

# Haplotype assembly from long reads

Roland Faure<sup>1,2</sup>

<sup>1</sup>Université Libre de Bruxelles (ULB) - Belgium

<sup>2</sup>Université de Rennes, IRISA - France

Genome Assembly Course - December 18th, 2024

“What species do you work on?”

- ▶ Kahoot n°1

# Using genome sequencing to study stuff

Organism, microbiota...

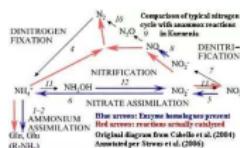


DNA extraction  
& preparation

sample



Sequencing



Understand things

Genomes or  
amplicons

...AACCTGCGTCACGTAGTCGAGG...  
...CGGGCCTGAGGCAGCAGTGCCA...

Assembly

Long reads

CGTAGCTAGGAT  
GTGCTAATCACGT  
TCCGAGCGATCAG  
AAAGCTAATCACTT  
CTCTGGGGTGACA

# Using genome sequencing to study stuff

Organism, microbiota...

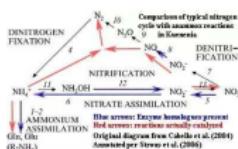


DNA extraction  
& preparation

sample



Sequencing



Understand things

...AACCTGCGTCACGTAGTCGAGG...  
...CGGGCCTGAGGCAGCAGTGCCA...

Genomes or  
amplicons

My Ph.D.

Assembly

CGTAGCTAGGAT  
GGGCTAATCACGT  
AAAGCTAATCACTT  
TCCGAGCGATCAG  
TGTCTGAAACCACAA  
CTCTGGGGTGACA

Long reads

# DNA sequencing

Extracted  
DNA



sequencer

CAGCATCAGTTTCGAGCACGT

TTACTCAGCAGATCGTCGATCAT

CCCGTAGCTTAGCAGGCATCAG

**Reads**

# DNA sequencing: difficulties



length: 1-20 kbp

CAGCATCAGTTTCGAGCACGT

TTACTCAGCAGATCGTCGATCAT

CCCGTAGCTTAGCAGGCATCAG

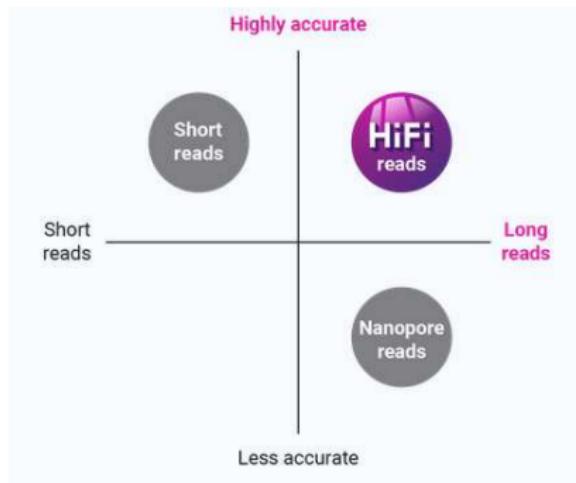
# DNA sequencing: difficulties



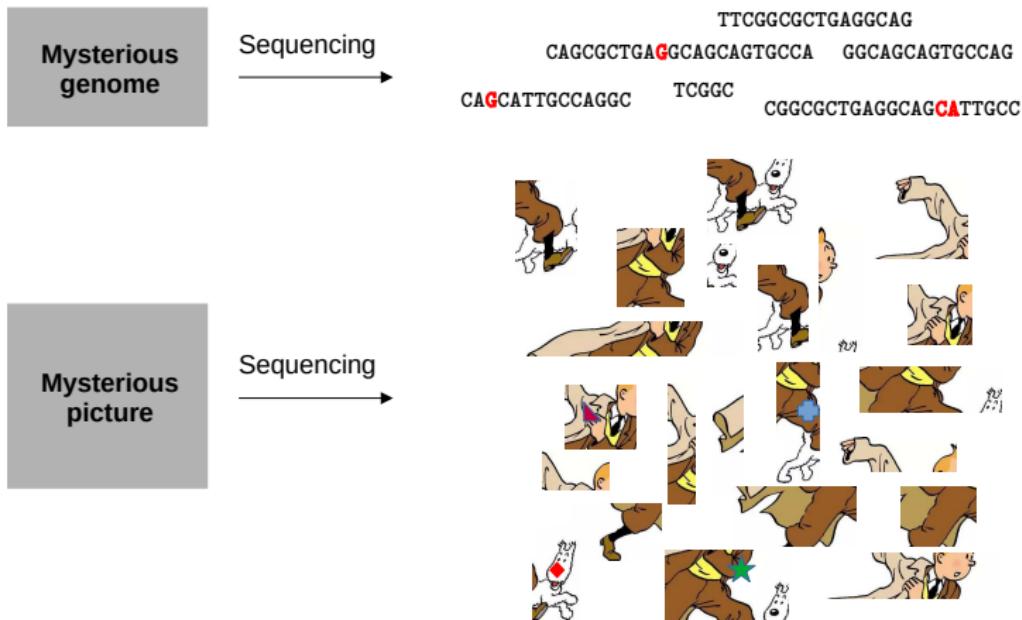
CAGCA**A**TCAAGTTTCGAGCACGT  
TTACTCAGCAGATCG**T**CGATCAT  
CCCGTAGC**TT**AGCAGGCATCAG

sequencing errors: 0.1 – 10 %

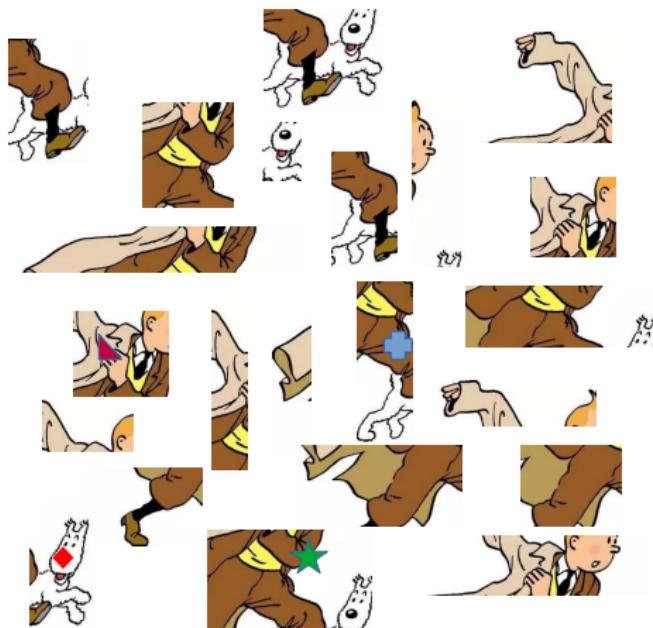
# The different sequencing technologies



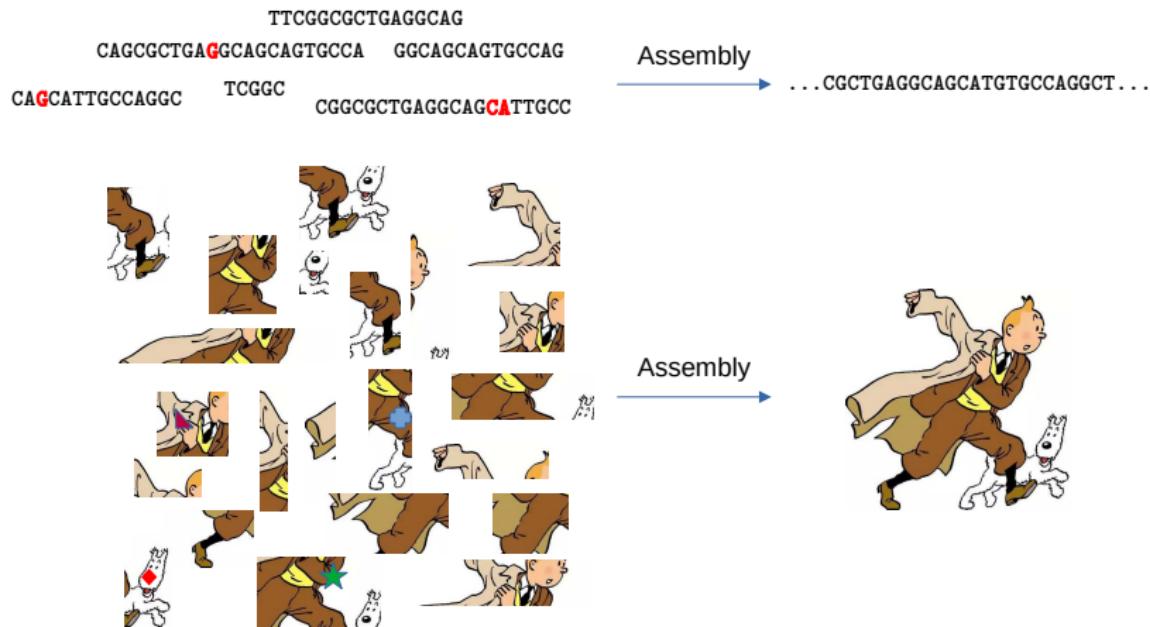
# Genome sequencing



Imagine you are an assembler: what is this picture?  
(Kahoot)



