

SNP

→

CAGACTA A GATATTTTGCTGACTCAT	→	AGAC	ATTT	CTCA
CAGACTACGAT ATTT TGCTGACTCAT	→	AGAC	ATTT	CTCA

- ▶ Is k-mer subsampling really the only way to sketch sequences ?
- ▶ Bassel, Luc & Medvedev, Paul & Chikhi, Rayan. (2022). *Mapping-friendly sequence reductions: Going beyond homopolymer compression.* iScience.

Generalizing Homopolymer Compression

CAT**TT**CGAGTA**AAGGGG**CAC**CT**G → CATCGAGTAGCAGCTG

Homopolymer compression

Generalizing Homopolymer Compression

Homopolymer compression

CAT**TT**CGAGTA**AAGGGG**CAC**CT**G → CATCGAGTAGCAGCTG

"Heteropolymer" compression

CATTTCGAGTA**AAGGGG**CAC**CT**G → TTAAGGGC

Generalizing Homopolymer Compression

Homopolymer compression

CAT**TT**CGAGTA**AAGGGG**CAC**CT**G → CATCGAGTAGCAGCTG

“Heteropolymer” compression

CATTTCGAGTA**AAGGGG**CAC**CT**G → TTAAGGGC

Turn As into Ts and Ts into As

CATTTCGAGTA**AAGGGG**CACCTG → CT**AAA**CGTGT**TTT**GGGG**TCCAG**

Generalizing Homopolymer Compression

compression		
Homopolymer compression	CATTCGAGTA AAGGGGCACCTG	→ CATCGAGTAGCAGTG 0.75
"Heteropolymer" compression	CATTTCGAGTA AAGGGGCACCTG	→ TTAAGGGC 0.25
Turn As into Ts and Ts into As	CATTTCGAGTA AAGGGGCACCTG	→ CTAAAACTGTGT TTTGGGGCTCCAG 1.0

Mapping-friendly Sequence Reductions: an example

$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**CAGTATGGAT**) = 0.0023

$f(\text{CAGTATGGAT}) = A$

sketch

A

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

C **AGTATGGATA** CAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**AGTATGGATA**) = 0.624

$f(\textbf{AGTATGGATA}) = \emptyset$

sketch

A

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**GTATGGATAC**) = 0.124

$f(\textbf{GTATGGATAC}) = G$

sketch

A G

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAG**TATGGATACA**GATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**TATGGATACA**) = 0.88

$f(\textbf{TATGGATACA}) = \emptyset$

sketch

A G

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGT**ATGGATACAG**ATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**ATGGATACAG**) = 0.32

$f(\textbf{ATGGATACAG}) = \emptyset$

sketch

A G

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGTA**TGGATACAGA**TGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**TGGATACAGA**) = 0.19

$f(\textbf{TGGATACAGA}) = T$

sketch

A G T

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGTAT**GGATACAGAT**GGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**GGATACAGAT**) = 0.214

$f(\textbf{GGATACAGAT}) = \emptyset$

sketch

A G T

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGTATG**GATACAGATG**GAGATATCATCGAGTAGGGCACTGTACCAGAG

$\text{hash}(\text{GATACAGATG}) = 0.678$

$f(\text{GATACAGATG}) = \emptyset$

sketch

A G T

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGTATGG**ATACAGATGG**AGATATCATCGAGTAGGGGCACTGTACCAGAG

$$\text{hash}(\textbf{ATACAGATGG}) = 0.669$$

$$f(\textbf{ATACAGATGG}) = \emptyset$$

sketch

A G T

Mapping-friendly Sequence Reductions: an example

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGCAC**TGTACCAGAG**

$$\text{hash}(\textbf{TGTACCAGAG}) = 0.06$$

$$f(\textbf{TGTACCAGAG}) = C$$

sketch

A G T

T C

C G T C

Mapping-friendly Sequence Reductions: homopolymer compression

$$f : \{A, C, G, T\}^2 \rightarrow \{A, C, G, T, \emptyset\}$$

$f(AA) \rightarrow \emptyset$	$f(CA) \rightarrow A$	$f(GA) \rightarrow A$	$f(TA) \rightarrow A$
$f(AC) \rightarrow C$	$f(CC) \rightarrow \emptyset$	$f(GC) \rightarrow C$	$f(TC) \rightarrow C$
$f(AG) \rightarrow G$	$f(CG) \rightarrow G$	$f(GG) \rightarrow \emptyset$	$f(TG) \rightarrow G$
$f(AT) \rightarrow T$	$f(CT) \rightarrow T$	$f(GT) \rightarrow T$	$f(TT) \rightarrow \emptyset$

