

A new sketching method for metagenome assembly

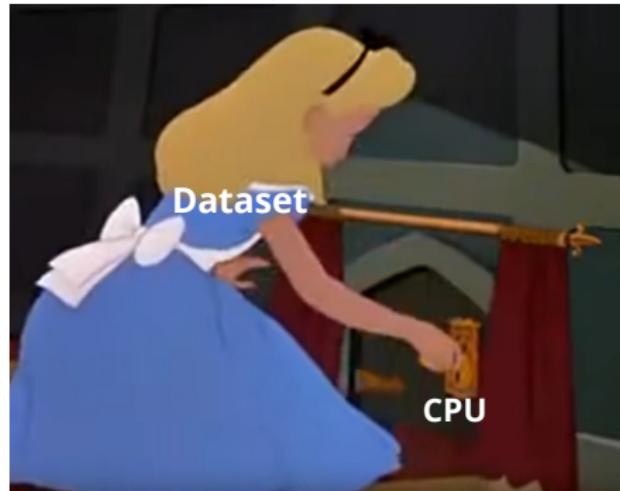
Roland Faure^{1,2}, Baptiste Hilaire², Jean-François Flot¹,
Dominique Lavenier²

¹Université libre de Bruxelles (ULB) - Belgium

²Université de Rennes, IRISA - France

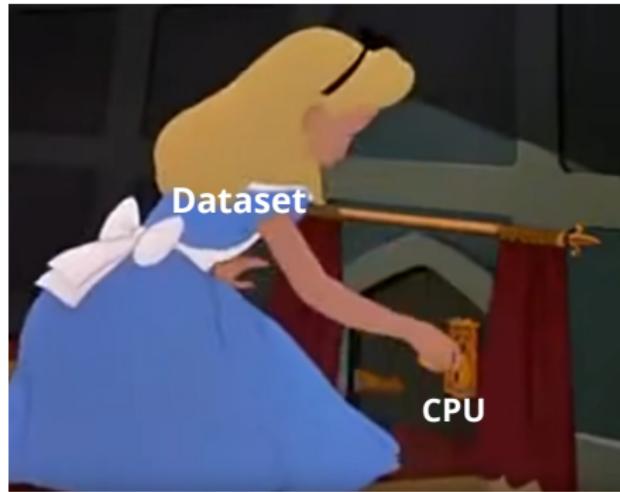
Genome Informatics 2024

Sequencing data is big



Credits: Alice in Wonderland, Lewis, Disney

Sequencing data is big



Credits: Alice in Wonderland, Lewis, Disney



Drink-me
potion

Sketching with sequence subsampling

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



k-mer subsampling

TACG

TGCT **TGCT**



expensive computation

...

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

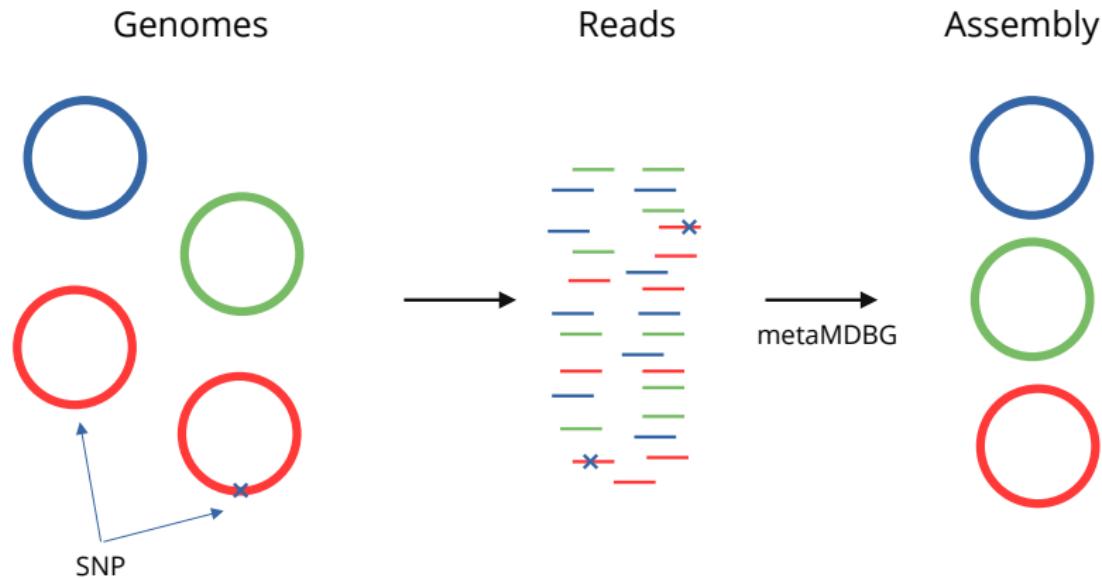
Sequence subsampling does not preserve SNPs