# Genome assembly



## Genome assembly

CGATGCTGGCTAGCATAGTCGATTATCT  
CTGGCTAGC**T**TAGTCGATTATCTGACAGT  
AGCATAGTCGATTATCTGACAGTCATAT  
AGTCGATTAT**A**TGACAGTCATATTGCT  
TTTATCTGACAGTC**A**GATTGCTACACAC

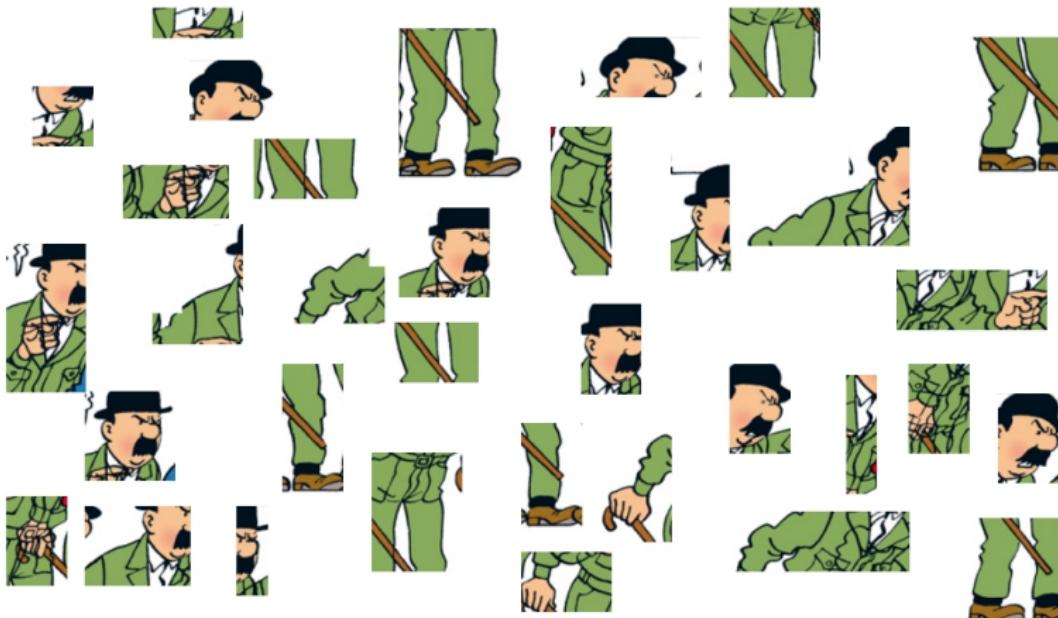
genome assembly: stitching reads  
correcting errors



**CGATGCTGGCTAGCATAGTCGATTATCTGACAGTCATATTGCTACACAC**

- ▶ Many software: Flye, wtdbg2, metaMDBG, hifiasm...

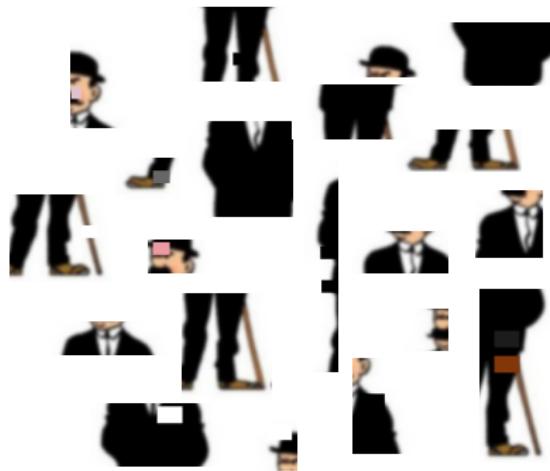
# Imagine you are an assembler (Kahoot)



## Imagine you are an assembler



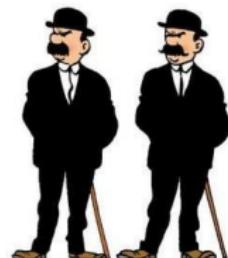
## Haplotype assembly



Collapsed assembly



Haplotype assembly



# Dupont & Dupond exist in genomes!



Maternal DNA

...ACACACCACACACCTCTACGA...

...ACACACTACACACACCTCTACGA...

Paternal DNA



# Dupont & Dupond exist in genomes!



...ACACACCACACACCT**G**TACGA...

...ACACACCACACACCTCTACGA...

...ACACACCACACACCTCTACGA...

...ACACACT**A**CACACACCTCTACGA...

...ACACACCACACACCTCTACGA...

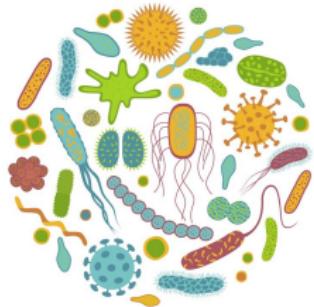
...ACACACCACACACCTCTACGA...

...ACA**A**CCACACACACCTCTACGA...

...ACACACCACACACCTCT**A**GA...



# Dupont & Dupond exist in genomes!



...ACACACCACACACACCTCTACGA...

...ACACACT**T**ACACACACCTCTACGA...

...ACACACCACACACACACCTCTACGA...

...ACA**A**ACCACACACACACCTCTACGA...

...ACACACCACACACACCTCTA**A**GA...

...

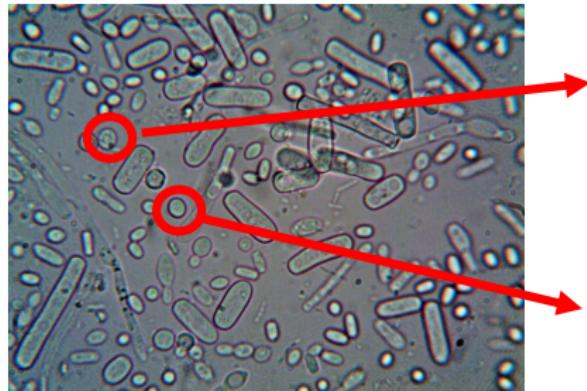


# Dupont & Dupond are important in genomes!

- ▶ Kahoot

# Dupont & Dupond are important in genomes!

## ► Kahoot



*Escherichia coli* Sakai



...ACACACCACACACCTCTACGA...

...ACACAC**T**ACACACACCTCTACGA...

*Escherichia coli* Nissle



## Problem: assembling several haplotypes

CGATGCTGGCTAGCATAGTCGATTATCT  
CTGGCTAGCTAGTCGATTATCTGACAGT  
AGCATAGTCGATTATCTGACAGTCATAT  
AGTCGATTATA**A**TGACAGTCATATTGCT  
TTTATA**A**TGACAGTCAGATTGCTACACAC

genome assembly: stitching reads  
correcting errors



**CGATGCTGGCTAGCATAGTCGATTATCTGACAGTCATATTGCTACACAC**  
**CGATGCTGGCTAGCATAGTCGATTATA**A**TGACAGTCATATTGCTACACAC**

## Problem: assembling several haplotypes

CGATGCTGGCTAGCATAGTCGATTATCT  
CTGGCTAGCTAGTCGATTATCTGACAGT  
AGCATAGTCGATTATCTGACAGTCATAT  
AGTCGATTATA**A**TGACAGTCATATTGCT  
TTTATA**A**TGACAGTCAGATTGCTACACAC

genome assembly: stitching reads  
correcting errors



**CGATGCTGGCTAGCATAGTCGATTATCTGACAGTCATATTGCTACACAC**  
**CGATGCTGGCTAGCATAGTCGATTATA**A**TGACAGTCATATTGCTACACAC**

- ▶ Not so many software!



