Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Haplotype assembly from long reads

## Roland Faure[1,2]

[1]Université Libre de Bruxelles (ULB) - Belgium

[2]Université de Rennes, IRISA - France

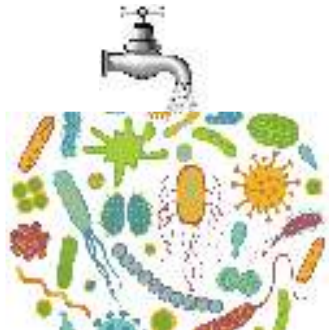Public Ph.D. defence - November the 27th, 2024

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## Microbiota



credits: microbiozindia.com

▶ Microbiota are mixes of bacteria, virus, archea and eukaryota

Roland Faure     Public defence

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Microbiota are everywhere

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## Microbiota are everywhere

Roland Faure    Public defence

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## Microbiota are everywhere

**Introduction**
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Microbiota are everywhere

**Introduction**
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Microbiota are important



The gut microbiome and mental health: advances in research and emerging priorities - Shoubridge et al.

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Microbiota are important



The gut microbiome and mental health: advances in research and emerging priorities - Shoubridge et al.



The Role of Soil Microorganisms in Plant Mineral Nutrition—Current Knowledge and Future Directions - Jacoby et al.

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Microbiota are important



The gut microbiome and mental health: advances in research and emerging priorities - Shoubridge et al.



The Role of Soil Microorganisms in Plant Mineral Nutrition—Current Knowledge and Future Directions - Jacoby et al.



Core microbiota drive functional stability of soil microbiome in reforestation ecosystems - Jiao et al.

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## How to study microbiota?



Anthonie van Leeuwenhoek
1673

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# How to study microbiota?

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# How to study microbiota?



▶ But many microorganisms are undistiguishable under a microscope

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## How to study microbiota?



Julius Petri
1887

Roland Faure    Public defence

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## How to study microbiota?



Julius Petri
1887

▶ But most microorganisms are not cultivable

Roland Faure     Public defence

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# How to study microbiota?



Ecosystem

DNA extraction

Sequencing

Frederick Sanger
1975

```
                                              CGTAGCTAGGAT
                                   GTGCTAATCACGT  AAAGCTAATCACTT
...AACGTGCGTCACGTAGTCGAGG...       TGCGAGCGATCAG   TGTCTGAAACCACA
...CGGCGCTGAGGCAGCAGTGCCA...       CTCTGGGGTGACA
```

Assembly

Understand the
microbiome

Genomes

Long reads

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# How to study microbiota?



Microbiome

DNA extraction
& preparation

Sequencing

Frederick Sanger
1975

sample

**My Ph.D.**

...AACGTGCGTCACGTAGTCGAGG...

...CGGCGCTGAGGCAGCAGTGCCA...

Assembly

Genomes

Understand the
microbiome

CGTAGCTAGGAT
GTGCTAATCACGT AAAGCTAATCACTT
TGCGAGCGATCAG
CTCTGGGGTGACA TGTCTGAAACCACA

Long reads

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## DNA sequencing

Extracted
DNA



CAGCATCAGTTTTCGAGCACGT

TTACTCAGCAGATCGTCGATCAT

CCCGTAGCTTAGCAGGCATCAG

**Reads**

sequencer

Roland Faure    Public defence

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# DNA sequencing: difficulties



**length: 1-20 kbp**

CAGCATCAGTTTTCGAGCACGT

TTACTCAGCAGATCGTCGATCAT

CCCGTAGCTTAGCAGGCATCAG

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# DNA sequencing: difficulties



sequencing errors: 0.1 – 10 %

**Introduction**
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Genome sequencing

Roland Faure     Public defence

**Introduction**
Distiguishing haplotypes with noisy reads – HairSplitter
Distinguishing haplotypes with high-fidelity reads – Alice
Conclusion

# Genome assembly



```
                    TTCGGCGCTGAGGCAG
    CAGCGCTGAGGCAGCAGTGCCA   GGCAGCAGTGCCAG                    Assembly
                                                         ─────────────→    ...CGCTGAGGCAGCATGTGCCAGGCT...
CAGCATTGCCAGGC    TCGGC
                    CGGCGCTGAGGCAGCATTGCC
```

Assembly

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Genome assembly

```
CGATGCTGGCTAGCATAGTCGATTTATCT
     CTGGCTAGCTTAGTCGATTTATCTGACAGT
            AGCATAGTCGATTTATCTGACAGTCATAT
                 AGTCGATTTATATGACAGTCATATTGCT
                    TTTATCTGACAGTCAGATTGCTACACAC
```

genome assembly: stitching reads
correcting errors

**CGATGCTGGCTAGCATAGTCGATTTATCTGACAGTCATATTGCTACACAC**

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

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## Imagine you are an assembler

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## *Haplotype* assembly



Collapsed assembly

Haplotype assembly

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Dupont & Dupond exist in microbiota!



*Escherichia coli* Sakai

...ACACACCACACACCTCTACGA...