Mapping-friendly Sequence Reductions: homopolymer compression

$$f : \{A, C, G, T\}^2 \rightarrow \{A, C, G, T, \emptyset\}$$

$f(AA) \rightarrow \emptyset$	$f(CA) \rightarrow A$	$f(GA) \rightarrow A$	$f(TA) \rightarrow A$
$f(AC) \rightarrow C$	$f(CC) \rightarrow \emptyset$	$f(GC) \rightarrow C$	$f(TC) \rightarrow C$
$f(AG) \rightarrow G$	$f(CG) \rightarrow G$	$f(GG) \rightarrow \emptyset$	$f(TG) \rightarrow G$
$f(AT) \rightarrow T$	$f(CT) \rightarrow T$	$f(GT) \rightarrow T$	$f(TT) \rightarrow \emptyset$

sequence

ACGTTG

sketch

Mapping-friendly Sequence Reductions: homopolymer compression

$$f : \{A, C, G, T\}^2 \rightarrow \{A, C, G, T, \emptyset\}$$

$f(AA) \rightarrow \emptyset$	$f(CA) \rightarrow A$	$f(GA) \rightarrow A$	$f(TA) \rightarrow A$
$f(AC) \rightarrow C$	$f(CC) \rightarrow \emptyset$	$f(GC) \rightarrow C$	$f(TC) \rightarrow C$
$f(AG) \rightarrow G$	$f(CG) \rightarrow G$	$f(GG) \rightarrow \emptyset$	$f(TG) \rightarrow G$
$f(AT) \rightarrow T$	$f(CT) \rightarrow T$	$f(GT) \rightarrow T$	$f(TT) \rightarrow \emptyset$

sequence **AC**GTTG

sketch C

Mapping-friendly Sequence Reductions: homopolymer compression

$$f : \{A, C, G, T\}^2 \rightarrow \{A, C, G, T, \emptyset\}$$

$f(AA) \rightarrow \emptyset$	$f(CA) \rightarrow A$	$f(GA) \rightarrow A$	$f(TA) \rightarrow A$
$f(AC) \rightarrow C$	$f(CC) \rightarrow \emptyset$	$f(GC) \rightarrow C$	$f(TC) \rightarrow C$
$f(AG) \rightarrow G$	$f(CG) \rightarrow G$	$f(GG) \rightarrow \emptyset$	$f(TG) \rightarrow G$
$f(AT) \rightarrow T$	$f(CT) \rightarrow T$	$f(GT) \rightarrow T$	$f(TT) \rightarrow \emptyset$

sequence ACGTTG

sketch CG

Mapping-friendly Sequence Reductions: homopolymer compression

$$f : \{A, C, G, T\}^2 \rightarrow \{A, C, G, T, \emptyset\}$$

$f(AA) \rightarrow \emptyset$	$f(CA) \rightarrow A$	$f(GA) \rightarrow A$	$f(TA) \rightarrow A$
$f(AC) \rightarrow C$	$f(CC) \rightarrow \emptyset$	$f(GC) \rightarrow C$	$f(TC) \rightarrow C$
$f(AG) \rightarrow G$	$f(CG) \rightarrow G$	$f(GG) \rightarrow \emptyset$	$f(TG) \rightarrow G$
$f(AT) \rightarrow T$	$f(CT) \rightarrow T$	$f(GT) \rightarrow T$	$f(TT) \rightarrow \emptyset$