Haplotype assembly from long reads

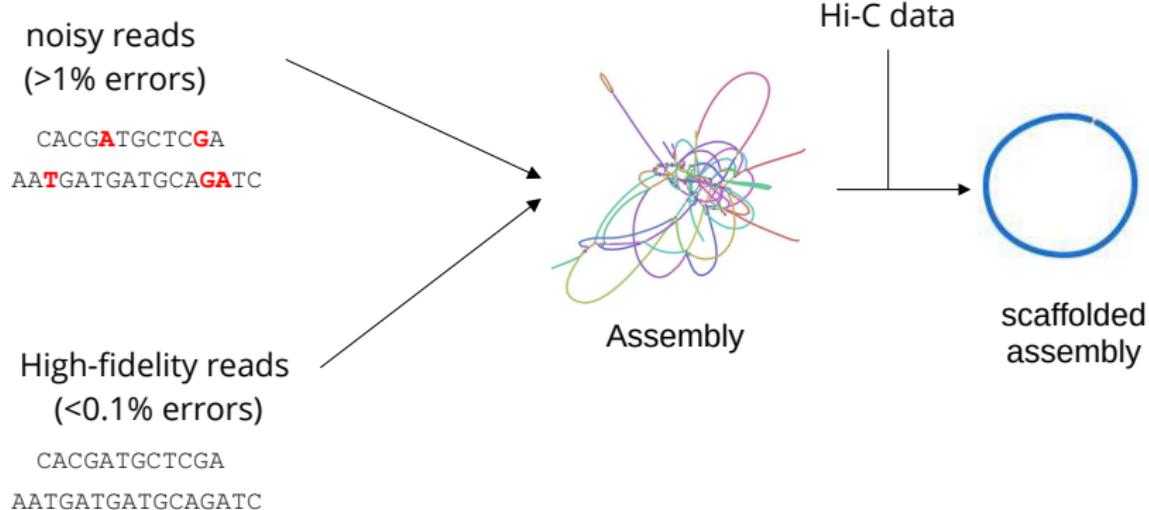
Roland Faure<sup>1,2</sup>

<sup>1</sup>Université Libre de Bruxelles (ULB) - Belgium

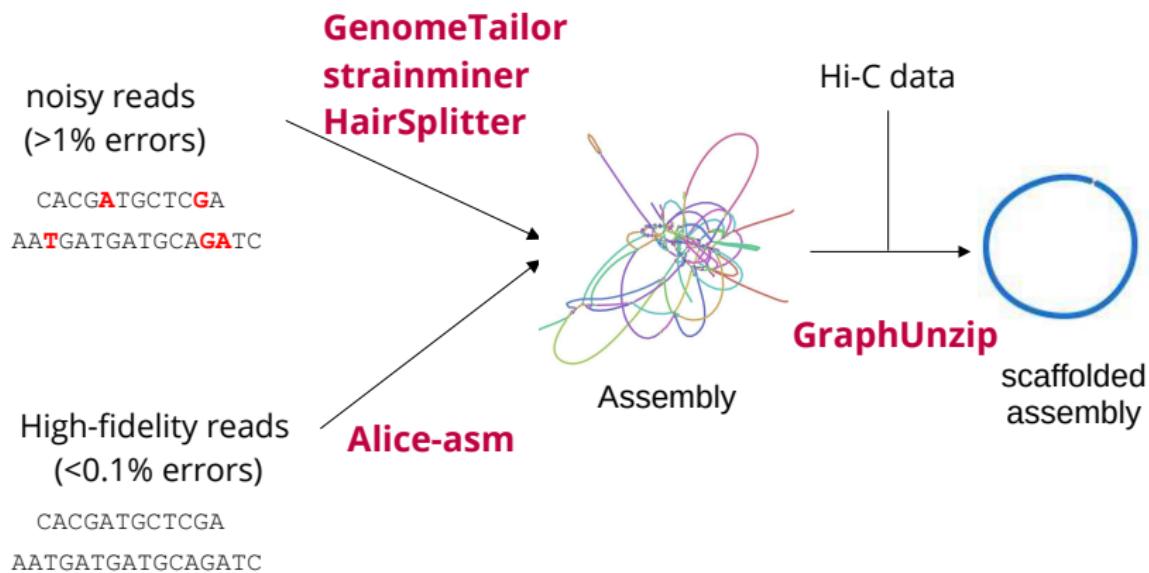
<sup>2</sup>Université de Rennes, IRISA - France

Public Ph.D. defence - November the 27th, 2024

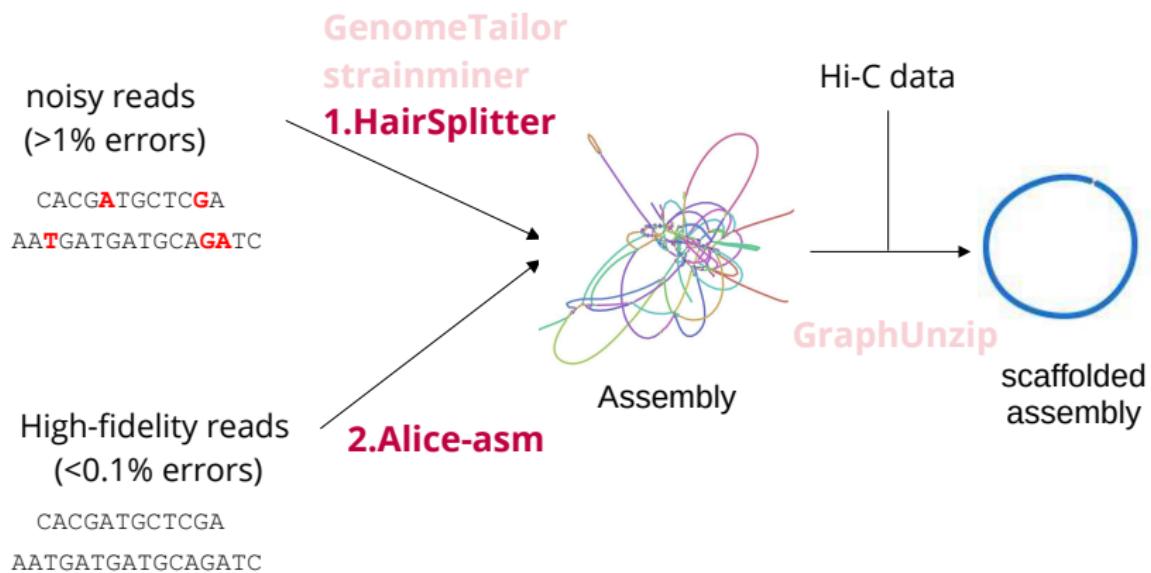
# Overview of (meta)genome assembly



# Overview of the Ph.D.



# Overview of the Ph.D.

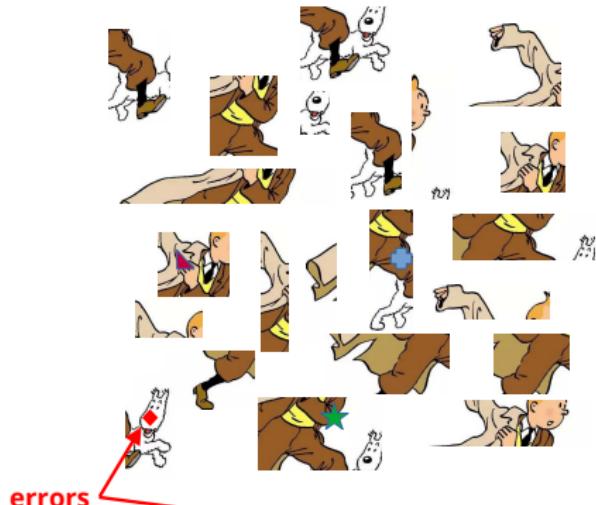


## Distinguishing haplotypes with noisy reads - HairSplitter

## Assemblers correct errors

- ▶ I want to correct sequencing errors but I don't know the solution
- ▶ Like correcting a text in an unknown language
- ▶ Kahoot

