

# Solving the haplotyping problem with linear programming

Tam K.M. Truong<sup>1</sup>, Roland Faure<sup>1,2</sup>, Rumen Andonov<sup>1</sup>

<sup>1</sup>Université de Rennes, IRISA, France

<sup>2</sup>Université libre de Bruxelles (ULB), Belgium

February 2023

# The ultimate goal: understand ecosystems



# The ultimate goal: understand ecosystems

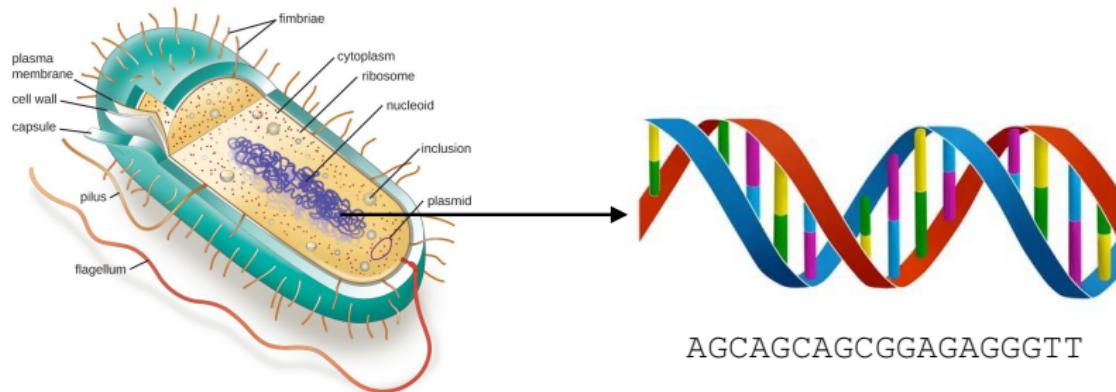


# The ultimate goal: understand ecosystems



## ► Genome sequencing

# Biology reminder



- ▶ DNA is the instruction manual of the bacteria
