

Mapping-friendly Sequence Reductions to process compressed genomic data

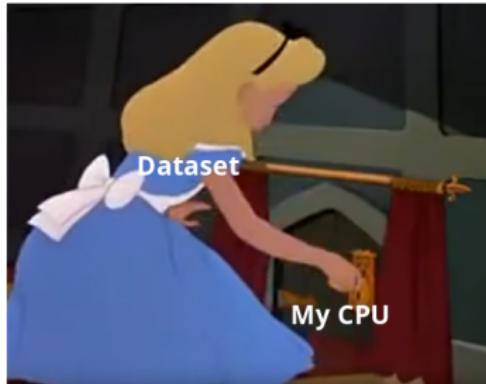
Roland Faure^{1,2}, Baptiste Hilaire², Dominique Lavenier²

¹Université libre de Bruxelles (ULB) - Belgium

²Université de Rennes, IRISA - France

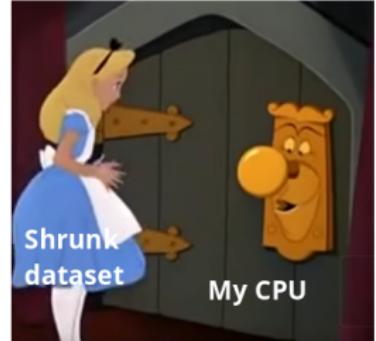
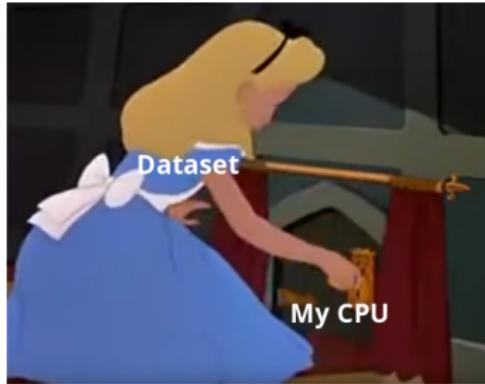
SeqBIM Lille 2023

Big data



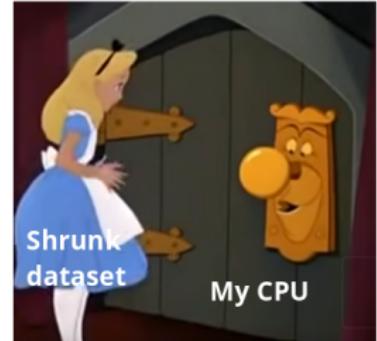
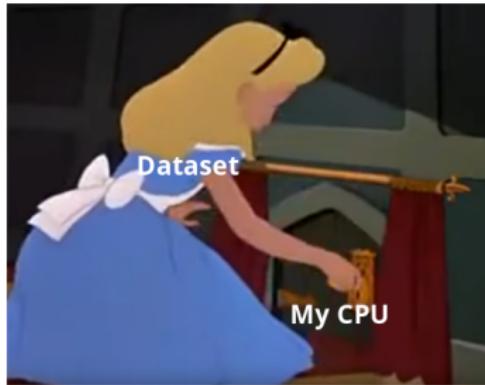
Credits: Alice in Wonderland, Lewis, Disney

Big data



Credits: Alice in Wonderland, Lewis, Disney

Big data



Credits: Alice in Wonderland, Lewis, Disney

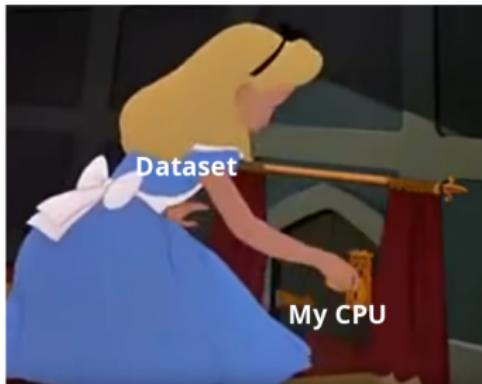
Shrinking sequence information to process shrunk datasets

Shrinking sequences

- ▶ gzip, xz...