SNP



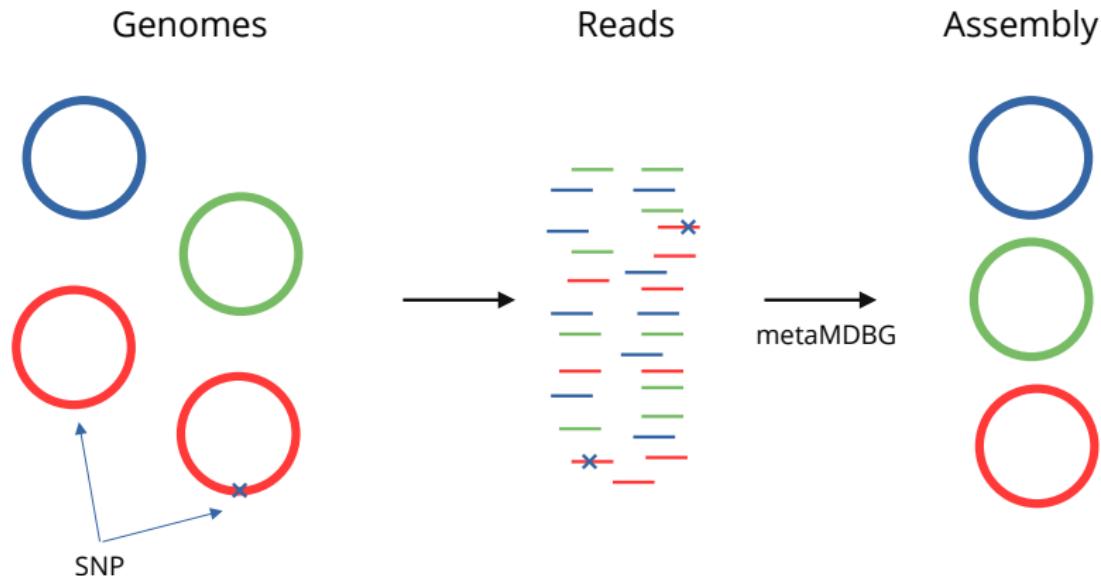
CAGACTA A GATATTTTGCTGACTCAT	→	AGAC	ATTT	CTCA
CAGACTACGAT ATTTTGCTGACTCAT	→	AGAC	ATTT	CTCA

My problem: complete metagenome assembly



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

My problem: complete metagenome assembly



- ▶ metaMDBG is very fast, but some variants are lost!
- ▶ Let me introduce **MSR sketching**

Introducing a MSR sketch

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

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

Introducing a MSR sketch

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

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

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Introducing a MSR sketch

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

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

sequence **CAGTATGGAT**ACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**CAGTATGGAT**) = 0.0023