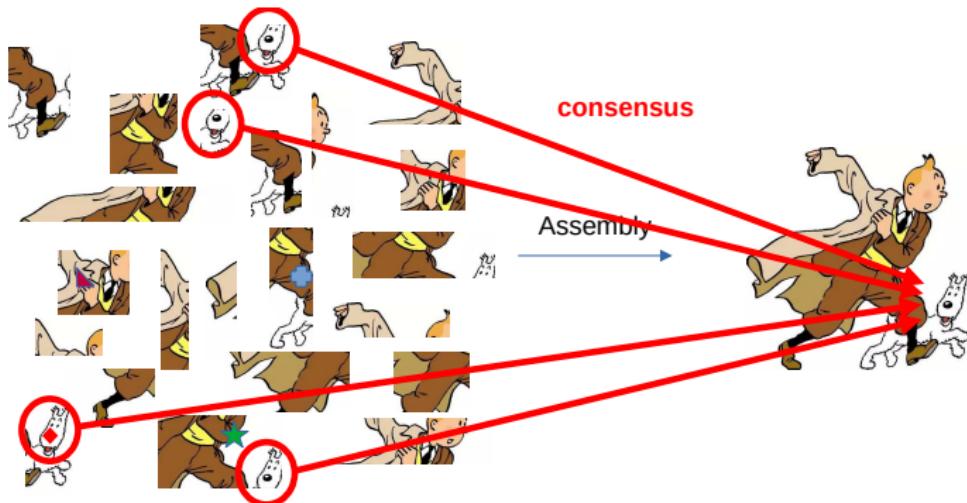
# Assembling noisy reads: correcting errors by consensus



errors

```
CAGCGCTGAGGCAGCAGTGCCA  
CAGCGCTGTGGCAGCAGTGCCA  
CAGCGCTGTGGCAGCAGTGCCA  
CAGCGCTGTGGCAGCAGTGCCA
```

# Assembling noisy reads: correcting errors by consensus

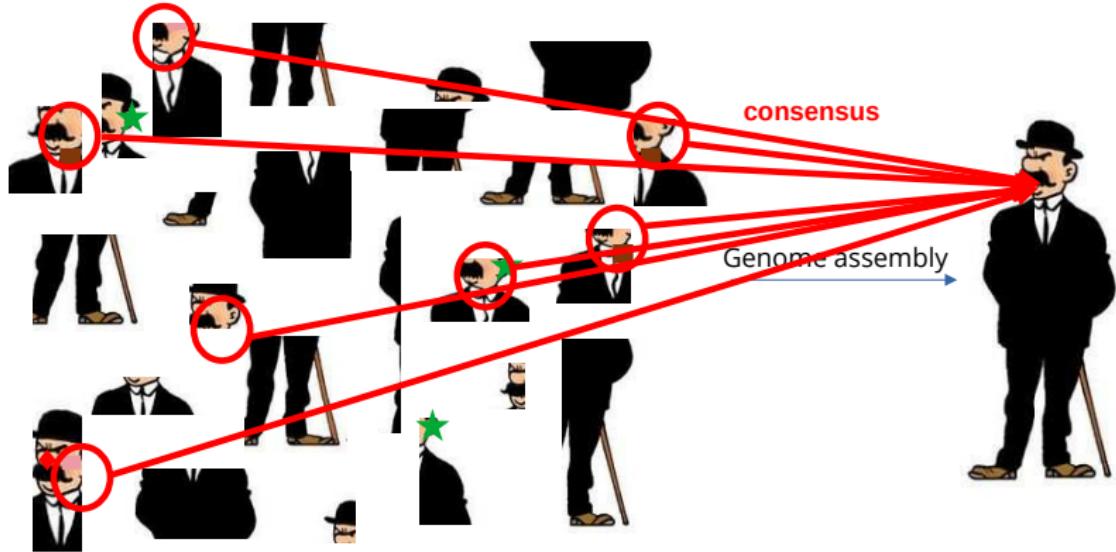


```
CAGCGCTGAGGCAGCAGTGCCA  
CAGCGCTGTGGCAGCAGTGCCA  
CAGCGCTGTGGCAGCAGTGCCA  
CAGCGCTGTGGCAGCAGTGCCA
```

Assembly

```
CAGCGCTGTGGCAGCAGTGCCA
```

## Consensus loses the variants



## Consensus loses the variants

r1 ACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGATAGCA  
r2 AATAAGATAGACGAGATAGACACAGCTTGGCGTTAGGAACAGATGATAGCA  
r3 ACAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACAGATGACAGATAGCA  
r4 ACAAGATCGACGAGATAGACACATCTTGGCGTTAGGAACATTGACAGATAGCA  
r5 ACAAGATCGACAAGATAGGCACATATTGGCGTTAGGAACAGTTGATAGATAGCA  
r6 ACAAGATCGACGAGATAGACACATATTGGCGTTAGGATCAGTTGACAGATAGCA



# Consensus loses the variants

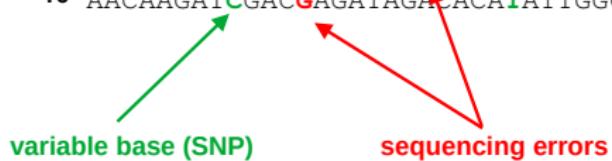
r1 ACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGATAGCA  
 r2 AATAAGATAGACGAGATAGACACAGCTTGGCGTTAGGAACAGATGATAGCA  
 r3 ACAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACAGATGACAGATAGCA  
 r4 ACAAGATCGACGAGATAGACACACTTGGCGTTAGGAACATTTGACAGATAGCA  
 r5 ACAAGATTCGACAAGATAGGCACATTTGGCGTTAGGAACAGTTGATAGATAGCA  
 r6 ACAAGATTCGACGAGATAGACACATTTGGCGTTAGGAATCAGTTGACAGATAGCA

AACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGACAGATAGCA



# How to distinguish errors and SNPs?

r1 ACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGATAGCA  
r2 AATAAGATAGACGAGATAGACACAGCTTGGCGTTAGGAACAGATGATAGCA  
r3 ACAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACAGATGACAGATAGCA  
r4 ACAAGATCGACGAGATAGACACATCTTGGCGTTAGGAACATTGACAGATAGCA  
r5 ACAAGATCGACAAGATAGGCACATATTGGCGTTAGGAACAGTTGATAGATAGCA  
r6 ACAAGATCGACGAGATAGACACATATTGGCGTTAGGATCAGTTGACAGATAGCA



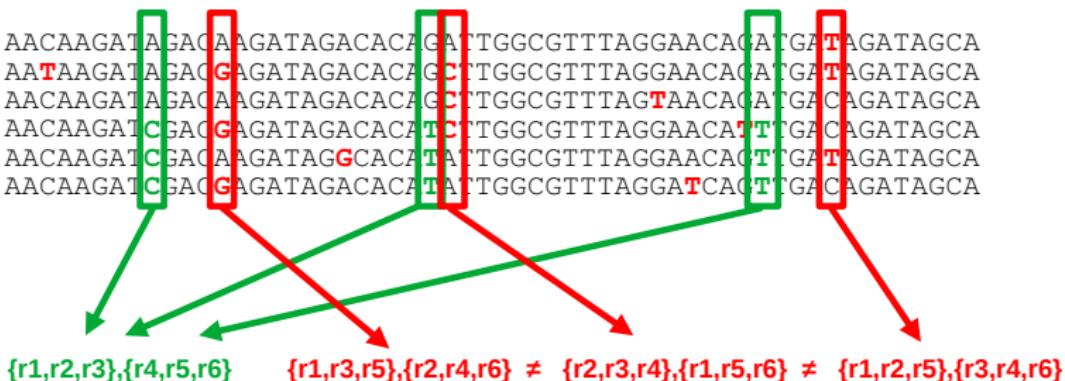
## My solution: looking at several positions simultaneously

**r1** AACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACACAGATGA**T**AGATAGCA  
**r2** AAT**A**GATAGAC**G**GAGATAGACACAGCTTGGCGTTAGGAACACAGATGA**T**AGATAGCA  
**r3** AACAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACACAGACAGATAGCA  
**r4** AACAAGAT**C**GAC**G**GAGATAGACACAT**T**CTTGGCGTTAGGAACAC**T****T**GACAGATAGCA  
**r5** AACAAGAT**C**GACAAGATAG**G**CACAT**T**TATTGGCGTTAGGAACAC**T****T**GATAGATAGCA  
**r6** AACAAGAT**C**GAC**G**GAGATAGACACAT**T**TATTGGCGTTAGGA**T**CACT**T**GACAGATAGCA

{r1,r2,r3},{r4,r5,r6}

## My solution: looking at several positions simultaneously

r1 AACAAGATAGACAAAGATAGACACAGATTGGCGTTAGGAACACAGATGATGAGATAGCA  
 r2 AATTAAGATAGACAAAGATAGACACAGCTTGGCGTTAGGAACACAGATGATGAGATAGCA  
 r3 AACAAGATAGACAAAGATAGACACAGCTTGGCGTTAGTAACACAGATGACAGATAGCA  
 r4 AACAAGATCGACAAAGATAGACACAGCTTGGCGTTAGGAACACATTGACAGATAGCA  
 r5 AACAAGATCGACAAAGATAGGCACACATTGGCGTTAGGAACACATTGATGAGATAGCA  
 r6 AACAAGATCGACAGATAGACACAGATTGGCGTTAGGAATCAGATTGACAGATAGCA