...ACACAC**T**ACACACCTCTACGA...

*Escherichia coli* Nissle

**Introduction**
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Problem: assembling several haplotypes

```
CGATGCTGGCTAGCATAGTCGATTTATCT
      CTGGCTAGCTTAGTCGATTTATCTGACAGT
            AGCATAGTCGATTTATCTGACAGTCATAT
                  AGTCGATTTATATGACAGTCATATTGCT
                     TTTATATGACAGTCAGATTGCTACACAC
```

genome assembly: stitching reads
correcting errors

**CGATGCTGGCTAGCATAGTCGATTTATCTGACAGTCATATTGCTACACAC**

**CGATGCTGGCTAGCATAGTCGATTTATATGACAGTCATATTGCTACACAC**

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## Problem: assembling several haplotypes

```
CGATGCTGGCTAGCATAGTCGATTTATCT
     CTGGCTAGCTTAGTCGATTTATCTGACAGT
          AGCATAGTCGATTTATCTGACAGTCATAT
               AGTCGATTTATATGACAGTCATATTGCT
                    TTTATATGACAGTCAGATTGCTACACAC
```

genome assembly: stitching reads
correcting errors

**CGATGCTGGCTAGCATAGTCGATTTATCTGACAGTCATATTGCTACACAC**

**CGATGCTGGCTAGCATAGTCGATTTATATGACAGTCATATTGCTACACAC**

▶ Not so many software!

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

Haplotype assembly from long reads

Roland Faure[1,2]

[1]Université Libre de Bruxelles (ULB) - Belgium

[2]Université de Rennes, IRISA - France

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

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Overview of (meta)genome assembly

noisy reads
(>1% errors)

CACG**A**TGCTC**G**A

AA**T**GATGATGCA**GA**TC

Hi-C data

Assembly

scaffolded
assembly

High-fidelity reads
(<0.1% errors)

CACGATGCTCGA

AATGATGATGCAGATC

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## Overview of the Ph.D.



noisy reads
(>1% errors)

CACG**A**TGCTC**G**A

AA**T**GATGATGCA**GA**TC

**GenomeTailor
strainminer
HairSplitter**

Hi-C data

Assembly

**GraphUnzip**

scaffolded
assembly

High-fidelity reads
(<0.1% errors)

CACGATGCTCGA

AATGATGATGCAGATC

**Alice-asm**

**Introduction**
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## Overview of the Ph.D.

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Distiguishing haplotypes with noisy reads - HairSplitter

Roland Faure     Public defence

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Assembling noisy reads: correcting errors by consensus



**errors**

```
CAGCGCTGAGGCAGCAGTGCCA
CAGCGCTGTGGCAGCAGTGCCA
CAGCGCTGTGGCAGCAGTGCCA
CAGCGCTGTGGCAGCAGTGCCA
```

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Assembling noisy reads: correcting errors by consensus



**consensus**

Assembly

```
CAGCGCTGAGGCAGCAGTGCCA
CAGCGCTGTGGCAGCAGTGCCA
CAGCGCTGTGGCAGCAGTGCCA
CAGCGCTGTGGCAGCAGTGCCA
```

Assembly

```
CAGCGCTGTGGCAGCAGTGCCA
```

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Consensus loses the variants



consensus

Genome assembly

Introduction
**Distinguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Consensus loses the variants

```
r1  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGATAGATAGCA
r2  AATAAGATAGACGAGATAGACACAGCTTGGCGTTTAGGAACAGATGATAGATAGCA
r3  AACAAGATAGACAAGATAGACACAGCTTGGCGTTTAGTAACAGATGACAGATAGCA
r4  AACAAGATCGACGAGATAGACACATCTTGGCGTTTAGGAACATTTGACAGATAGCA
r5  AACAAGATCGACAAGATAGGCACATTATTGGCGTTTAGGAACAGTTGATAGATAGCA
r6  AACAAGATCGACGAGATAGACACATTATTGGCGTTTAGGATCAGTTGACAGATAGCA
```

**variable base (SNP)**          **sequencing errors**

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Consensus loses the variants

```
r1  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGATAGATAGCA
r2  AATAAGATAGACGAGATAGACACAGCTTGGCGTTTAGGAACAGATGATAGATAGCA
r3  AACAAGATAGACAAGATAGACACACTTGGCGTTTAGTAACAGATGACAGATAGCA
r4  AACAAGACCGACGAGATAGACACTTTGGCGTTTAGGAACTTGGACAGATAGCA
r5  AACAAGACCGACAAGATAGGCACTTTGGCGTTTAGGAACTTGATAGATAGCA
r6  AACAAGACCGACGAGATAGACACTTTGGCGTTTAGGATCATTGGACAGATAGCA
```

```
AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGACAGATAGCA
```