Shrinking sequences

- ▶ gzip, xz...

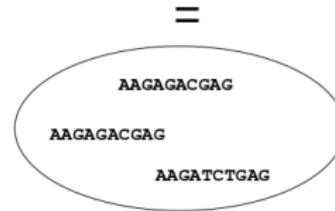


Credits: Alice in Wonderland, Lewis, Disney

Shrinking sequence information **to process shrunk datasets**

State-of-the-art: shrinking sequences with k-mers

CCAAGAGACGAGCAATACGAGCTTTTCAGAGCAGATAATAAATAGCGCTAGCTTACGACAGAGAGGAGACTTGAGAAGATCTGAGATCGGCATAAGCA



- ▶ subsets of k-mers are compressed representation of sequences

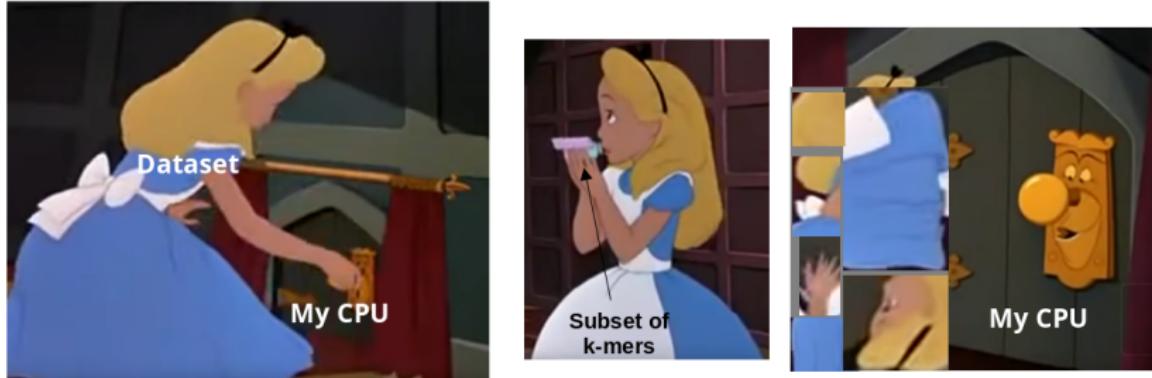
State-of-the-art: shrinking sequences with k-mers

CCAAGAGACGAGCAATACGAGCTTTTCAGAGCAGATAATAAATAGCGCTAGCTTACGACAGAGAGGAGACTTGAGAAGATCTGAGATCGGCATAAGCA



- ▶ subsets of k-mers are compressed representation of sequences
- ▶ minimap2 (alignment), mDBG (assembly), mash (sequence similarity), PebbleScout (sequence querying)...

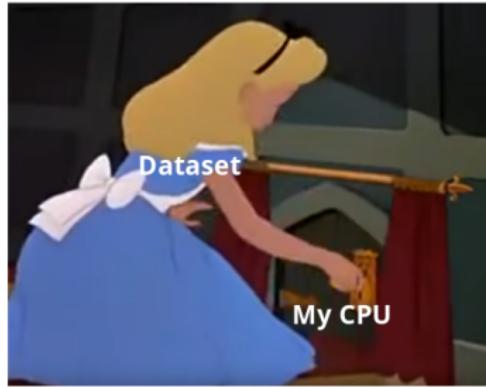
Going beyond shrinking sequences with k-mers



Credits: Alice in Wonderland, Lewis, Disney

- ▶ No structural information
- ▶ Big chunks of data missing
- ▶ Awkward tasks: SVs, gene detection, assembly...