## Algorithm: 1) looking for variant patterns

r1 ACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGATAGCA  
 r2 AATAAGATAGACGAGATAGACACAGCTTGGCGTTAGGAACAGATGATAGCA  
 r3 ACAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACAGATGACAGATAGCA  
 r4 ACAAGATCAGACGAGATAGACACATCTTGGCGTTAGGAACATTTGACAGATAGCA  
 r5 ACAAGATCAGACAAGATAGGCACATATTGGCGTTAGGAACACTTGATAGATAGCA  
 r6 ACAAGATCAGACGAGATAGACACATATTGGCGTTAGGATCACATTGACAGATAGCA

**variant pattern: subset of reads and positions containing alternative bases**  
 size: 3x3

## Algorithm: 1) looking for variant patterns

r1 ACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGATAGCA  
 r2 AATAGATAGAGAGATAGACACAGCTTGGCGTTAGGAACAGATGATAGCA  
 r3 ACAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACAGATGACAGATAGCA  
 r4 ACAAGATCGAGATAGACACAGCTTGGCGTTAGGAACATTGACAGATAGCA  
 r5 ACAAGATCGACAAGATAGGCACATATTGGCGTTAGGAACAGTTGATAGATAGCA  
 r6 ACAAGATCGACGAGATAGACACATATTGGCGTTAGGATCAGTTGACAGATAGCA

**variant pattern: subset of reads and positions containing alternative bases  
 size: 2x2**

# Is this variant too big to be due to sequencing errors?

```

r1 ACAAAGATAGACAAGATAGACAGATGGCGTTAGGAACAGATGATAGATAGCA
r2 AATAGATAGACGAGATAGACACAGCTTGGCGTTAGGAACAGATGATAGATAGCA
r3 ACAAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACAGATGACAGATAGCA
r4 ACAAAGATCGACGAGATAGACACATCTTGGCGTTAGGAACATTTGACAGATAGCA
r5 ACAAAGATCGACAAGATAGGCACATTATTGGCGTTAGGAACAGTTGATTAGATAGCA
r6 ACAAAGATCGACGAGATAGACACATATTGGCGTTAGGATCAGTTGACAGATAGCA
  
```

**variant pattern: subset of reads and positions containing alternative bases  
size: 2x2**

- ▶ I need a statistical test

# Is this variant too big to be due to sequencing errors?

```

r1 ACAAAGATAGACAAGATAGACAGATGGCGTTAGGAACAGATGATAGATAGCA
r2 AATAGATAGACGAGATAGACACAGCTTGGCGTTAGGAACAGATGATAGATAGCA
r3 ACAAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACAGATGACAGATAGCA
r4 ACAAAGATCGACGAGATAGACACATCTTGGCGTTAGGAACATTTGACAGATAGCA
r5 ACAAAGATCGACAAGATAGGCACATTATTGGCGTTAGGAACAGTTGATTAGATAGCA
r6 ACAAAGATCGACGAGATAGACACATATTGGCGTTAGGATCAGTTGACAGATAGCA
  
```

**variant pattern: subset of reads and positions containing alternative bases  
size: 2x2**

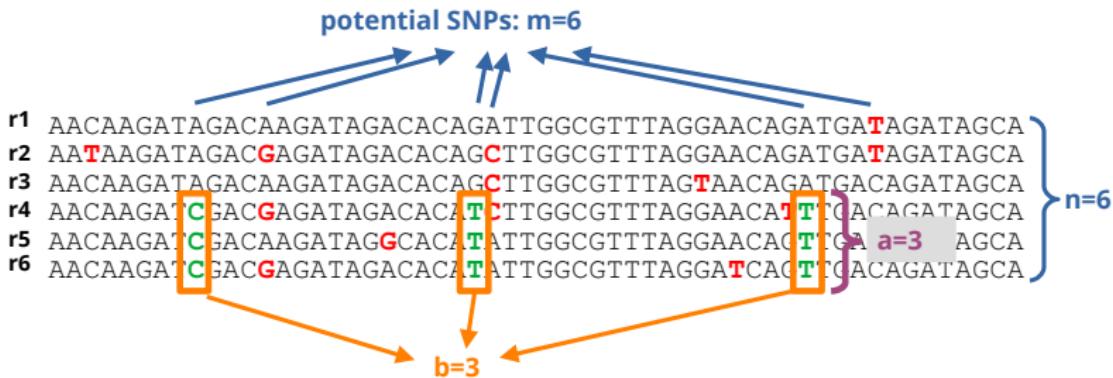
- ▶ I need a statistical test
- ▶ I need help: Kahoot

## Algorithm: 2) Statistical test

- ▶ Null hypothesis: there are no haplotypes, just sequencing errors

## Algorithm: 2) Statistical test

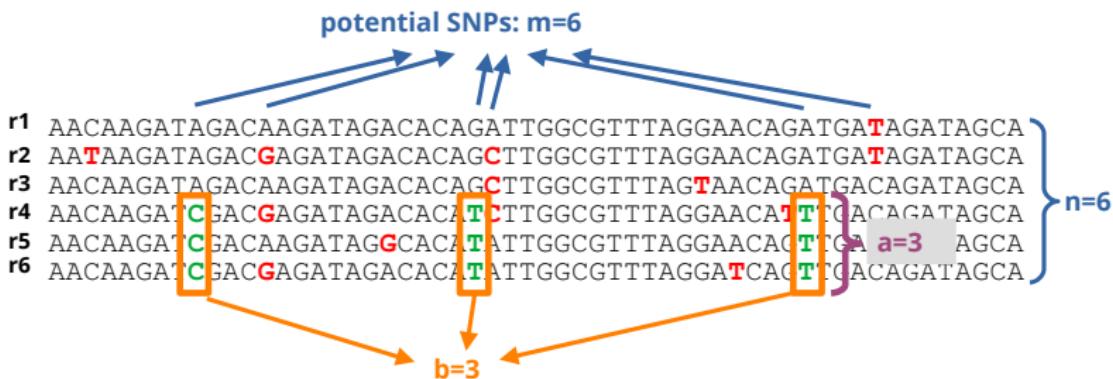
- ▶ Null hypothesis: there are no haplotypes, just sequencing errors



$$P(\text{errors produce pattern of size } ab) \leq \binom{n}{a} \binom{m}{b} * \frac{a^{ab}}{n^{ab}} = 0.006$$

## Algorithm: 2) Statistical test

- ▶ Null hypothesis: there are no haplotypes, just sequencing errors



$$P(\text{errors produce pattern of size } ab) \leq \binom{n}{a} \binom{m}{b} * \frac{a^{ab}}{n^{ab}} = 0.006$$

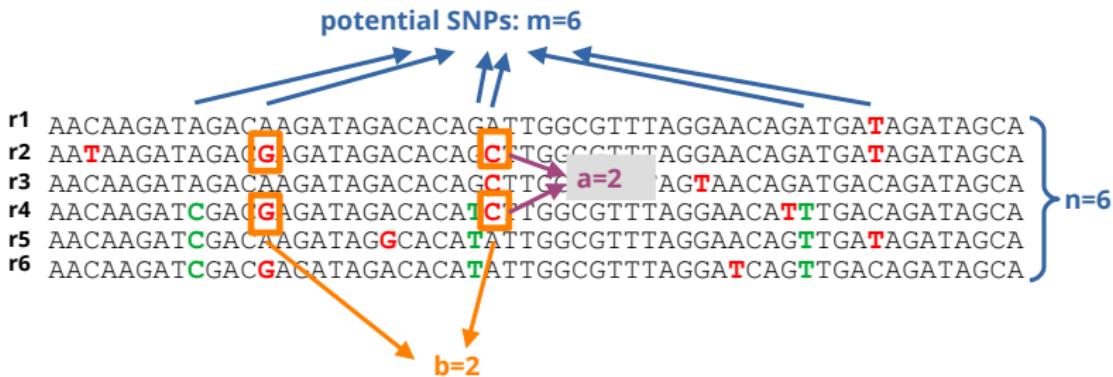
- ▶ THIS LOOKS WEIRD

## Algorithm: 2) Statistical test

- ▶ Null hypothesis: there are no haplotypes, just sequencing errors

## Algorithm: 2) Statistical test

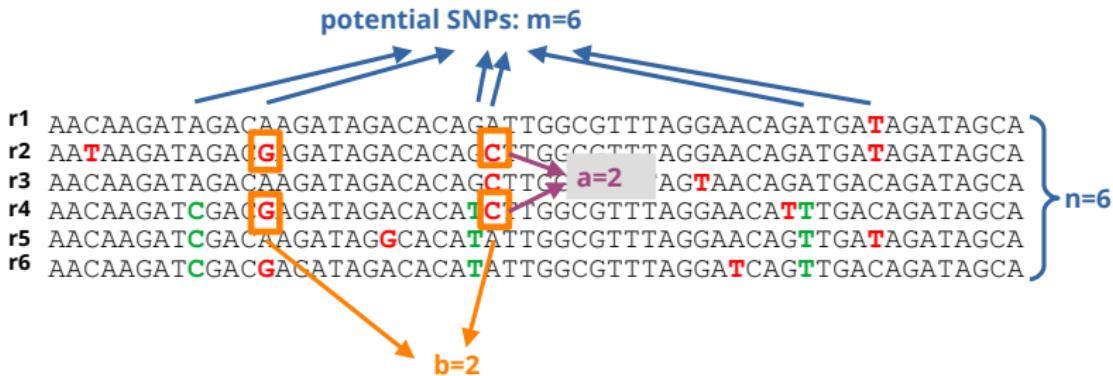
- ▶ Null hypothesis: there are no haplotypes, just sequencing errors



$$P(\text{errors produce pattern of size } ab) \leq \binom{n}{a} \binom{m}{b} * \frac{a^{ab}}{n^{ab}} = 0.30$$