**lost SNPs**

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
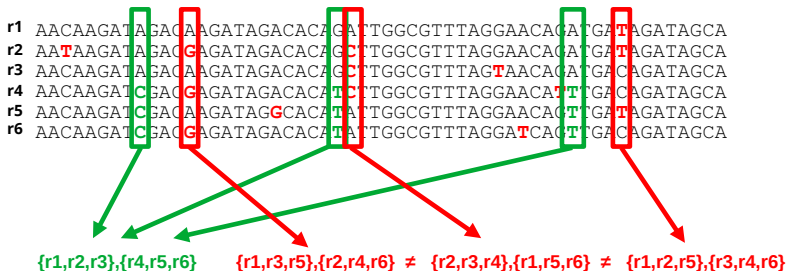Conclusion

# Haplotype separation: state of the art [1]

```
r1  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGATAGATAGCA
r2  AATAAGATAGACGAGATAGACACAGCTTGGCGTTTAGGAACAGATGATAGATAGCA
r3  AACAAGATAGACAAGATAGACACAGCTTGGCGTTTAGTAACAGATGACAGATAGCA
r4  AACAAGATCGACGAGATAGACACATCTTGGCGTTTAGGAACATTTGACAGATAGCA
r5  AACAAGATCGACAAGATAGGCACATATTGGCGTTTAGGAACAGTTGATAGATAGCA
r6  AACAAGATCGACGAGATAGACACATATTGGCGTTTAGGATCAGTTGACAGATAGCA
```

3 diffs

5 diffs

**Reads from the same haplotype are more similar than reads from different haplotypes**

---

[1]WhatsHap, HapCut, Strainberry, stRainy...

Introduction
**Distinguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Haplotype separation: state of the art [1]



**Reads from the same haplotype are more similar than reads from
different haplotypes
on average**

---

[1]WhatsHap, HapCut, Strainberry, stRainy...

Roland Faure    Public defence

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# How to distinguish errors and SNPs?

```
r1  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGATAGATAGCA
r2  AATAAGATAGACGAGATAGACACAGCTTGGCGTTTAGGAACAGATGATAGATAGCA
r3  AACAAGATAGACAAGATAGACACAGCTTGGCGTTTAGTAACAGATGACAGATAGCA
r4  AACAAGATCGACGAGATAGACACATCTTGGCGTTTAGGAACATTTGACAGATAGCA
r5  AACAAGATCGACAAGATAGGCACATATTGGCGTTTAGGAACAGTTGATTAGATAGCA
r6  AACAAGATCGACGAGATAGACACATATTGGCGTTTAGGATCAGTTGACAGATAGCA
```

**variable base (SNP)**          **sequencing errors**

Introduction
**Distinguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# My solution: looking at several positions simultaneously

Introduction
**Distinguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# My solution: looking at several positions simultaneously
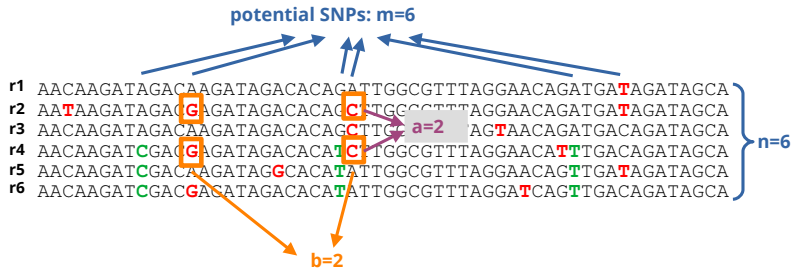


```
r1  AACAAGATAGACAAGATAGACACAGACTGGCGTTTAGGAACAGATGATAGATAGCA
r2  AATTAAGATAGACGCAGATAGACACAGCTGGCGTTTAGGAACAGATGATAGATAGCA
r3  AACAAGATAGACAAGATAGACACAGCTTGGCGTTTAGTTAACAGATGACAGATAGCA
r4  AACAAGATCGACGCAGATAGACACATTGTGGCGTTTAGGAACATTGACAGATAGCA
r5  AACAAGATCGACAAGATAGGCACATTATGGCGTTTAGGAACAGTTGATTAGATAGCA
r6  AACAAGATCGACGCAGATAGACACATTTGGCGTTTAGGATTCAGTTGACAGATAGCA
```

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

Introduction
**Distinguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Algorithm: 1) looking for variant patterns



```
r1   AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGATAGATAGCA
r2   AATAAGATAGACGAGATAGACACAGCTTGGCGTTTAGGAACAGATGATAGATAGCA
r3   AACAAGATAGACAAGATAGACACAGCTTGGCGTTTAGTAACAGATGACAGATAGCA
r4   AACAAGATCACGGAGATAGACACTTTTGGCGTTTAGGAACATTTGACAGATAGCA
r5   AACAAGATCACAAGATAGGCACTTTTGGCGTTTAGGAACACTTTGATAGATAGCA
r6   AACAAGATCACGGAGATAGACACTTTTGGCGTTTAGGATCATTTGACAGATAGCA
```
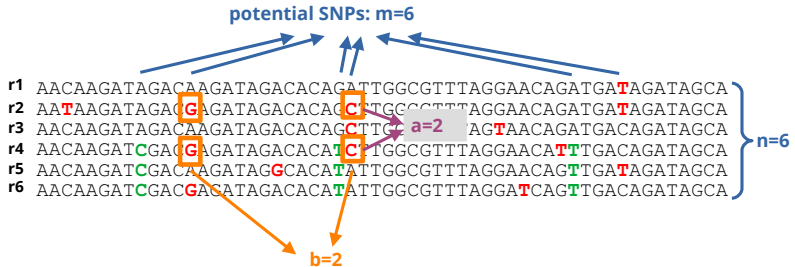
**variant pattern: subset of reads and positions containing minority bases**
**size: 3x3**