MSR compression to shrink sequences



Credits: Alice in Wonderland, Lewis, Disney

- ▶ Different properties compared to subset of k-mers
- ▶ Different problems

Mapping-friendly Sequence Reduction (MSR)

- ▶ Introduced in iScience 2022¹ to improve mapping accuracy

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



Mapping-friendly Sequence Reduction (MSR)

- ▶ Introduced in iScience 2022¹ to improve mapping accuracy

ACAA**A**GACG**GGGATTT**CGCG**GAT** → ACAGACGATCGCGAT

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



Mapping-friendly Sequence Reduction (MSR)

- ▶ Introduced in iScience 2022¹ to improve mapping accuracy

ACA**A**GACG**GGGATTT**CGCG**GAT**
ACG**GGGATTT**CGCG**GATCGGTCA**

→ ACAGACGATCGCGAT
ACGATCGCGATCGTCA

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



Mapping-friendly Sequence Reduction (MSR)

- ▶ Introduced in iScience 2022¹ to improve mapping accuracy

ACA**A**GACG**GGGATTT**CGCG**GAT**CG**GTCA** → ACAGACGATCGCGAT
ACG**GGGATTT**CGCG**GAT**CG**GTCA** → ACGATCGCGATCGTCA

ACAAGACGGGG**ATTT**CGCG**GAT**
ACGGGG**ATTT**CGCG**GAT**CG**GTCA** → AGGGTTG
GGGTTGG

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



Mapping-friendly Sequence Reduction (MSR)

- ▶ Introduced in iScience 2022¹ to improve mapping accuracy

ACA**A**GACG**GGGATTT**CGCG**GATCGGTCA** → ACAGACGATCGCGAT
ACG**GGGATTT**CGCG**GATCGGTCA** → ACGATCGCGATCGTCA

ACAAGACGGGGATTTCGCG**GAT**
ACGGGGATTTCGCG**GATCGGTCA** → AGGGTTG
GGGTTGG

- ▶ Mapping-friendly
- ▶ Computed in a streaming fashion
- ▶ Reverse-complement insensitive

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



Computing a MSR

$$f : \{A, C, G, T\}^3 \longrightarrow \{A, C, G, T, _\}$$

AGT → G	CAG → A
ACT → C	CCG → C
ATT → T	CGG → G
AAT → A	CTG → T
else → _	

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence

A

Computing a MSR

$$f : \{A, C, G, T\}^3 \longrightarrow \{A, C, G, T, _\}$$

AGT → G	CAG → A
ACT → C	CCG → C
ATT → T	CGG → G
AAT → A	CTG → T
else → _	

sequence

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence

AG

Computing a MSR

$$f : \{A, C, G, T\}^3 \longrightarrow \{A, C, G, T, _\}$$

AGT → G	CAG → A
ACT → C	CCG → C
ATT → T	CGG → G
AAT → A	CTG → T
else → _	

sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence **AG**_

Computing a MSR

$$f : \{A, C, G, T\}^3 \longrightarrow \{A, C, G, T, _\}$$

AGT → G	CAG → A
ACT → C	CCG → C
ATT → T	CGG → G
AAT → A	CTG → T
else → _	

sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence **AG**_

Computing a MSR

$$f : \{A, C, G, T\}^3 \longrightarrow \{A, C, G, T, _\}$$

AGT → G	CAG → A
ACT → C	CCG → C
ATT → T	CGG → G
AAT → A	CTG → T
else → _	

sequence CAGT**ATG**GATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence **AG**_____