sequence ACGTTG

sketch CGT

Mapping-friendly Sequence Reductions: homopolymer compression

$$f : \{A, C, G, T\}^2 \rightarrow \{A, C, G, T, \emptyset\}$$

$f(AA) \rightarrow \emptyset$	$f(CA) \rightarrow A$	$f(GA) \rightarrow A$	$f(TA) \rightarrow A$
$f(AC) \rightarrow C$	$f(CC) \rightarrow \emptyset$	$f(GC) \rightarrow C$	$f(TC) \rightarrow C$
$f(AG) \rightarrow G$	$f(CG) \rightarrow G$	$f(GG) \rightarrow \emptyset$	$f(TG) \rightarrow G$
$f(AT) \rightarrow T$	$f(CT) \rightarrow T$	$f(GT) \rightarrow T$	$f(TT) \rightarrow \emptyset$

sequence ACG**TT**G

sketch CGT

Mapping-friendly Sequence Reductions: homopolymer compression

$$f: \{A, C, G, T\}^2 \rightarrow \{A, C, G, T, \emptyset\}$$

$$\begin{array}{llll} f(AA) \rightarrow \emptyset & f(CA) \rightarrow A & f(GA) \rightarrow A & f(TA) \rightarrow A \\ f(AC) \rightarrow C & f(CC) \rightarrow \emptyset & f(GC) \rightarrow C & f(TC) \rightarrow C \\ f(AG) \rightarrow G & f(CG) \rightarrow G & f(GG) \rightarrow \emptyset & f(TG) \rightarrow G \\ f(AT) \rightarrow T & f(CT) \rightarrow T & f(GT) \rightarrow T & f(TT) \rightarrow \emptyset \end{array}$$

sequence ACGTTG

sketch CGT G

Mapping-friendly Sequence Reductions: definition

$l \in \mathbb{N}$

$$f: \{A, C, G, T\}^l \rightarrow \{A, C, G, T, \emptyset\}$$

$$\forall kmer \in \{A, C, G, T\}^l, rc(f(kmer)) = f(rc(kmer))$$

Mapping-friendly Sequence Reductions: key parameters

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

order (l)

 $f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$ $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$ $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$ $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$ $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

compression ratio (c)

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch

A G T

T C

C G T C

MSRs=Mapping-friendly Sequence Reductions

- ▶ MSR reductions are **mapping-friendly**

ATCATCGAGTAGGGGCACTGTACCAAGAGCGCTTTAATGTAC

CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAAGAG

A G T C C G T C

- ▶ original sequences align \iff reduced sequences align

MSRs=Mapping-friendly Sequence Reductions

- ▶ MSR reductions are **mapping-friendly**

The diagram shows two DNA sequences aligned vertically. A red rectangular box highlights a segment of the top sequence: 'C G T C'. Below the top sequence, the labels 'CC' and 'A' are positioned under the last two bases. Below the bottom sequence, the labels 'A G T T C' are positioned under the first four bases. The bottom sequence itself is: 'ATCATCGAGTAGGGGCACTGTACCAAGAGCGCTTTAATGTAC'. The aligned portion of the bottom sequence is: 'ATCATCGAGTAGGGGCACTGTACCAAGAG'. The labels 'A G T T C' are aligned under the first four bases of the bottom sequence.

- ▶ original sequences align \iff reduced sequences align
- ▶ Key property of MSRs

MSRs=Mapping-friendly Sequence Reductions

- ▶ MSR reductions are **mapping-friendly**

Diagram illustrating the mapping-friendliness of MSR reductions. Two DNA sequences are shown:

Top sequence: ATCATCGAGTAGGGGCACTGTACCAGAGCGCTTTAATGTAC

Bottom sequence: CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

A red box highlights a 4-base window from positions 5 to 8 of the top sequence (TGTC) and its corresponding 4-base window in the bottom sequence (CAGA).

Below the top sequence, the labels A G T C are aligned under the first four bases, and C G T C are aligned under the next four bases.

- ▶ original sequences align \iff reduced sequences align
- ▶ Key property of MSRs
- ▶ Let's compress massively and assemble!

Assembling using MSR sketches

AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
GAGATATCATCGAGTAGGGGCACTGTACCAGAGCCGG
GATATCATCGAGTAGGGGCACTGTACCAGAGGCCGGTATAC

MSR sketching

AGTTCCGT

TCCGTCAA

CGTCAATG

Assembly

AGTTCCGT
TCCGTCAA
CGTCAATG
AGTTCCGTCAAATG

Assembling using MSR sketches

AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
GAGATATCATCGAGTAGGGGCACTGTACCAGAGCCGG
GATATCATCGAGTAGGGGCACTGTACCAGAGCCGGTTATAC

MSR sketching

AGTTCCGT

TCCGTCAA

CGTCAATG

Assembly

AGTTCCGT
TCCGTCAA
CGTCAATG
AGTTCCGTCAAATG

Inflating | Inverse sketching ??

AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAGCCGGTTATAC

Inflating a reduced assembly

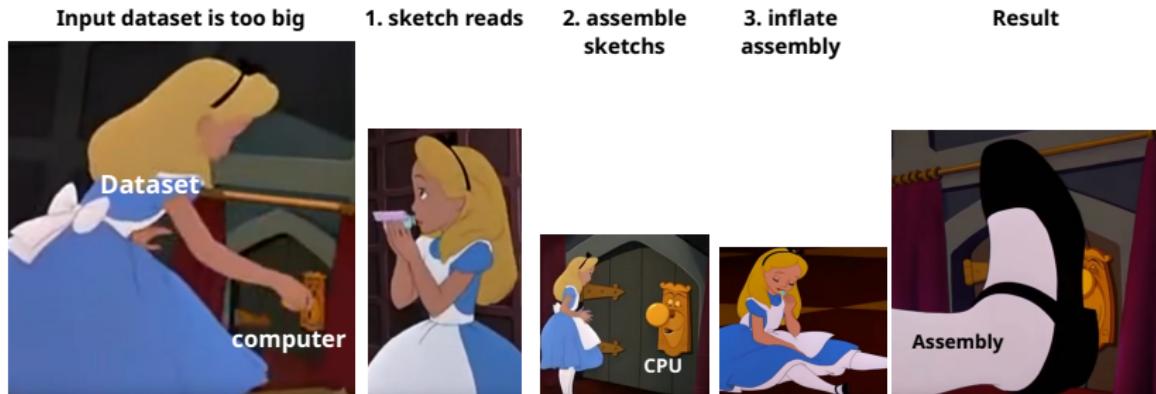
- ▶ Keep a record while compressing

sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

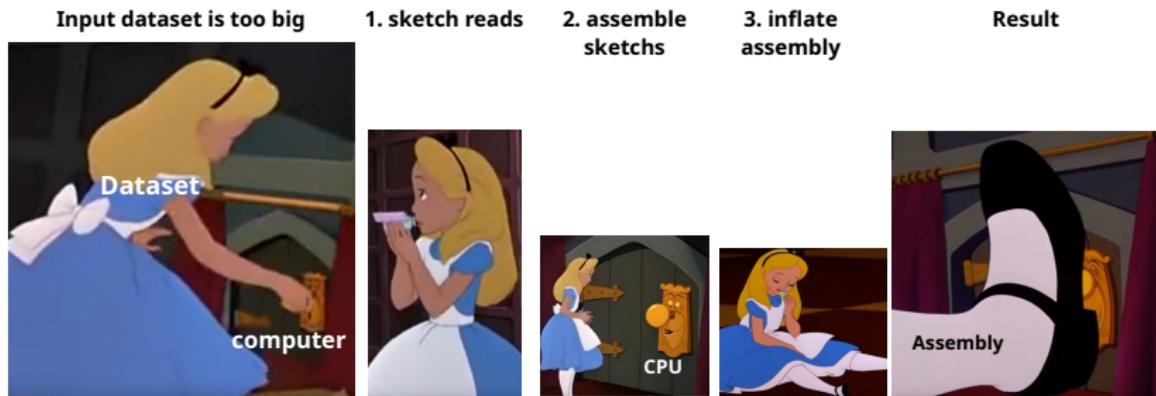
sketch _____ **A****G** _____ **A** _____ **G** _____ **C** _____ **A** _____

record { AGAG → ATGGATAACAGATGGAGATATCATCG,
 GAGC → AGTATGGATAACAGATGGAGATATCATCGAGTAGGGC,
 AGCA → GATGGAGATATCATCGAGTAGGGGCACTGTAC }

The Alice assembler: assembling with MSR



The Alice assembler: assembling with MSR



- ▶ Function f : chaotic hash function
- ▶ Any assembler for step 2., by default BCALM2+tip-clipping
- ▶ github.com/rolandfaure/alice-asm