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Algorithm: 1) looking for variant patterns



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

Introduction
**Distinguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Algorithm: 2) Statistical test

```
r1  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGATAGATAGCA
r2  AATAAGATAGAGGAGATAGACACACTGGCGTTTAGGAACAGATGATAGATAGCA
r3  AACAAGATAGACAAGATAGACACACCTTGGCGTTTAGTAACAGATGACAGATAGCA
r4  AACAAGATCGACGAGATAGACACCTGGCGTTTAGGAACATTTGACAGATAGCA
r5  AACAAGATCCGACAAGATAGGCACATATTGGCGTTTAGGAACAGTTGATAGATAGCA
r6  AACAAGATCGACGAGATAGACACATATTGGCGTTTAGGATCAGTTGACAGATAGCA
```

Is this pattern too big to be due to errors ?

Introduction
**Distinguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Algorithm: 2) Statistical test



Is this pattern too big to be due to errors ?

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Algorithm: 2) Statistical test



**potential SNPs: m=6**

```
r1  AACAAGATAGACAAGATAGACACAGTTGGCGTTTAGGAACAGATGATAGATAGCA
r2  AATAAGATAGACGAGATAGACACACCTTGGCGTTTAGGAACAGATGATAGATAGCA
r3  AACAAGATAGACAAGATAGACACACCTTCGTGAGTAACAGATGACAGATAGCA
r4  AACAAGATCGACGAGATAGACACATCTGGCGTTTAGGAACATTTGACAGATAGCA
r5  AACAAGATCGACAGATAGGCACATTTTGGCGTTTAGGAACAGTTGATAGATAGCA
r6  AACAAGATCGACGACATAGACACATTTTGGCGTTTAGGATCAGTTGACAGATAGCA
```
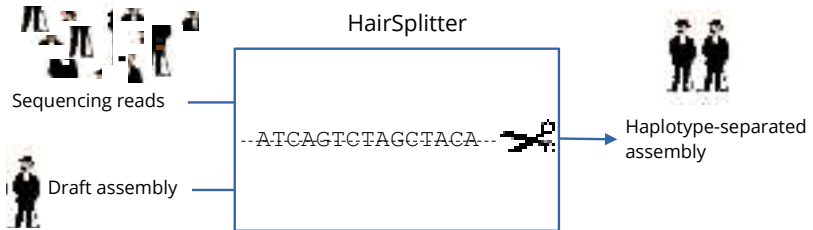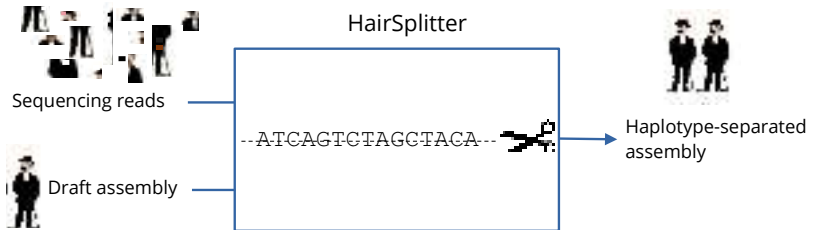
**a=2**

**b=2**

**n=6**

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

Introduction
**Distinguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Algorithm: 2) Statistical test



potential SNPs: m=6

r1 AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGA**T**AGATAGCA
r2 AA**T**AAGATAGAC**G**GAGATAGACACAG**C**TTGGCGTTTAGGAACAGATGA**T**AGATAGCA
r3 AACAAGATAGACAAGATAGACACAG**C**TTGGCGTTTAG**T**AACAGATGACAGATAGCA
r4 AACAAGA**C**AC**G**GAGATAGACAC**T**TTGGCGTTTAGGAACA**T**TGACAGATAGCA
r5 AACAAGA**C**ACAAGATAG**G**CAC**T**TTGGCGTTTAGGAACA**T**TG AGCA
r6 AACAAGA**C**AC**G**GAGATAGACAC**T**TTGGCGTTTAGGA**T**CA**T**TACAGATAGCA

n=6

a=3

b=3

P(errors produce pattern of size **ab**) $\leq \binom{n}{a}\binom{m}{b} * \dfrac{a^{ab}}{n^{ab}} = 0.07$

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## 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