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

sketch A

Introducing a MSR sketch

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

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

sequence C**AGTATGGATA**CAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**AGTATGGATA**) = 0.624
 $f(\textbf{AGTATGGATA}) = \emptyset$

sketch A

Introducing a MSR sketch

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

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

sequence CA**GTATGGATAC**AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

$$\text{hash}(\textbf{GTATGGATAC}) = 0.124$$
$$f(\textbf{GTATGGATAC}) = G$$

sketch A G

Introducing a MSR sketch

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

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

sequence CAG**TATGGATACA**GATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**TATGGATACA**) = 0.88
 $f(\textbf{TATGGATACA}) = \emptyset$

sketch A G

Introducing a MSR sketch

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

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

sequence CAGT**ATGGATACAG**ATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**ATGGATACAG**) = 0.32
 $f(\textbf{ATGGATACAG}) = \emptyset$

sketch A G

Introducing a MSR sketch

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

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

sequence CAGTA**TGGATACAGA**TGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**TGGATACAGA**) = 0.19
 $f(\textbf{TGGATACAGA}) = T$

sketch A G T

Introducing a MSR sketch

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

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

sequence CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**GGATACAGAT**) = 0.214
 $f(\textbf{GGATACAGAT}) = \emptyset$

sketch A G T

Introducing a MSR sketch

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

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

sequence CAGTATG**GATACAGATG**GAGATATCATCGAGTAGGGGCACTGTACCAGAG

$$\text{hash}(\textcolor{red}{GATACAGATG}) = 0.678$$
$$f(\textcolor{red}{GATACAGATG}) = \emptyset$$

sketch A G T

Introducing a MSR sketch

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

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

sequence

CAGTATGG**ATACAGATGG**AGATATCATCGAGTAGGGGCACTGTACCAGAG

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

sketch

A G T

Introducing a MSR sketch

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

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

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCAC **TGTACCAGAG**

$$\text{hash}(\textcolor{red}{TGTACCAGAG}) = 0.06$$

$$f(\textcolor{red}{TGTACCAGAG}) = C$$

sketch

A G T T C C G T C

Introducing a MSR sketch

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

order (l) $\xrightarrow{\hspace{1cm}}$

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

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

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

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

$f(10-mer) \rightarrow \emptyset \text{ if } \text{hash}(10-mer) > 0.2$ $\xrightarrow{\hspace{1cm}}$ compression ratio (c)

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch

A G T

T C

C G T C

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	<u>CAGTATGGAT</u> ACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
sketch1	A
sequence2	<u>CAGTATGGAT</u> ACAGATGGAGATAT <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	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	CA <u>GTATGGATA</u> CAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
sketch1	A G
sequence2	CA <u>GTATGGATA</u> CAGATGGAGATAT <u>G</u> ATCGAGTAGGGGCACTGTACCAGAG
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	CAGTATGGATACAGA	TGGAGATATC	ATCGAGTAGGGGCACTGTACCAGAG
sketch1	A G T		T
sequence2	CAGTATGGATACAGA	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	<u>GGAGATATCA</u>	TCGAGTAGGGGCACTGTACCAGAG
sketch1	A G T		T
sequence2	CAGTATGGATAACAGAT	<u>GGAGATATGA</u>	TCGAGTAGGGGCACTGTACCAGAG
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	CGAGTAGGGGCACTGTACCAGAG
sketch1	A G T	T C	
sequence2	CAGTATGGATAACAGATG	GAGATATGAT	CGAGTAGGGGCACTGTACCAGAG
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	AGATATCATC	GAGTAGGGGCACTGTACCAGAG
sketch1	A G T		T C
sequence2	CAGTATGGATAACAGATGG	AGATATGATC	GAGTAGGGGCACTGTACCAGAG
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 **TGTACCAGAG**