- ▶ Typically 5Mb long

# Studying an ecosystem through sequencing

Ecosystem



# Studying an ecosystem through sequencing

Ecosystem



DNA extraction



# Studying an ecosystem through sequencing

Ecosystem



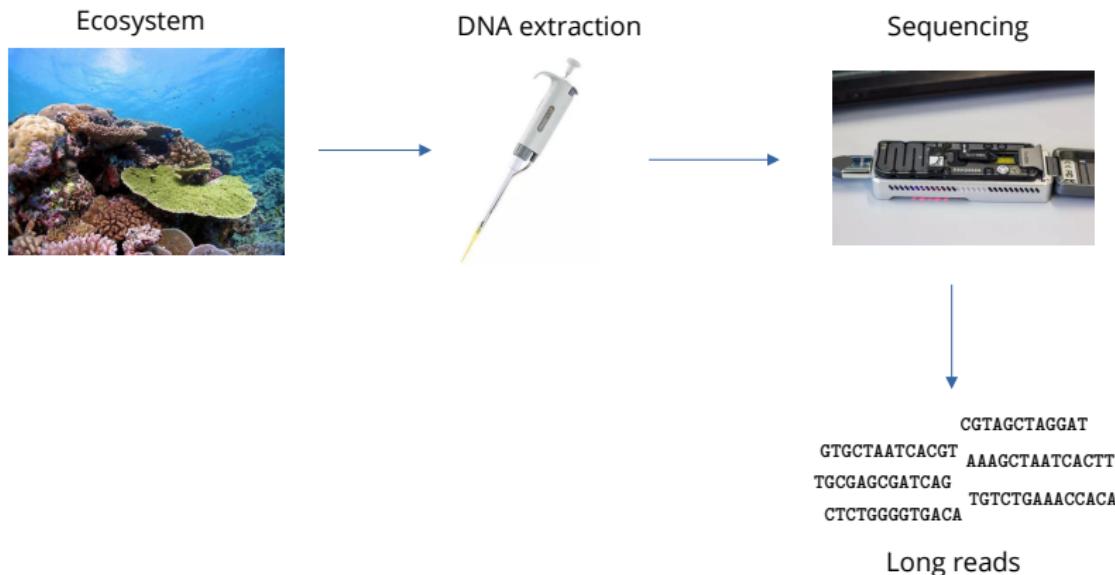
DNA extraction



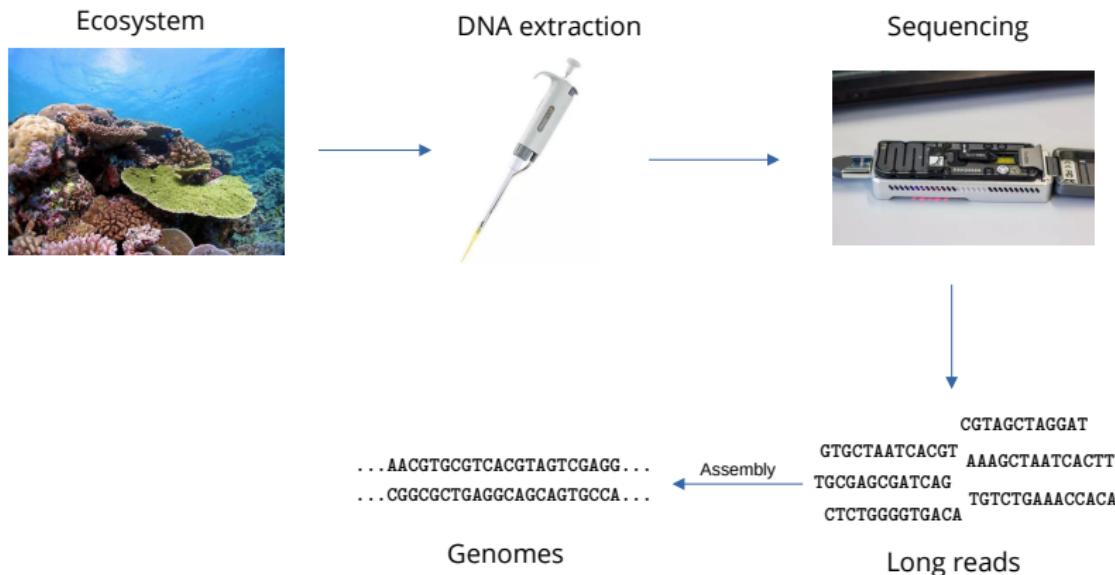
Sequencing



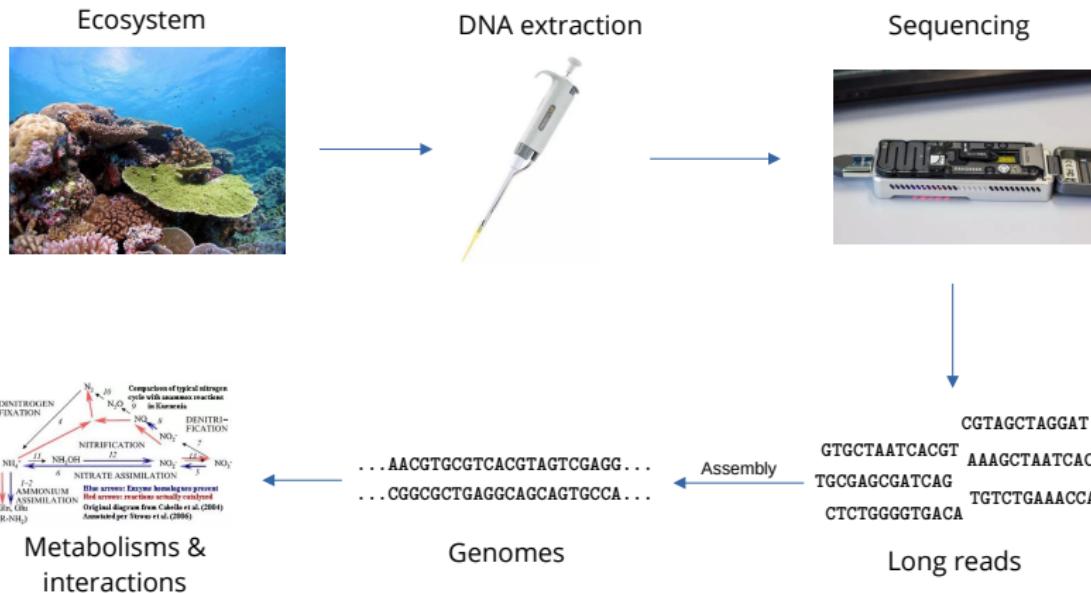
# Studying an ecosystem through sequencing