Roland Faure    Public defence

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Algorithm: 3) Group reads by haplotype

```
r1  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGATAGATAGCA
r2  AATAAGATAGACGAGATAGACACAGCTTGGCGTTTAGGAACAGATGATAGATAGCA
r3  AACAAGATAGACAAGATAGACACACCTTGGCGTTTAGTAACAGATGACAGATAGCA
r4  AACAAGATCACGAGATAGACACTCTTGGCGTTTAGGAACATTGACAGATAGCA
r5  AACAAGATCACAAGATAGGCACTTTTGGCGTTTAGGAACATTTATAGATAGCA
r6  AACAAGATCACGAGATAGACACTTTTGGCGTTTAGGATTCATTGACAGATAGCA
```

**Passed the test**

↓ **group reads by haplotypes**

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

Roland Faure    Public defence

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## The HairSplitter program



HairSplitter

Sequencing reads

--ATCAGTCTAGCTACA--

Draft assembly

Haplotype-separated
assembly

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## The HairSplitter program



Sequencing reads

Draft assembly

HairSplitter

--ATCAGTCTAGCTACA--

Haplotype-separated assembly

▶ *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*

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## Let's evaluate HairSplitter - k-mer completeness

Assembly

Solution

**ACGCAGCTAGTACGCAT**

**GCAGCTAGTACGCATAA**

ACGCAGCTAG
 CGCAGCTAGT
  GCAGCTAGTA
   CAGCTAGTAC
    AGCTAGTACG
     GCTAGTACGC
      CTAGTACGCA
       TAGTACGCAT

**GCAGCTAGTA**
 **CAGCTAGTAC**
  **AGCTAGTACG**
   **GCTAGTACGC**
    **CTAGTACGCA**
     **TAGTACGCAT**
      **AGTACGCATA**
       **GTACGCATAA**

10-mer completeness:
6 out of 8 (75%)

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## 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                | **0.957**             |
| Nanopore Q20 | 0.7524   | 0.9527               | **0.961**             |
| PacBio HiFi  | 0.9589   | 0.9793               | **0.9895**            |

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

▶ Improves over the state of the art on complex assemblies

Introduction
**Distiguishing haplotypes with noisy reads - HairSplitter**
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

## The HairSplitter project

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

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# Distinguishing haplotypes with high-fidelity reads - Alice

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# New technology: high-fidelity long reads



pacb.com

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

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# Assembly with high-fidelity long reads: easy!



```
r1  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGACAGATAGCA
r2  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGACAGATAGCA
r3  AACAAGATAGACAAGATAGACACAGATTGGCGTTTAGGAACAGATGACAGATAGCA
r4  AACAAGATCGACAAGATAGACACATCTTGGCGTTTAGGAACAGTTGACAGATAGCA
r5  AACAAGATCGACAAGATAGGCACATTATTGGCGTTTAGGAACAGTTGACAGATAGCA
r6  AACAAGATCGACAAGATAGACACATTATTGGCGTTTAGGAACAGTTGACAGATAGCA
```

**variable base (SNP)**          **sequencing error**

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# Assembly with high-fidelity long reads: easy!



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

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

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

Table: CPU time

|                                          | hifiasm | metaFlye+HairSplitter |
| ---------------------------------------- | ------- | --------------------- |
| Zymobiomics Gut Microbiome Standard      | 20 days | 4 days                |

Roland Faure     Public defence

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## 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               |

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# 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[1] | $\geq$ 60 days | $\geq$ 60 days |

---

[1]Highly accurate metagenome-assembled genomes from human gut microbiota using long-read assembly, binning, and consolidation methods - BiorXiv

Roland Faure    Public defence

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## How to perform fast assembly?



Credits: Alice in Wonderland, Lewis, Disney

Introduction
Distiguishing haplotypes with noisy reads – HairSplitter
**Distinguishing haplotypes with high-fidelity reads – Alice**
Conclusion

# Solution for fast assembly: sketching the reads



Credits: Alice in Wonderland, Lewis, Disney

**Drink-me potion**

**sketched dataset**

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# Sketching: reducing the size of the data

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## Sketching: reducing the size of the data

Roland Faure    Public defence

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## Sketching: reducing the size of the data

State-of-the art
sketch

What we
want

...CGACGTATGCATCATGCAG... ⟶ ?

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# My contribution: MSR sketching

sequence    CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

sequence      CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Roland Faure    Public defence