Computing a MSR

$$f : \{A, C, G, T\}^3 \longrightarrow \{A, C, G, T, _\}$$

AGT → G	CAG → A
ACT → C	CCG → C
ATT → T	CGG → G
AAT → A	CTG → T
else → _	

sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCA**GAG**

Reduced sequence **AG** _____ **A** _____ **G** _____ **C** _____ **A** _____

Computing a MSR

$$f : \{A, C, G, T\}^3 \longrightarrow \{A, C, G, T, _\}$$

AGT → G	CAG → A
ACT → C	CCG → C
ATT → T	CGG → G
AAT → A	CTG → T
else → _	

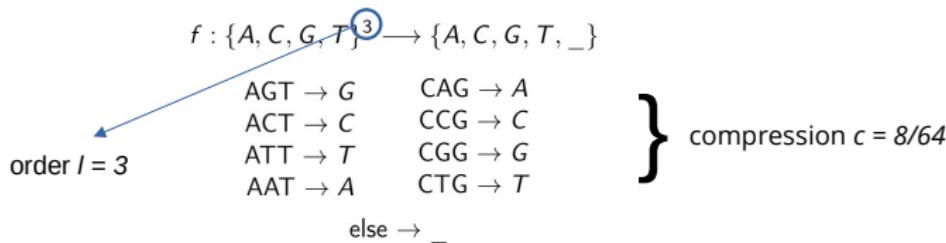
}

compression $c = 8/64$

sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCA**GAG**

Reduced sequence **AG** _____ **A** _____ **G** _____ **C** _____ **A** _____

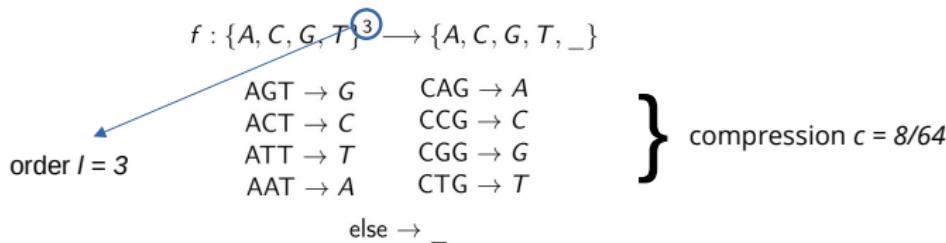
Computing a MSR



sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCA**GAG**

Reduced sequence **AG**_____ **A**_____ **G**_____ **C**_____ **A**_____

Computing a MSR



sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCA **GAG**

Reduced sequence **AG** _____ **A** _____ **G** _____ **C** _____ **A** _____

- ▶ Let's use this to shrink the sequences of our datasets

Demo: an experimental HiFi assembler

Demo: an experimental HiFi assembler

- ▶ Step 1: Compress the reads $c=0.05$ $l=201$

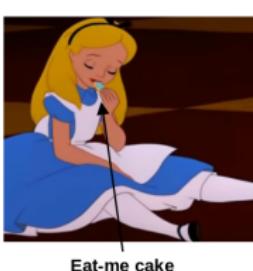
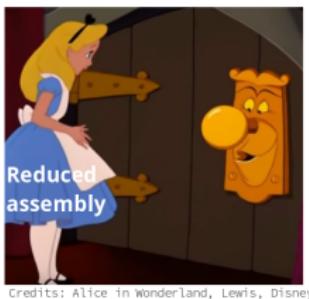
Demo: an experimental HiFi assembler

- ▶ Step 1: Compress the reads $c=0.05$ $l=201$
- ▶ Step 2: Assemble the compressed reads BCALM2² $k=41$

²Rayan Chikhi, Antoine Limasset and Paul Medvedev, Compacting de Bruijn graphs from sequencing data quickly and in low memory, Proceedings of ISMB 2016, Bioinformatics, 32 (12): i201-i208.

Demo: an experimental HiFi assembler

- ▶ Step 1: Compress the reads $c=0.05$ $l=201$
- ▶ Step 2: Assemble the compressed reads BCALM2² $k=41$
- ▶ Step 3: Inflate the contigs



Credits: Alice in Wonderland, Lewis, Disney

²Rayan Chikhi, Antoine Limasset and Paul Medvedev, Compacting de Bruijn graphs from sequencing data quickly and in low memory, Proceedings of ISMB 2016, Bioinformatics, 32 (12): i201-i208.

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

Reduced sequence **A**G_____ **A**_____ **G**_____ **C**_____ **A**_____

Record { AGAG → CAGTATGGATAACAGATGGAGATATCATCGAGT,
 GAGC → AGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACT,
 AGCA → CAGATGGAGATATCATCGAGTAGGGGCACTGTACCAG }

Why use MSR ?

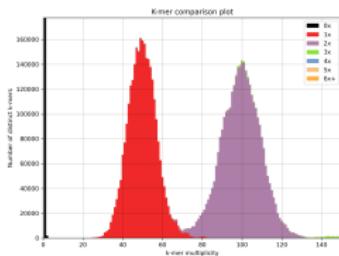
- ▶ Dataset: mix of *E. coli* strains 12009 and HS, simulated HiFi
- ▶ hifiasm: perfect assembly, 22351s
- ▶ mDBG: not-so-good assembly, 30s

Why use MSR ?