# Studying an ecosystem through sequencing

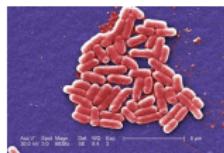


# Studying an ecosystem through sequencing



# Strains are important

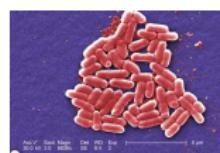
*Escherichia coli K12*



...ACGCTGAGGCAGCATGTGCCA...



*Escherichia coli Sakai*

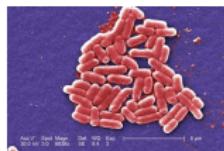


...GCGCTGAGGCAGCATGTGCGA...



# Strains are important

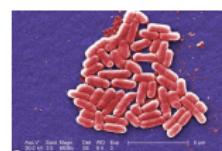
*Escherichia coli K12*



...ACGCTGAGGCAGCATGTGCCA...



*Escherichia coli Sakai*

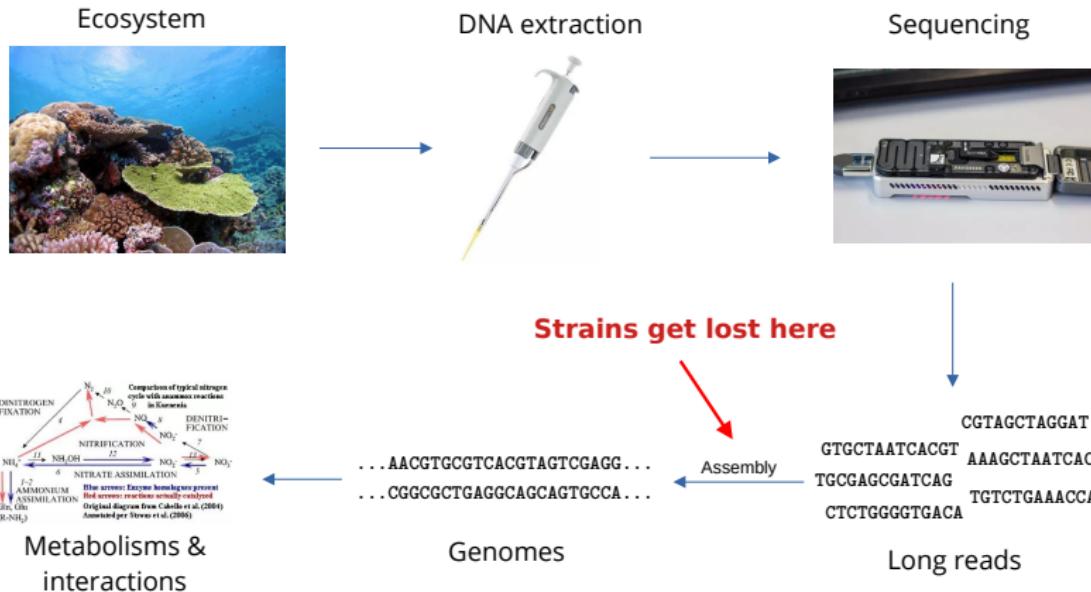


...GCGCTGAGGCAGCATGTGCGA...