sketch1 A G T T C C G T C

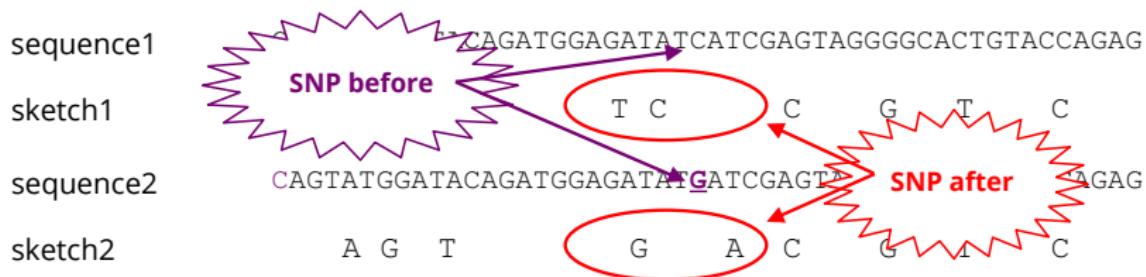
sequence2 CAGTATGGATAACAGATGGAGATAT**G**ATCGAGTAGGGGCAC **TGTACCAGAG**

sketch2 A G T G A C G T C

MSRs keep and amplify SNPs

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

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



- If $l \cdot c > 1$, SNPs are **amplified**

MSR = Mapping-friendly Sequence Reductions¹

- ▶ Sketches are reduced sequences
- ▶ Mapping-friendly

ATCATCGAGTAGGGGCACTGTACCAGAGCGCTTTAATGTAC
CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

A G T T C C G T C

¹Bassel, Luc & Medvedev, Paul & Chikhi, Rayan. (2022). Mapping-friendly sequence reductions: Going beyond homopolymer compression. iScience.

MSR = Mapping-friendly Sequence Reductions¹

- ▶ Sketches are reduced sequences
- ▶ Mapping-friendly

ATCATCGAGTAGGGGCACTGTACCAGAGCGCTTTAATGTAC
CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
A G T T C C G T C

- ▶ Reverse-complement property not shown

¹Bassel, Luc & Medvedev, Paul & Chikhi, Rayan. (2022). Mapping-friendly sequence reductions: Going beyond homopolymer compression. iScience.

An MSR assembler

MSR sketching

AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
GAGATATCATCGAGTAGGGGCACTGTACCAGAGCCGG
GATATCATCGAGTAGGGGCACTGTACCAGAGGCCGGTTATAC

↓

AGTTCCGT TCCGTCAA CGTCAATG

An MSR assembler



An MSR assembler

MSR sketching

AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
GAGATATCATCGAGTAGGGGCACTGTACCAGAGCCGG
GATATCATCGAGTAGGGGCACTGTACCAGAGGCCGGTTATAC



AGTTCCGT

TCCGTCAA

CGTCAATG

Assembly



AGTTCCGT
TCCGTCAA
CGTCAATG
AGTTCCGTCAAATG



Inflating

AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAGGCCGGTTATAC

Going back to uncompressed space

- ▶ MSR is lossy → not strictly reversible

Going back to uncompressed space

- ▶ MSR is lossy → not strictly reversible
- ▶ Keep a record while compressing

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch

AG _____ **A** _____ **G** _____ **C** _____ **A** _____

record

{ AGAG → ATGGATAACAGATGGAGATATCATCG,
GAGC → AGTATGGATAACAGATGGAGATATCATCGAGTAGGGG,
AGCA → GATGGAGATATCATCGAGTAGGGGCACTGTAC }

Going back to uncompressed space

- ▶ MSR is lossy → not strictly reversible
- ▶ Keep a record while compressing

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch

AG _____ **A** _____ **G** _____ **C** _____ **A** _____

record

{ AGAG → ATGGATAACAGATGGAGATATCATCG,
GAGC → AGTATGGATAACAGATGGAGATATCATCGAGTAGGGG,
AGCA → GATGGAGATATCATCGAGTAGGGGCACTGTAC }