## Algorithm: 2) Statistical test

- ▶ Null hypothesis: there are no haplotypes, just sequencing errors



$$P(\text{errors produce pattern of size } ab) \leq \binom{n}{a} \binom{m}{b} * \frac{a^{ab}}{n^{ab}} = 0.30$$

- ▶ This can happen

## Statistical test: main result

$$\binom{n}{a} \binom{m}{b} * \frac{a^{ab}}{n^{ab}}$$

- ▶ No assumption on the number of haplotypes
- ▶ No assumption on balanced coverage
- ▶ No assumption on the error pattern of the reads
- ▶ Assumption: errors are independent

## Algorithm: 3) Group reads by haplotype

r1 ACAAGATAGACAAGATAGACAGATTGGCGTTAGGAACAGATGATAGATAGCA  
 r2 AATAGATAGACGAGATAGACACAGCTTGGCGTTAGGAACAGATGATAGATAGCA  
 r3 ACAAGATAGACAAGATAGACACAGCTTGGCGTTAGTAACAGATGACAGATAGCA  
 r4 ACAAGATCAGACGAGATAGACACATCTTGGCGTTAGGAACACATGACAGATAGCA  
 r5 ACAAGATCAGACAAGATAGGCACATATTGGCGTTAGGAACACATGATAGATAGCA  
 r6 ACAAGATCAGACGAGATAGACACATATTGGCGTTAGGAATCACATGACAGATAGCA

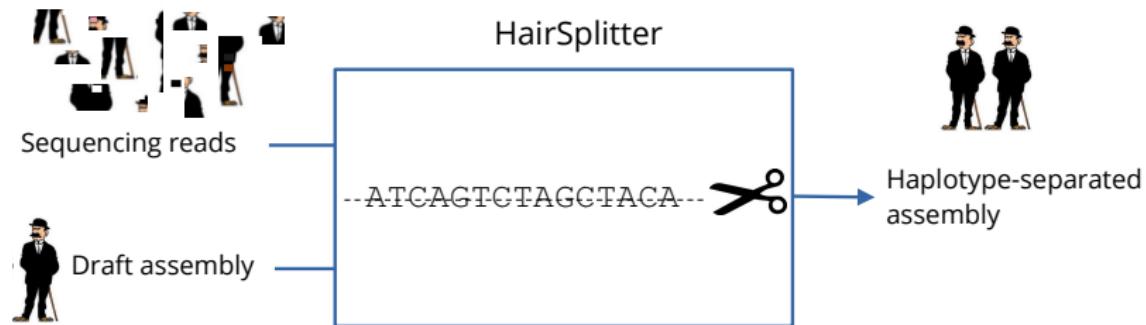
Passed the test



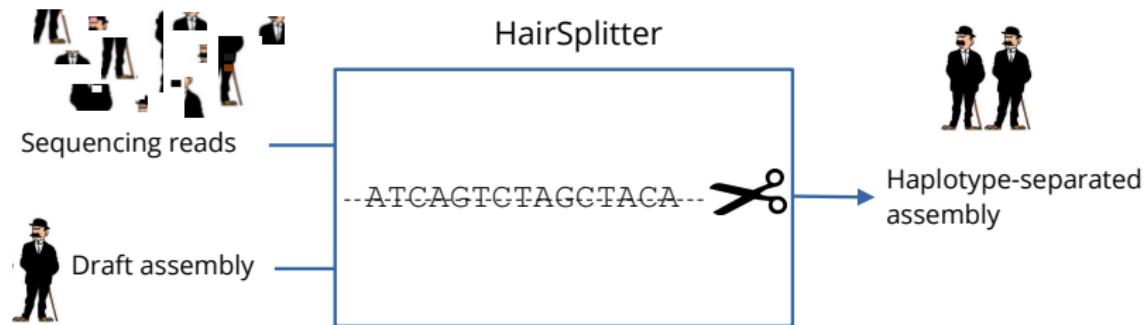
group reads by haplotypes

{r1,r2,r3} {r4,r5,r6}

# The HairSplitter program



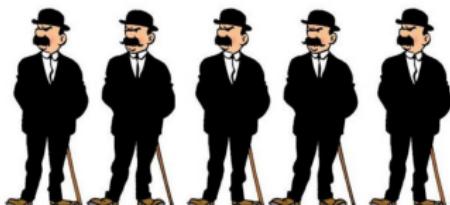
# The HairSplitter program



- ▶ *Hairsplitter*: A person who makes extremely, possibly excessively, fine distinctions (who would separate something as fine as a hair into two pieces and distinguish them) - *Wiktionary*

## Evaluating HairSplitter - results

- ▶ ZymoBIOMICS gut microbiome standard: contains a mix of 5 *E. coli* strains



	metaFlye	metaFlye+Strainberry	metaFlye+HairSplitter
Nanopore Q9	0.586	0.749	<b>0.957</b>
Nanopore Q20	0.7524	0.9527	<b>0.961</b>
PacBio HiFi	0.9589	0.9793	<b>0.9895</b>

Table: 31-mer completeness of assemblies compared to the solution

- ▶ Improves over the state of the art on complex microbiota

# The HairSplitter project

- ▶ Presented at JOBIM, SeqBIM, ISMB/ECCB
- ▶ Published in *Peer Community Journal*

[bioconda / packages / hairsplitter](#) 1.9.10



Recovers collapsed haplotypes from a draft assembly and long reads

Conda

Files

Labels

Badges

License: GPL-3.0-or-later

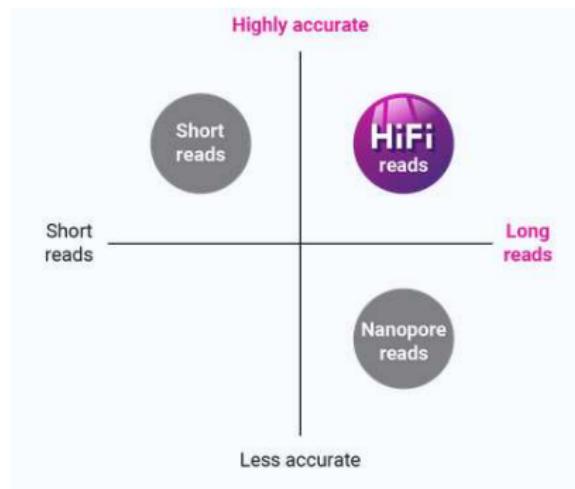
Home: <https://github.com/RolandFaure/HairSplitter>

2169 total downloads

Last upload: 3 months and 6 days ago

# Distinguishing haplotypes with high-fidelity reads - Alice

## New technology: high-fidelity long reads



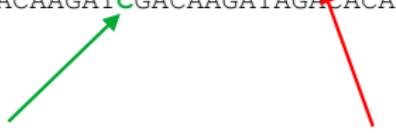
pacb.com

- ▶ Emerged recently and are still emerging
- ▶ << 1% sequencing errors

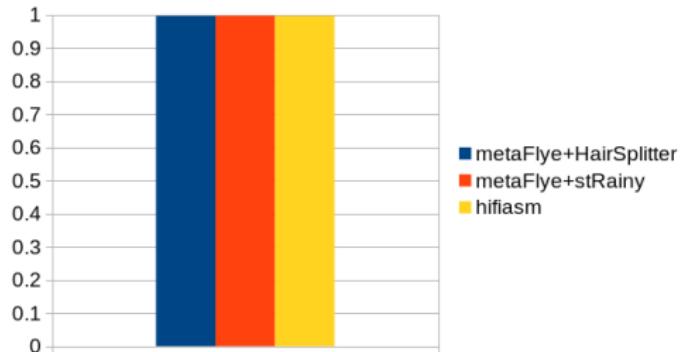
## Assembly with high-fidelity long reads: easy!