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

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

sequence        CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

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

sequence      **CAGTATGGAT**ACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**CAGTATGGAT**)= 0.0023
f(**CAGTATGGAT**)= A

sketch         A

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching

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

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

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

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

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

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

sequence      C**AGTATGGATA**CAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**AGTATGGATA**)= 0.624

f(**AGTATGGATA**)= Ø

sketch      A

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

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

sequence      CA**GTATGGATAC**AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**GTATGGATAC**)= 0.124
f(**GTATGGATAC**)= G

sketch      A G

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching

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

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

sequence    CAG**TATGGATACA**GATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**TATGGATACA**)= 0.88
f(**TATGGATACA**)= ∅

sketch    A  G

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

$$f : \{A, C, G, T\}^{10} \to \{A, C, G, T, \varnothing\}$$

$f(10-mer) \to A$   $if$   $hash(10-mer) \in [0, 0.05]$
$f(10-mer) \to C$   $if$   $hash(10-mer) \in [0.05, 0.1]$
$f(10-mer) \to G$   $if$   $hash(10-mer) \in [0.1, 0.15]$
$f(10-mer) \to T$   $if$   $hash(10-mer) \in [0.15, 0.2]$
$f(10-mer) \to \varnothing$   $if$   $hash(10-mer) > 0.2$

sequence      CAGT**ATGGATACAG**ATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**ATGGATACAG**)= 0.32
    f(**ATGGATACAG**)= ∅

sketch      A  G

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

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

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

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

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

$f(10-mer) \rightarrow \emptyset \quad if \quad hash(10-mer) > 0.2$

sequence      CAGTA**TGGATACAGA**TGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**TGGATACAGA**)= 0.19

f(**TGGATACAGA**)= T

sketch      A  G   T

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

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

sequence      CAGTAT**GGATACAGAT**GGAGATATCATCGAGTAGGGGCACTGTACCAGAG

      hash(**GGATACAGAT**)= 0.214
         f(**GGATACAGAT**)= Ø

sketch       A   G    T

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

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

sequence      CAGTATG**GATACAGATG**GAGATATCATCGAGTAGGGGCACTGTACCAGAG

       hash(**GATACAGATG**)= 0.678
         f(**GATACAGATG**)= ∄

sketch        A   G    T

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

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

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

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

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

$f(10-mer) \rightarrow \emptyset \quad if \quad hash(10-mer) > 0.2$

sequence      CAGTATGG**ATACAGATGG**AGATATCATCGAGTAGGGGCACTGTACCAGAG

hash(**ATACAGATGG**)= 0.669

f(**ATACAGATGG**)= ∅

sketch      A   G    T

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching

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

$f(10-mer) \to A \quad if \quad hash(10-mer) \in [0, 0.05]$

$f(10-mer) \to C \quad if \quad hash(10-mer) \in [0.05, 0.1]$

$f(10-mer) \to G \quad if \quad hash(10-mer) \in [0.1, 0.15]$

$f(10-mer) \to T \quad if \quad hash(10-mer) \in [0.15, 0.2]$

$f(10-mer) \to \emptyset \quad if \quad hash(10-mer) > 0.2$

sequence        CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCAC**TGTACCAGAG**

hash(**TGTACCAGAG**)= 0.06

f(**TGTACCAGAG**)= C

sketch            A  G   T          T  C        C      G     T      C

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching

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

**order (l)**

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

**compression ratio (c)**

sequence    CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch         A  G   T            T  C        C      G      T      C

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR=**Mapping-friendly** Sequence Reductions

▶ If two reads align, their sketchs align too

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR=**Mapping-friendly** Sequence Reductions

▶ If two reads align, their sketchs align too

Introduction
Distiguishing haplotypes with noisy reads – HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## Assembling using MSR sketches

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# Very fast assembly: the Alice assembler

AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG
GAGATATCATCGAGTAGGGGCACTGTACCAGAGCCGG
GATATCATCGAGTAGGGGCACTGTACCAGAGCCGGTTATAC

**MSR sketching** ↓

AGTTCCGT          TCCGTCAA          CGTCAATG

**Assembly** ↓

AGTTCCGT
   TCCGTCAA
      CGTCAATG
AGTTCCGTCAATG

**Inflating** ↓

AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAGCCGGTTATAC

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching keeps SNPs

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

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

sequence1        CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sequence2        CAGTATGGATACAGATGGAGATAT**G**ATCGAGTAGGGGCACTGTACCAGAG

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching keeps SNPs

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

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

sequence1       **CAGTATGGAT**ACAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch1         A

sequence2       **CAGTATGGAT**ACAGATGGAGATAT**G**ATCGAGTAGGGGCACTGTACCAGAG

sketch2         A

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching keeps SNPs

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

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

sequence1      C**AGTATGGATA**CAGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch1         A

sequence2      C**AGTATGGATA**CAGATGGAGATAT**G**ATCGAGTAGGGGCACTGTACCAGAG

sketch2         A

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching keeps SNPs

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

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

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

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

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

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

sequence1       CA**GTATGGATAC**AGATGGAGATATCATCGAGTAGGGGCACTGTACCAGAG

sketch1             A  G

sequence2       CA**GTATGGATAC**AGATGGAGATAT**G**ATCGAGTAGGGGCACTGTACCAGAG

sketch2             A  G

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching keeps SNPs

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

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

sequence1    CAGTATGGATACAG**ATGGAGATAT**CATCGAGTAGGGGCACTGTACCAGAG

sketch1        A  G    T

sequence2    CAGTATGGATACAG**ATGGAGATATG**ATCGAGTAGGGGCACTGTACCAGAG

sketch2        A  G    T

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching keeps SNPs

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

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

| sequence1 | CAGTATGGATACAGA**TGGAGATATC**ATCGAGTAGGGGCACTGTACCAGAG |
|---|---|
| sketch1 | A G  T        T |
| sequence2 | CAGTATGGATACAGA**TGGAGATATG**ATCGAGTAGGGGCACTGTACCAGAG |
| sketch2 | A G  T |