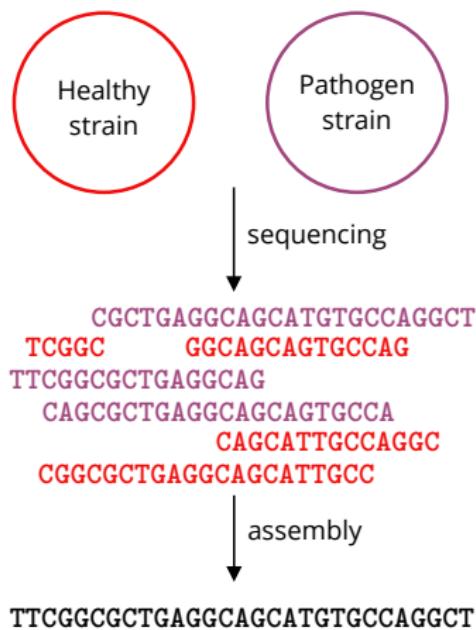


- ▶ “Knowledge gaps remain: strain diversity at a microscale and between species” R. Caldwell *et al.*, doi.org/10.1016/j.mib.2022.102222, 2022

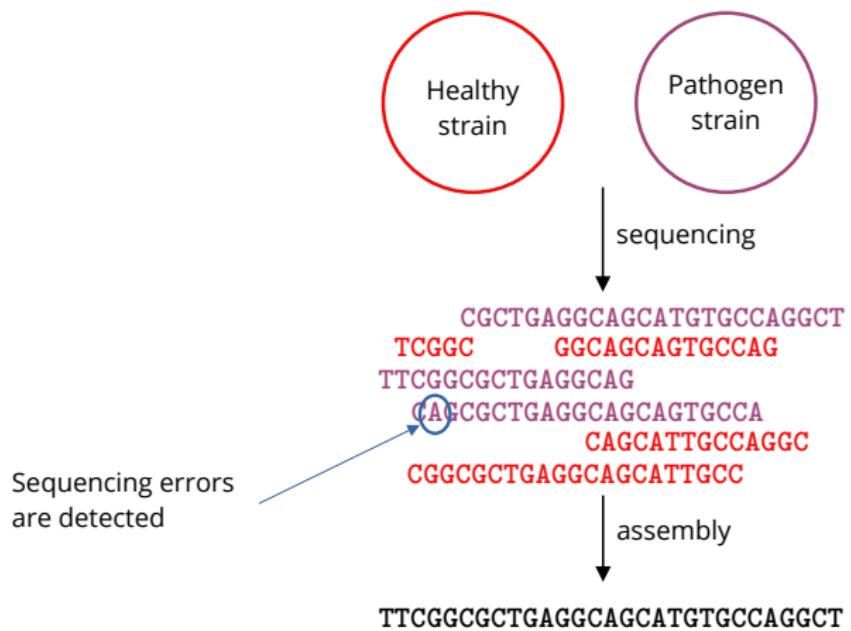
# Problème: différencier les souches



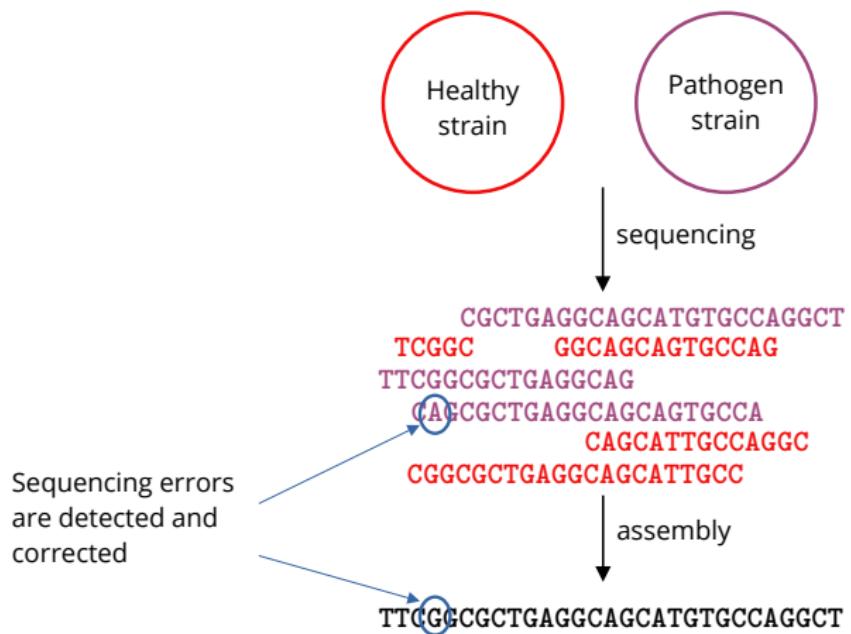
# Problem: differentiating strains



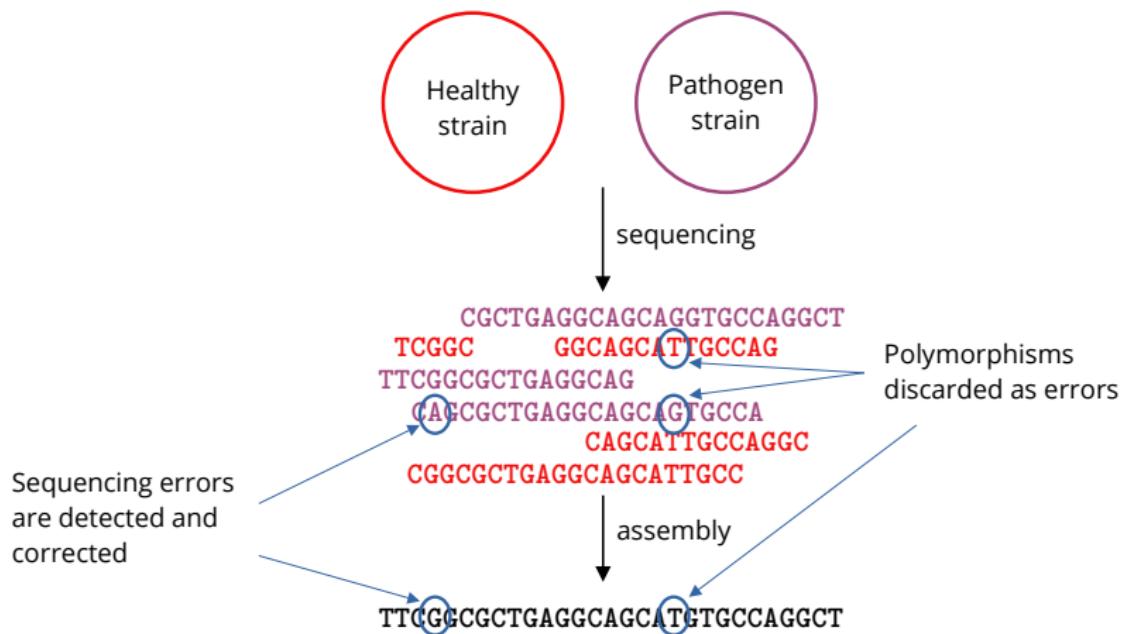
# Problème: différencier les souches



# Problème: différencier les souches



# Problem: differentiating strains



# Goal

```
AACTGTGTCCCT-TAGAGCGATT CGCGAGCGTA  
AACGGTGTCCCTATGGAGCG--TCGCGACCGTA  
AACTGTGTCCCTATAGAGCGATACCGGACCGTA  
AACTGTGTCCCT-TAGAGCGATT CGCGAGCGTA  
AACGGTGTCCCTATAGAGCGATT CGCGACCGTA  
AACGGTGTCCCTATAGAGCGATT CGCGACCGTA  
AACTGTGACCCCTATAGAGCGATACCGGACCGTA  
AACTGTGTCCCT-TAGAGCGATT CGC AAGCGTA  
AACTCGGTCCCTATAGAGCGATACCGGACCGTA
```