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

- $f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
- $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
- $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
- $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
- $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1

CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch1

A

sequence2

CAGTATGGATACAGATGGAGATAT**G**ATCGAGTAGGGGCACTGTACCAGAG

sketch2

A

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

- $f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
- $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
- $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
- $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
- $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1	C <u>AGTATGGATA</u> CAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
sketch1	A
sequence2	C <u>AGTATGGATA</u> CAGATGGAGATAT <u>G</u> ATCGAGTAGGGGCACTGTACCAGAG
sketch2	A

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

- $f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
- $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
- $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
- $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
- $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1

CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch1

A G

sequence2

CAGTATGGATACAGATGGAGATATGATCGAGTAGGGGCACTGTACCAGAG

sketch2

A G

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

- $f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
- $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
- $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
- $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
- $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1	CAGTATGGATAACAG	ATGGAGATAT	CATCGAGTAGGGGCACTGTACCAGAG
sketch1	A G T		
sequence2	CAGTATGGATAACAG	ATGGAGATATG	CATCGAGTAGGGGCACTGTACCAGAG
sketch2	A G T		

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

- $f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
- $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
- $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
- $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
- $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1	CAGTATGGATAACAGA	TGGAGATATC	ATCGAGTAGGGGCACTGTACCAGAG
sketch1	A G T		T
sequence2	CAGTATGGATAACAGA	TGGAGATATG	ATCGAGTAGGGGCACTGTACCAGAG
sketch2	A G T		

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
 $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
 $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
 $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
 $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1	CAGTATGGATAACAGAT	GGAGATATCA	TCGAGTAGGGCACTGTACCAGAG
sketch1	A G T		T
sequence2	CAGTATGGATAACAGAT	GGAGATATGA	TCGAGTAGGGCACTGTACCAGAG
sketch2	A G T		G

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

- $f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
- $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
- $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
- $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
- $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1	CAGTATGGATAACAGATG	GAGATATCAT	CGAGTAGGGCACTGTACCAGAG
sketch1	A G T		T C
sequence2	CAGTATGGATAACAGATG	GAGATATGAT	CGAGTAGGGCACTGTACCAGAG
sketch2	A G T		G

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
 $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
 $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
 $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
 $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1	CAGTATGGATAACAGATGG	AGATATC ATC	GAGTAGGGCACTGTACCAGAG
sketch1	A G T	T C	
sequence2	CAGTATGGATAACAGATGG	AGATATG ATC	GAGTAGGGCACTGTACCAGAG
sketch2	A G T	G	

MSRs keep SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

- $f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$
- $f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$
- $f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$
- $f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$
- $f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$

sequence1

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCAC**TGTACCAAGAG**

sketch1

A G T T C C G T C

sequence2

CAGTATGGATAACAGATGGAGATAT**G**ATCGAGTAGGGGCAC**TGTACCAAGAG**

sketch2

A G T G A C G T C

MSRs keep and amplify SNPs

$$f : \{A, C, G, T\}^{10} \rightarrow \{A, C, G, T, \emptyset\}$$

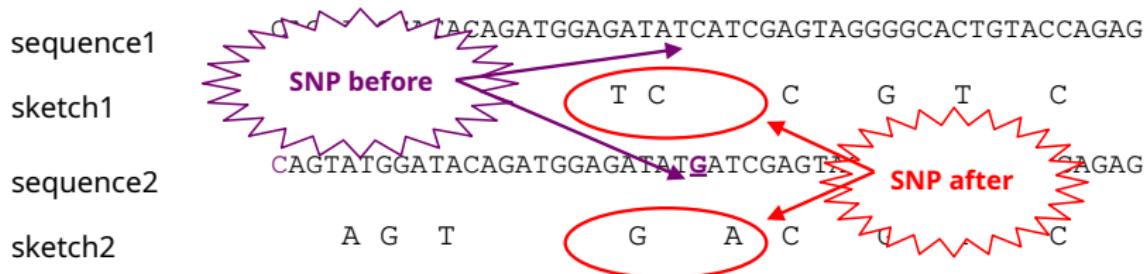
$f(10-mer) \rightarrow A \text{ if } \text{hash}(10-mer) \in [0, 0.05]$