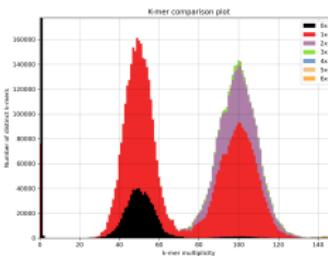
- ▶ Dataset: mix of *E. coli* strains 12009 and HS, simulated HiFi
- ▶ hifiasm: perfect assembly, 22351s
- ▶ mDBG: not-so-good assembly, 30s
- ▶ MSR assembler: 90s

Why use MSR ?

- ▶ Dataset: mix of *E. coli* strains 12009 and HS, simulated HiFi
- ▶ hifiasm: perfect assembly, 22351s
- ▶ mDBG: not-so-good assembly, 30s
- ▶ MSR assembler: 90s



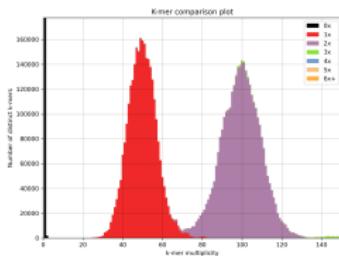
KAT of hifiasm



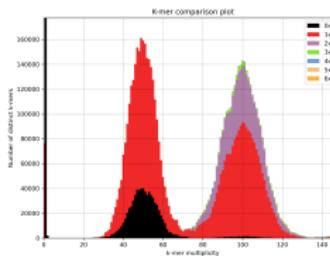
KAT of mDBG

Why use MSR ?

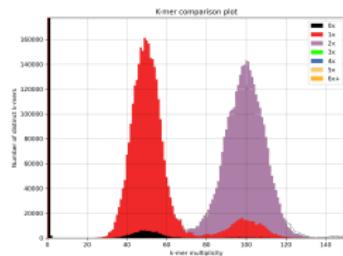
- ▶ Dataset: mix of *E. coli* strains 12009 and HS, simulated HiFi
- ▶ hifiasm: perfect assembly, 22351s
- ▶ mDBG: not-so-good assembly, 30s
- ▶ MSR assembler: 90s



KAT of hifiasm



KAT of mDBG



KAT of MSR

Difference between minimizers and MSRs

CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

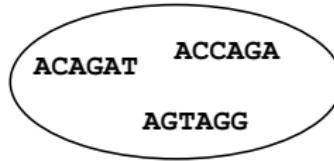
CAGTATGGATAACAGATGGAGATATGATCGAGTAGGGGCACTGTACCAGAG

Difference between minimizers and MSRs

CAGTATGGAT**ACAGAT**GGAGATATCATCG**AGTAGG**GGCACTGT**ACCAGA**G



CAGTATGGAT**ACAGAT**GGAGATAT**G**ATCG**AGTAGG**GGCACTGT**ACCAGA**G



Difference between minimizers and MSRs

- ▶ Example $c = 0.2, l = 15$

sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence

sequence CAGTATGGATAACAGATGGAGATAT**G**ATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence

Difference between minimizers and MSRs

- ▶ Example $c = 0.2, l = 15$

sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence —

sequence CAGTATGGATAACAGATGGAGATATGATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence —

Difference between minimizers and MSRs

- ▶ Example $c = 0.2, l = 15$

sequence CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A

sequence CAGTATGGATACAGATGGAGATATGATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A

Difference between minimizers and MSRs

- ▶ Example $c = 0.2, l = 15$

sequence CAGTATGGATAACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence A

sequence CAGTATGGATAACAGATGGAGATATGATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence A

Difference between minimizers and MSRs

- ▶ Example $c = 0.2, l = 15$

sequence CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_

sequence CAGTATGGATACAGATGGAGATATGATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_