**Input:** All reads and the draft assembly  
**Output:** Reads split into groups

```
AACTGTGTCCCT-TAGAGCGATT CGCGAGCGTA  
AACTGTGTCCCT-TAGAGCGATT CGCGAGCGTA  
AACTGTGTCCCT-TAGAGCGATT CGCA AAGCGTA
```

```
AACGGTGTCCCTATGGAGCG--TCGCGACCGTA  
AACGGTGTCCCTATAGAGCGATT CGCGACCGTA  
AACGGTGTCCCTATAGAGCGATT CGCGACCGTA
```

```
AACTGCGTCCCTATAGAGCGATACCGGACCGTA  
AACTGTGACCCCTATAGAGCGATACCGGACCGTA  
AACTGTGTCCCTATAGAGCGATACCGGACCGTA
```

# Differentiating errors from signal

- r1 AACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGATAGATAGCA
- r2 AACAAGATAGAC**G**AGATAGACACAG**C**TTGGCGTTAGGAACAGATGATAGATAGCA
- r3 AACAAGATAGACAAGATAGACACAG**C**TTGGCGTTAG**T**AACAGATGACAGATAGCA
- r4 AACAAGAT**C**GAC**G**AGATAGACACAT**T**CTTGGCGTTAGGAACAT**TT**GACAGATAGCA
- r5 AACAAGAT**C**GACAAAGATAG**G**CACAT**A**TTGGCGTTAGGAACAG**T**TGATAGATAGCA
- r6 AACAAGAT**C**GAC**G**AGATAGACACAT**A**TTGGCGTTAGGAT**C**AG**T**GACAGATAGCA

## Differentiating errors from signal

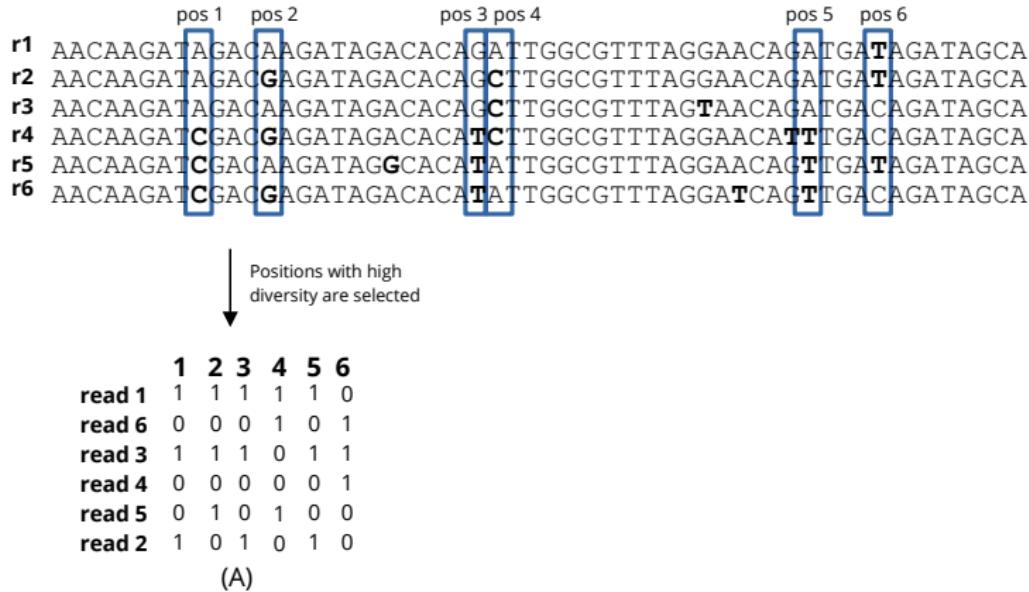
r1 AACAAGATAGACAAGATAGACACAGATTGGCGTTAGGAACAGATGATAGATAGCA  
r2 AACAAGATAGAC**G**GAGATAGACACAG**C**TTGGCGTTAGGAACAGATGATAGATAGCA  
r3 AACAAGATAGACAAGATAGACACAG**C**TTGGCGTTAG**T**AACAGATGACAGATAGCA  
r4 AACAAGAT**C**GAC**G**GAGATAGACACAT**T**CTTGGCGTTAGGAACAT**T**TGACAGATAGCA  
r5 AACAAGAT**C**GACACAAGATAG**G**CACAT**T**ATTGGCGTTAGGAACAG**T**TGATAGATAGCA  
r6 AACAAGAT**C**GAC**G**GAGATAGACACAT**T**ATTGGCGTTAGGAT**C**AG**T**TGACAGATAGCA