$f(10-mer) \rightarrow C \text{ if } \text{hash}(10-mer) \in [0.05, 0.1]$

$f(10-mer) \rightarrow G \text{ if } \text{hash}(10-mer) \in [0.1, 0.15]$

$f(10-mer) \rightarrow T \text{ if } \text{hash}(10-mer) \in [0.15, 0.2]$

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$



MSRs keep and amplify SNPs

- ▶ A SNP affects l l -mers
- ▶ Each l -mer outputs a base with probability c
- ▶ Probability that a SNP disappears in the sketch:

$$\sum_{i=0}^l \left(\binom{l}{i} c^i \cdot (1-c)^{l-i} \right)^2 \cdot \frac{1}{4^i} \approx (1-c)^{2l}$$

MSRs keep and amplify SNPs

- ▶ A SNP affects l l -mers
- ▶ Each l -mer outputs a base with probability c
- ▶ Probability that a SNP disappears in the sketch:

$$\sum_{i=0}^l \left(\binom{l}{i} c^i \cdot (1-c)^{l-i}\right)^2 \cdot \frac{1}{4^i} \approx (1-c)^{2l}$$

	k-mer subsampling	MSR
$c=0.1$	0.81	10^{-10}
$c=0.01$	0.98	0.13

Table: Probability that a SNP disappears in sketch, using $l=101$

The Alice assembler: results

- ▶ Zymobiomics Gut Microbiome Standard with 5 strains of *E.coli*

The Alice assembler: results

- ▶ Zymobiomics Gut Microbiome Standard with 5 strains of *E.coli*

	Genome fraction (%)	
	metamdbq	alice
Escherichia_coli_B1109	78.408	92.039
Escherichia_coli_B3008	36.411	99.968
Escherichia_coli_B766	95.647	95.641
Escherichia_coli_JM109	38.211	96.334
Escherichia_coli_b2207	37.335	95.495

Measured using metaQUAST

- ▶ Strains are not collapsed

The Alice assembler: results

- ▶ Assembling a human gut metagenome (HiFi sequencing)

Flye
4d, 256G RAM



hifiasm_meta
11d, 454G RAM



metaMDBG
19h, 10G RAM



The Alice assembler: results

- ▶ Assembling a human gut metagenome (HiFi sequencing)

Flye
4d, 256G RAM



hifiasm_meta
11d, 454G RAM



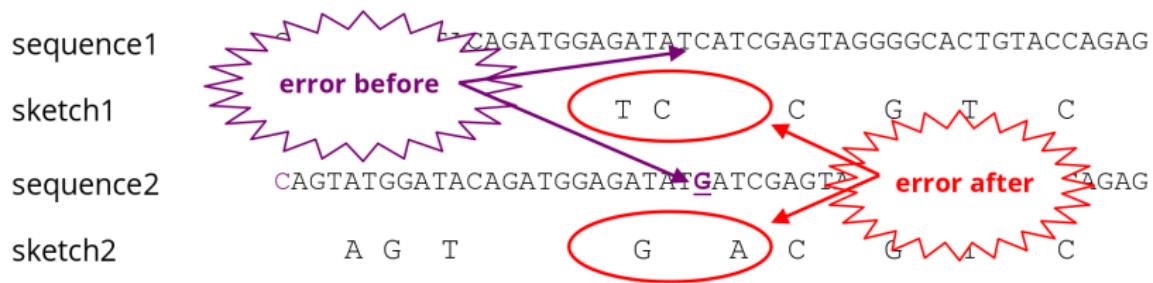
metaMDBG
19h, 10G RAM



Alice
5h, 10G RAM



The dark side of MSR: errors



- ▶ Distance between errors \approx Original distance * compression ratio c
- ▶ original sequences align \iff reduced sequences align **not completely true**

Potential applications and future MSRs



- ▶ MSRs are wild and unexplored

Potential applications and future MSRs



- ▶ MSRs are wild and unexplored
- ▶ Alignment? SNP calling? Indexing? Whole genome operations (e.g. pangenome graph building)?

Potential applications and future MSRs



- ▶ MSRs are wild and unexplored
- ▶ Alignment? SNP calling? Indexing? Whole genome operations (e.g. pangenome graph building)?
- ▶ Changing the MSR itself: error rate? Biology-informed MSR?