{ r1 AACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGACAGATAGCA  
r2 AACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGACAGATAGCA  
r3 AACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGACAGATAGCA  
r4 AACAAGATC GACAAGATAGACACATCTTGGCGTTAGGAACAGTTGACAGATAGCA  
r5 AACAAGATC GACAAGATAGG CACATATTGGCGTTAGGAACAGTTGACAGATAGCA  
r6 AACAAGATC GACAAGATAGACACATATTGGCGTTAGGAACAGTTGACAGATAGCA

variable base (SNP) sequencing error



## Assembly with high-fidelity long reads: easy!



27-mer completeness of the assemblies of the ZymoBIOMICS Gut Microbiome Standard

## Assembly with high-fidelity long reads: slow!

Table: CPU time

	hifiasm	metaFlye+HairSplitter
ZymoBIOMICS Gut Microbiome Standard	20 days	4 days

## Assembly with high-fidelity long reads: slow!

Table: CPU time

	hifiasm	metaFlye+HairSplitter
ZymoBIOMICS Gut Microbiome Standard	20 days	4 days
human genome	34 days	25 days

# Assembly with high-fidelity long reads: slow!

Table: CPU time

	hifiasm	metaFlye+HairSplitter
ZymoBIOMICS Gut Microbiome Standard	20 days	4 days
human genome	34 days	25 days
human gut microbiome <sup>1</sup>	$\geq 60$ days	$\geq 60$ days

<sup>1</sup>Highly accurate metagenome-assembled genomes from human gut microbiota using long-read assembly, binning, and consolidation methods - BiorXiv - Portik et al.

## Assembly with high-fidelity long reads: slow!

Table: CPU time

	hifiasm	metaFlye+HairSplitter
ZymoBIOMICS Gut Microbiome Standard	20 days	4 days
human genome	34 days	25 days
human gut microbiome	$\geq 60$ days	$\geq 60$ days

- ▶ Also takes memory, energy and money

# Why are assemblers so slow?? (Kahoot)

**Read 1**

ATGCATCGAGTAGGGGCACTGTACC

**Read 2**

GAGTAGGGGCACTGTACCAGAGCCAGTAGCAT

**Read 3**

CAGATGGAGAATGCATCGAGTAGG

compute overlaps

**Read 3** CAGATGGAGAATGCATCGAGTAGG

**Read 1** ATGCATCGAGTAGGGGCACTGTACC

**Read 2** GAGTAGGGGCACTGTACCAGAGCCAGTAGCAT

stitch and consensus reads

CAGATGGAGAATGCATCGAGTAGGGGCACTGTACCAGAGCCAGTAGCAT

# Why are assemblers so slow??

**Read 1**

ATGCATCGAGTAGGGGCACTGTACC

**Read 2**

GAGTAGGGGCACTGTACCAGAGCCAGTAGCAT

**Read 3**

CAGATGGAGAATGCATCGAGTAGG

compute overlaps

slow !

**Read 3** CAGATGGAGAATGCATCGAGTAGG**Read 1** ATGCATCGAGTAGGGGCACTGTACC**Read 2** GAGTAGGGGCACTGTACCAGAGCCAGTAGCAT

stitch and consensus reads

CAGATGGAGAATGCATCGAGTAGGGGCACTGTACCAGAGCCAGTAGCAT

The reads and the genome are too long -> let's compress!

- ▶ No compression: CCCGGTTTAA
- ▶ Lossless compression: 3C 2G 3T 2A
- ▶ Lossy compression: CGTA

The reads and the genome are too long -> let's compress!

- ▶ No compression: CCCGGTTTAA
- ▶ Lossless compression: 3C 2G 3T 2A
- ▶ Lossy compression: CGTA
  
- ▶ We want **lossy compression**

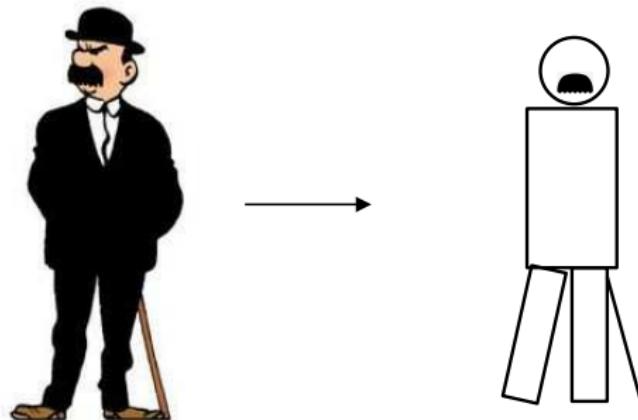
The reads and the genome are too long -> let's compress!

- ▶ No compression: CCCGGTTTAA
- ▶ Lossless compression: 3C 2G 3T 2A
- ▶ Lossy compression: CGTA
  
- ▶ We want **lossy compression**
- ▶ You all use lossy compression formats (kahoot)

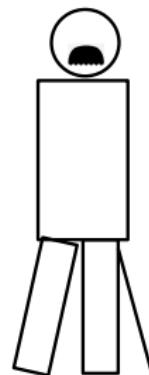
## Lossy compression can be extremely good



# Lossy data representation = sketching



# Lossy data representation = sketching



...CGACGTATGCATCATGCAG...



?

# My contribution: MSR sketching

**Sketching rule: “take all letters after an A”**

sequence

CAGATGGAGAATGCATCGAGTAGGGGCACTGTACCAGAG

# MSR sketching

**Sketching rule: “take all letters after an A”**

sequence

C**A**G**A**TGG**A**GA**AA**TGC**A**TCG**A**GT**A**GGGGC**A**CTGT**A**CC**A**GA**G**

# MSR sketching

**Sketching rule: “take all letters after an A”**

sequence

C **A**G**A**TGG**A**GA**A**TGC**A**TCG**A**GT**A**GGGGC**A**CTGT**A**CC**A**GA**G**  
↓    ↓    ↓    ↓↓    ↓    ↓    G    G    ↓    C    C    G    G  
G    T    G    AT    T    G    G    C    C    G    G

# MSR sketching

**Sketching rule: “take all letters after an A”**

sequence

CAGA**TGGAGA**A TGC**A**TCG**A**GT**A**GGGGC**A**CTGT**A**CC**A**GA**G**

↓    ↓    ↓    ↓    ↓    ↓    ↓    ↓    ↓    ↓    ↓    ↓    ↓

G    T    G    AT    T    G    G    C    C    G    G



sketch

GTGATTGGCCGG

# MSR sketching

**Sketching rule: “take all letters after an A”**

sequence

CAGATGGAGAATGCATCGAGTAGGGGCACTGTACCAGAG

↓      ↓      ↓    ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓  
G    T      G    AT      T      G    G      C      C    G    G



sketch

GTGATTGGCCGG

- ▶ The sketch is 4 times smaller than the original sequence!

# MSR sketching

**Sketching rule: “take all letters after an A”**

sequence

CAGATGGAGAATGCATCGAGTAGGGGCACTGTACCAGAG

↓      ↓      ↓    ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓      ↓  
G    T      G    AT      T      G    G      C      C    G    G



sketch

GTGATTGGCCGG

- ▶ The sketch is 4 times smaller than the original sequence!
- ▶ I can choose another rule if I want

# MSR=Mapping-friendly Sequence Reductions

- ▶ If two reads align, their sketches align too

	CCAGCATCATGCATACGTGCTACTAG	<b>original sequence</b>
	G T T T C C G	<b>sketch</b>
<b>sketch</b>	T T C C G G T	
<b>original sequence</b>	TCATGCATACGTGCTACTAGCGCAGTATG	