# The model

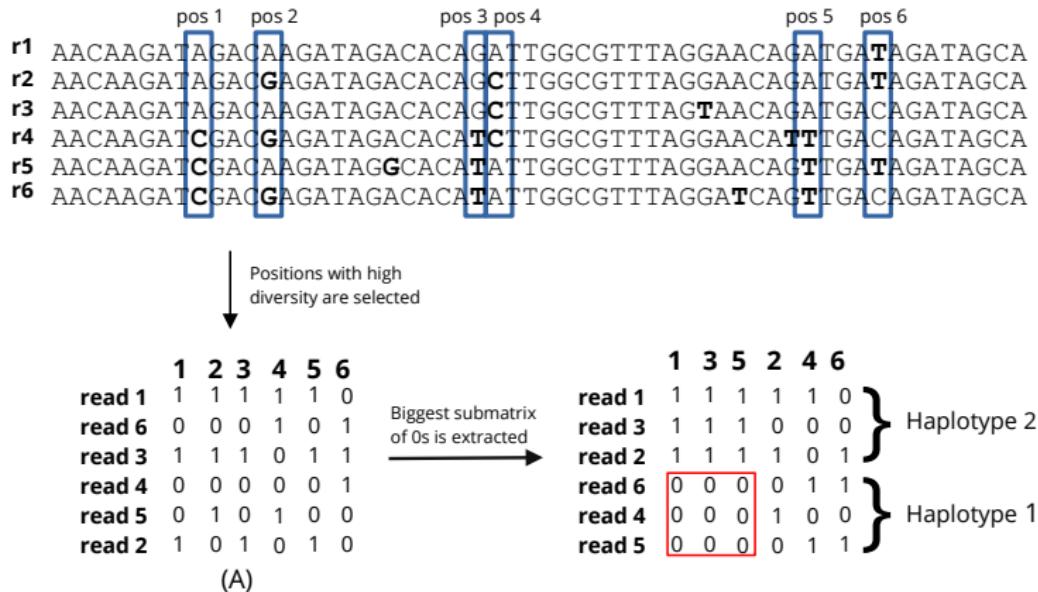
	pos 1	pos 2	pos 3	pos 4	pos 5	pos 6	
r1	AACAAGATAGACAAGATAGACACAGATGGCGTTAGGAACAGATGA			T	AGATAGCA		
r2	AACAAGATAGACGAGATAGACACAG	C	T	GGCGTTAGGAACAGATGA	T	AGATAGCA	
r3	AACAAGATAGACAAGATAGACACAG	C	T	GGCGTTAGTAACAGATGACAGATAGCA			
r4	AACAAGATC GACG GAGATAGACACAT	T	C	T	GGCGTTAGGAACAC	TT	IGACAGATAGCA
r5	AACAAGATC GACAAGATAGGCACAT	T	AT	GGCGTTAGGAACAGT	GAT	T	AGATAGCA
r6	AACAAGATC GACG GAGATAGACACAT	T	AT	GGCGTTAGGA	T	CAGT	GACAGATAGCA