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching keeps SNPs

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

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

sequence1    CAGTATGGATACAGAT**GGAGATATCA**TCGAGTAGGGGCACTGTACCAGAG

sketch1          A  G    T              T

sequence2    CAGTATGGATACAGAT**GGAGATATGA**TCGAGTAGGGGCACTGTACCAGAG

sketch2          A  G    T                   G

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching keeps SNPs

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

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

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

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

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

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

| sequence1 | CAGTATGGATACAGATG**GAGATATCAT**CGAGTAGGGGCACTGTACCAGAG |
|---|---|
| sketch1 | A G T   T C |
| sequence2 | CAGTATGGATACAGATG**GAGATATGAT**CGAGTAGGGGCACTGTACCAGAG |
| sketch2 | A G T   G |

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching keeps SNPs

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

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

| | |
|---|---|
| sequence1 | CAGTATGGATACAGATGG**AGATATCATC**GAGTAGGGGCACTGTACCAGAG |
| sketch1 | A G  T        T C |
| sequence2 | CAGTATGGATACAGATGG**AGATATGATC**GAGTAGGGGCACTGTACCAGAG |
| sketch2 | A G  T          G |

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## MSR sketching keeps SNPs

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

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

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

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

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

$f(10-mer) \rightarrow \emptyset \quad if \quad hash(10-mer) > 0.2$

| | |
|---|---|
| sequence1 | CAGTATGGATACAGATGGAGATATCATCGAGTAGGGGCAC**TGTACCAGAG** |
| sketch1 |     A G  T         T C      C      G     T     C |
| sequence2 | CAGTATGGATACAGATGGAGATAT**G**ATCGAGTAGGGGCAC**TGTACCAGAG** |
| sketch2 |     A G  T         G     A C     G     T     C |

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching keeps and amplify SNPs

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

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

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## 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 | 99.965 |
| Escherichia_coli_B766 | 95.641 |
| Escherichia_coli_JM109 | 96.334 |
| Escherichia_coli_b2207 | 95.495 |

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

## 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[1] | $\geq$ 60 days | $\geq$ 60 days | 5h00 |

N.B. only assemblers that distinguish strains are shown

---

[1]Highly accurate metagenome-assembled genomes from human gut microbiota using long-read assembly, binning, and consolidation methods - BiorXiv

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# The dark side: MSR sketching keeps errors

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

$f(10-mer) \to A \quad if \quad hash(10-mer) \in [0, 0.05]$

$f(10-mer) \to C \quad if \quad hash(10-mer) \in [0.05, 0.1]$

$f(10-mer) \to G \quad if \quad hash(10-mer) \in [0.1, 0.15]$

$f(10-mer) \to T \quad if \quad hash(10-mer) \in [0.15, 0.2]$

$f(10-mer) \to \emptyset \quad if \quad hash(10-mer) > 0.2$

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching: conclusion & perspectives

▶ mapping-friendly and keeps SNPs: perfectly adapted to
  haplotype assembly with high-fidelity reads

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching: conclusion & perspectives

▶ mapping-friendly and keeps SNPs: perfectly adapted to haplotype assembly with high-fidelity reads

▶ Still a lot to explore on MSR sketching: changing the function, changing the use case...

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
**Distinguishing haplotypes with high-fidelity reads - Alice**
Conclusion

# MSR sketching: conclusion & perspectives
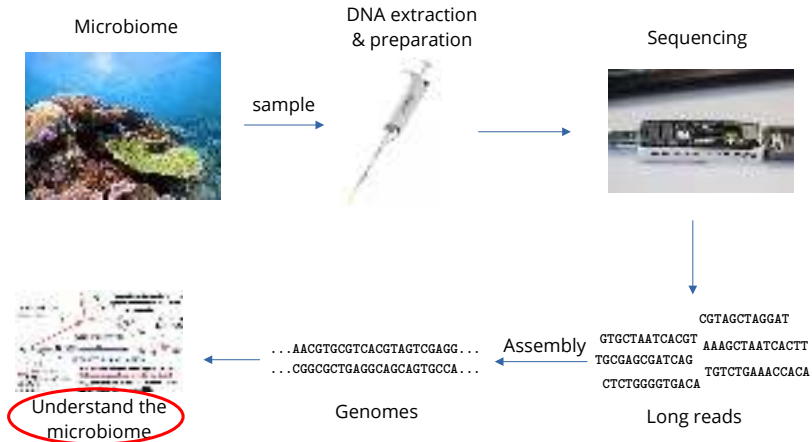
▶ mapping-friendly and keeps SNPs: perfectly adapted to haplotype assembly with high-fidelity reads

▶ Still a lot to explore on MSR sketching: changing the function, changing the use case...