- ▶ Actually record 31-mers

The Alice assembler

1.
sketch the reads



2.
assemble the
sketches



3.
inflate the
assembly



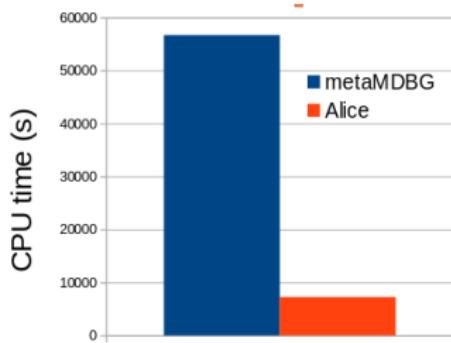
Credits: Alice in Wonderland, Lewis, Disney

- ▶ Any assembler, by default BCALM2+tip-clipping
- ▶ github.com/rolandfaure/alice-asm
(warning: immature code)



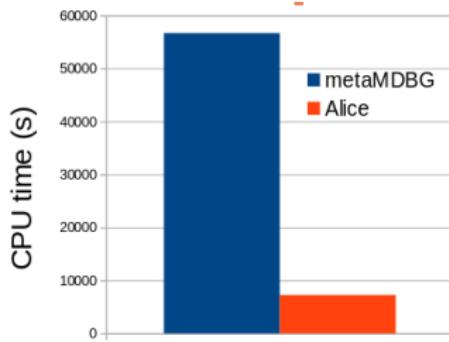
The Alice assembler: results

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



The Alice assembler: results

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

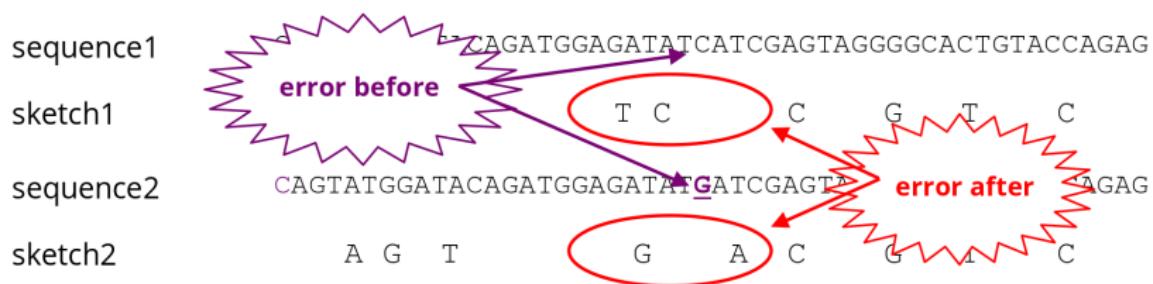


	Genome fraction (%)	
	metamdbg	alice
Escherichia_coli_B1109	78.408	92.039
Escherichia_coli_B3008	36.411	99.968
Escherichia_coli_B766	95.647	95.641
Escherichia_coli_JM109	38.211	96.334
Escherichia_coli_b2207	37.335	95.495

Measured using metaQUAST

- ▶ Strains are not collapsed

The dark side of MSR: errors



- ▶ Errors are amplified: Alice only works on highly accurate reads

MSR sketching: take-home messages

1.
sketch the reads



2.
assemble the sketches



3.
inflate the assembly



Credits: Alice in Wonderland, Lewis, Disney

- ▶ MSR sketches are sequences

MSR sketching: take-home messages

1.
sketch the reads



2.
assemble the sketches



3.
inflate the assembly



Credits: Alice in Wonderland, Lewis, Disney

- ▶ MSR sketches are sequences
- ▶ Differences between sequences are **kept & amplified**

MSR sketching: take-home messages

1.
sketch the reads



2.
assemble the sketches



3.
inflate the assembly



Credits: Alice in Wonderland, Lewis, Disney

- ▶ MSR sketches are sequences
- ▶ Differences between sequences are **kept & amplified**
- ▶ Alice HiFi assembler out soon
(github.com/rolandfaure/alice-asm)