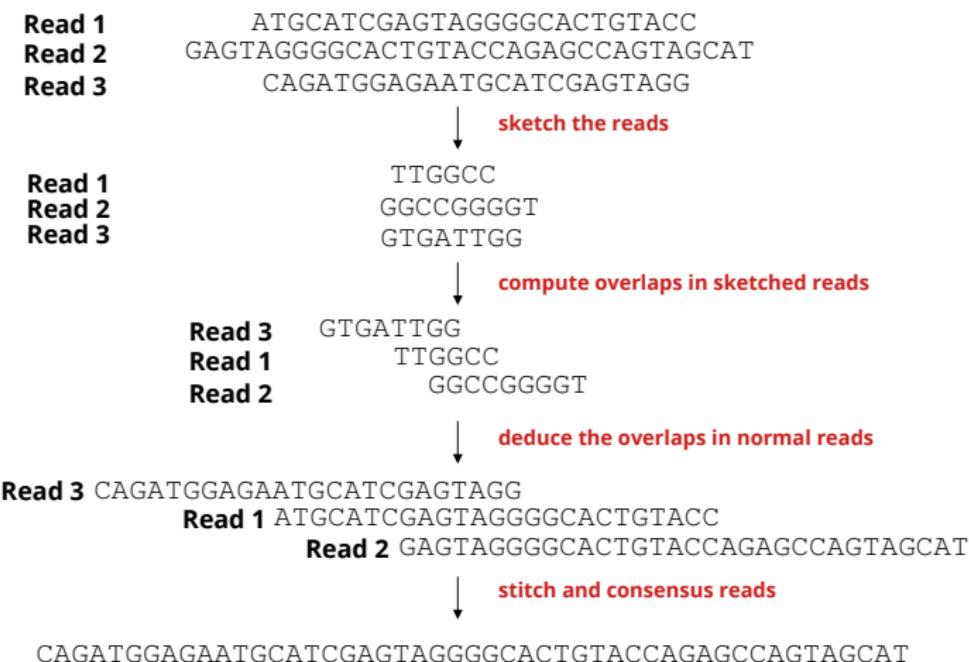
# MSR=Mapping-friendly Sequence Reductions

- ▶ If two reads align, their sketches align too

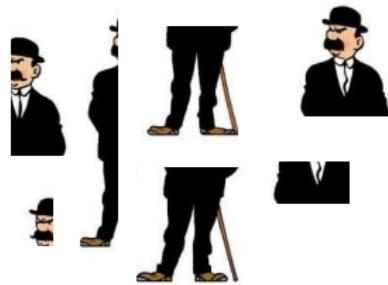
	CCAGCATCATGCATACGTGCTACTAG	<b>original sequence</b>
	G T T T C C G	<b>sketch</b>
<b>sketch</b>	T T C C G G T	
<b>original sequence</b>	TCATGCATACGTGCTACTAGCGCAGTATG	

- ▶ Let's try to compute the overlaps again: Kahoot

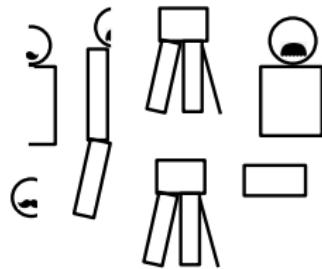
# The new assembly strategy



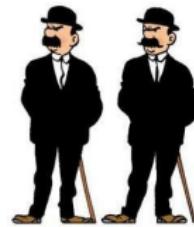
## Solution for fast assembly: sketching the reads



↓ ⌚ €

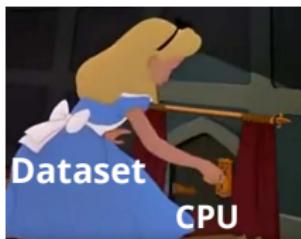
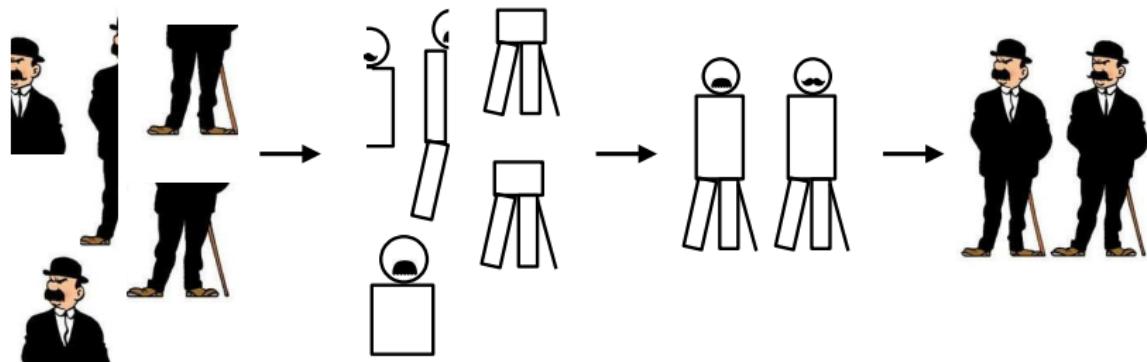


⌚ € →



⌚ € ↑

## Assembling using MSR sketches: the Alice assembler



## Results: Alice assemblies are complete

- ▶ Assembly of the Zymobiomic Gut Microbiome Standard containing 5 strains of *E. coli*



Genome fraction (%)

alice

Escherichia_coli_B1109	92.039
Escherichia_coli_B3008	99.968
Escherichia_coli_B766	95.641
Escherichia_coli_JM109	96.334
Escherichia_coli_b2207	95.495

## Results: Alice assemblies are fast

	hifiasm	metaFlye +HairSplitter	Alice-asm
ZymoBIOMICS Gut Microbiome Standard	20 days	4 days	1h20

## Results: Alice assemblies are fast

	hifiasm	metaFlye +HairSplitter	Alice-asm
ZymoBIOMICS Gut Microbiome Standard	20 days	4 days	1h20
human genome	34 days	25 days	8h40

## Results: Alice assemblies are fast

	hifiasm	metaFlye +HairSplitter	Alice-asm
ZymoBIOMICS Gut Microbiome Standard	20 days	4 days	1h20
human genome	34 days	25 days	8h40
human gut microbiome	≥ 60 days	≥ 60 days	5h00

## Conclusion

## Conclusion: achievements

- ▶ **Noisy reads:** assemble a mix of haplotypes of unprecedented complexity
- ▶ **High-fidelity reads:** assemble very fast while keeping haplotypes with MSR sketching
- ▶ **Hi-C data:** improved the scaffolding of haploid and multiploid assemblies

# Was this thesis really useful?

Organism, microbiota...

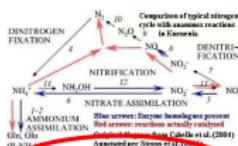


DNA extraction  
& preparation

sample



Sequencing



Understand things

...AACCTGCGTCACGTAGTCGAGG...  
...CGGGCCTGAGGCAGCAGTGCCA...

Genomes or  
amplicons

Assembly

GTGCTAATCACGT  
TCCGAGCGATCAG  
TGTCTGAAACCACA  
CTCTGGGGTGACA

Long reads

# What is the future of assembly?

Organism, microbiota...

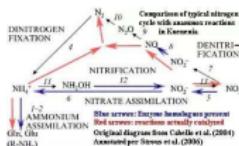


DNA extraction & preparation

sample



Sequencing



Understand things

...AACCTGCGTCACGTAGTCGAGG...  
...CGGGCCTGAGGCAGCAGTGCCA...

Genomes or amplicons

Assembly

CGTAGCTAGGAT  
GTGCTAATCACGT  
TCCGAGCGATCAG  
TGTCTGAAACCACA  
CTCTGGGGTGACA

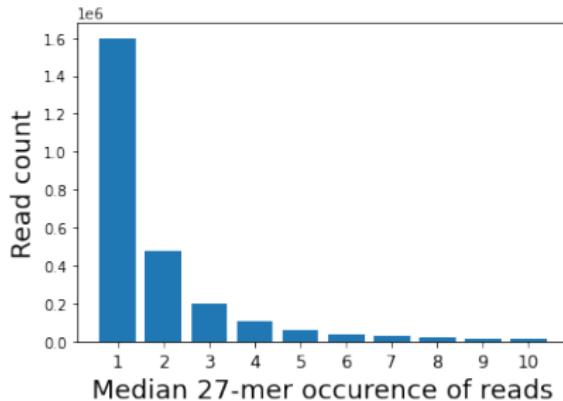
Long reads

# All DNA is not captured by sequencing



Nicolas Maurice

- ▶ Example of a sequencing of the soil microbiota



Adapted from the work of Nicolas Maurice

- ▶ Low-coverage sequencing
- ▶ Missing DNA

# What is the future of assembly?

Organism, microbiota...

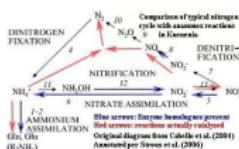


DNA extraction  
& preparation

sample



Sequencing



Understand things

...AACCTGGTGTACGTAGTCGAGG...  
...CGGGCCTGAGGCAGCAGTGCCA...

Genomes or  
amplicons

Assembly

CGTAGCTAGGAT  
GTGCTAATCACGT  
TCCGAGCGATCAG  
TGTCTGAAACCACA  
CTCTGGGGTGACA

Long reads

## Conclusion: let's finish this kahoot

- ▶ Three general questions
- ▶ Double points