Difference between minimizers and MSRs

- ▶ Example $c = 0.2, l = 15$

sequence CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T__

sequence CAGTATGGATACAGATGGAGATATGATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T__

Difference between minimizers and MSRs

- ▶ Example $c = 0.2, l = 15$

sequence CAGTATGGATAAC**AGATGGAGATATCA**TCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_____G

sequence CAGTATGGATAAC**AGATGGAGATATGA**TCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_____

Difference between minimizers and MSRs

- ▶ Example $c = 0.2, l = 15$

sequence CAGTATGGATAAC**AGATGGAGATATCAT**CGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_____G_

sequence CAGTATGGATAACA**GATGGAGATATGAT**CGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_____T

Difference between minimizers and MSRs

- ▶ Example $c = 0.2$, $l = 15$

sequence CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_____G_____A_C_A_T_G_____AC__T

sequence CAGTATGGATACAGATGGAGATTATGATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_____TA_____G_T_G_____AC__T

Difference between minimizers and MSRs

- ▶ Example $c = 0.2$, $l = 15$

sequence CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_____G_____A_C_A_T_G_____AC__T

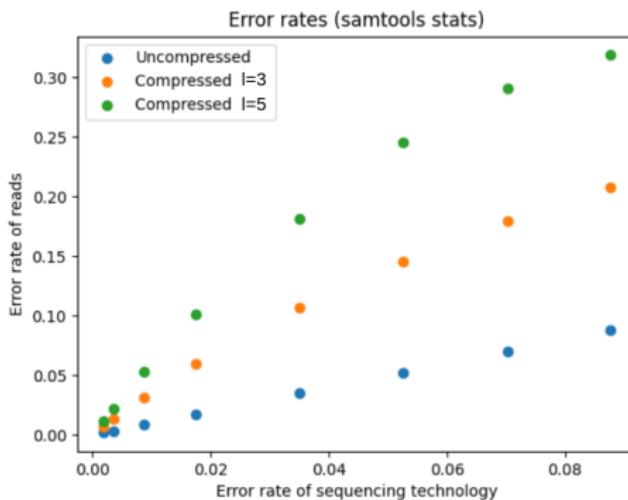
sequence CAGTATGGATACAGATGGAGATTATCGAGTAGGGGCACTGTACCAGAG

Reduced sequence _A_____G_T_____TA_____G_T_G_____AC__T

- ▶ With high l , all bases are used to generate the compression

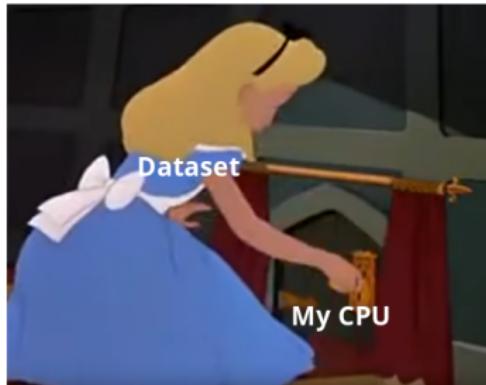
Error rates

- ▶ Sequencing errors get amplified too



- ▶ Rule of thumb for low l : $\text{new_error} = l * \text{old_error}$

Summary: MSR as a compression technique



Credits: Alice in Wonderland, Lewis, Disney

- ▶ Retains contiguity
- ▶ No big chunks of data missing
- ▶ MSRs preserve and amplify differences
- ▶ Compatible with existing software

Perspectives

- ▶ SNP/SV calling
- ▶ pangenome graph building
- ▶ data indexing
- ▶ taxonomic assignment

Perspectives

- ▶ SNP/SV calling
- ▶ pangenome graph building
- ▶ data indexing
- ▶ taxonomic assignment

- ▶ Inflating the data: open question
- ▶ Using different MSRs
- ▶ Apply to more noisy data
- ▶ ...

Acknowledgments



Parameters

- ▶ How I chose my parameters ($c=0.05$ $l=201$, $B_{calm}=31$)

$$\frac{31}{0.05} + 201 = 821$$