▶ Tune to what extent we want to keep variation

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion

# Conclusion

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
**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

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
**Conclusion**

# Why is this thesis useful?



Microbiome

DNA extraction
& preparation

Sequencing

sample

CGTAGCTAGGAT

...AACGTGCGTCACGTAGTCGAGG...    GTGCTAATCACGT  AAAGCTAATCACTT
Assembly
...CGGCGCTGAGGCAGCAGTGCCA...    TGCGAGCGATCAG
TGTCTGAAACCACA
CTCTGGGGTGACA

Understand the
microbiome

Genomes

Long reads

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
**Conclusion**

## Why is this thesis useful?

Roland Faure        Public defence

Introduction
Distiguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
Conclusion
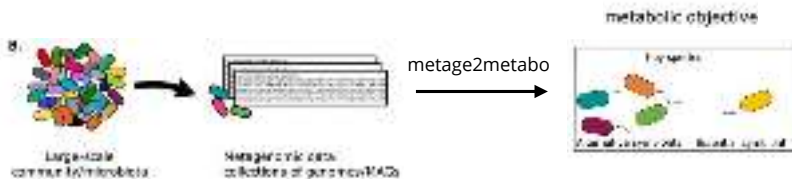
## Example of an application: metage2metabo[1]

Arnaud Belcour

---

[1]Metage2Metabo, microbiota-scale metabolic complementarity for the identification of key species - Belcour et al., 2020

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
**Conclusion**

## Example of an application: metage2metabo[1]



Arnaud Belcour



metabolic objective

metage2metabo

Large-scale
community/microbiota

Metagenomic data
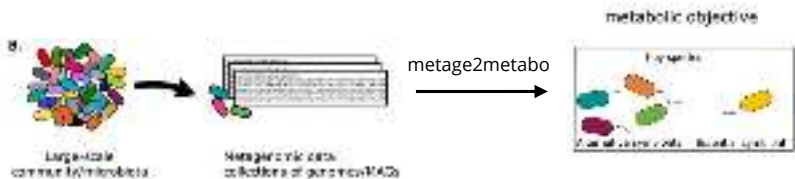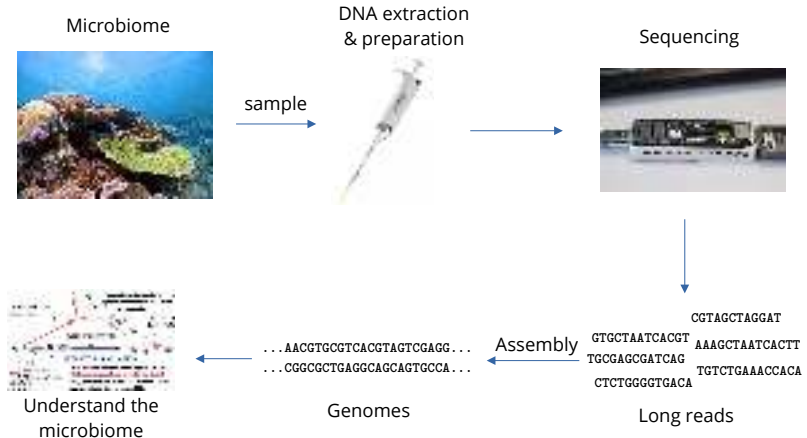collections of genomes-MAGs

▶ Predictions for human health, soil fertility, ecology...

---

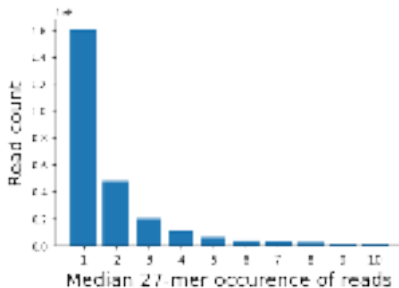[1]Metage2Metabo, microbiota-scale metabolic complementarity for the
identification of key species - Belcour et al., 2020

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
**Conclusion**

# What is the future of assembly?



Microbiome

DNA extraction
& preparation

Sequencing

sample

...AACGTGCGTCACGTAGTCGAGG...
...CGGCGCTGAGGCAGCAGTGCCA...

Assembly

CGTAGCTAGGAT
GTGCTAATCACGT   AAAGCTAATCACTT
TGCCGAGCGATCAG   TGTCTGAAACCACA
CTCTGGGGTGACA

Understand the
microbiome

Genomes

Long reads

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
**Conclusion**

## All DNA is not captured by sequencing

▶ Example of the sequencing of the soil microbiote

Nicolas Maurice



Adapted from the work of Nicolas Maurice

▶ Low-coverage assembly

▶ Missing DNA

Roland Faure          Public defence

Introduction
Distinguishing haplotypes with noisy reads - HairSplitter
Distinguishing haplotypes with high-fidelity reads - Alice
**Conclusion**

# What is the future of assembly?