Positions with high diversity are selected

# The model



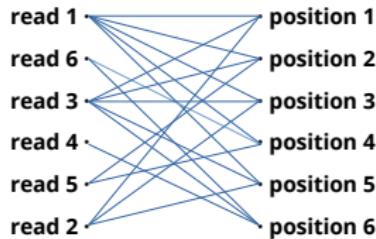
# The model



# From the matrix to the graph

	1	2	3	4	5	6
read 1	1	1	1	1	1	0
read 6	0	0	0	1	0	1
read 3	1	1	1	0	1	1
read 4	0	0	0	0	0	1
read 5	0	1	0	1	0	0
read 2	1	0	1	0	1	0

(A)



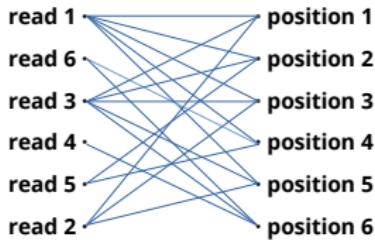
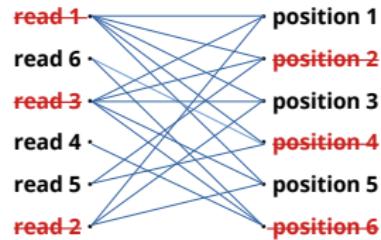
## From the matrix to the graph

	1	2	3	4	5	6
read 1	1	1	1	1	1	0
read 6	0	0	0	1	0	1
read 3	1	1	1	0	1	1
read 4	0	0	0	0	0	1
read 5	0	1	0	1	0	0
read 2	1	0	1	0	1	0

(A)

Bigest submatrix  
of 0s is extracted

	1	3	5	2	4	6
read 1	1	1	1	1	1	0
read 3	1	1	1	0	0	0
read 2	1	1	1	1	0	1
read 6	0	0	0	0	1	1
read 4	0	0	0	1	0	0
read 5	0	0	0	0	1	1

Minimum cover  
vertex problem

# Formulation of the constraints

$$\begin{aligned} \min \quad & \sum_{i \in U \cup V} \deg(i)v_i \\ v_i + v_j \geq 1, \quad & \forall (i, j) \in E \\ v_i \in \{0, 1\}, \quad & \forall i \in U \cup V \end{aligned}$$

Weighted minimal vertex cover problem

$$\begin{aligned} \min \quad & \sum_{i \in U \cup V} \deg(i)v_i \\ 1 - v_i \geq c_{ij}, \quad & \forall i \in U, \forall j \in V \\ 1 - v_j \geq c_{ij}, \quad & \forall i \in U, \forall j \in V \\ 1 - v_i - v_j \leq c_{ij}, \quad & \forall i \in U, \forall j \in V \\ \epsilon \times \sum_{i \in U} \sum_{j \in V} (1 - A_{i,j})c_{ij} \geq \sum_{i \in U} \sum_{j \in V} A_{i,j}c_{ij} \\ v_i \in \{0, 1\}, c_{ij} \in \{0, 1\} \quad & \forall i \in U \cup V \end{aligned}$$

Weighted minimal vertex cover problem, tolerating an  $\epsilon$  ratio of 1 in the submatrix of 0s

# Formulation of the constraints

$$\begin{aligned} \min \quad & \sum_{i \in U \cup V} \deg(i)v_i \\ v_i + v_j \geq 1, \quad & \forall (i, j) \in E \\ v_i \in \{0, 1\}, \quad & \forall i \in U \cup V \end{aligned}$$

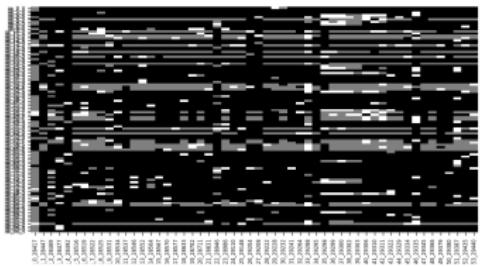
Weighted minimal vertex cover problem

$$\begin{aligned} \min \quad & \sum_{i \in U \cup V} \deg(i)v_i \\ 1 - v_i \geq c_{ij}, \quad & \forall i \in U, \forall j \in V \\ 1 - v_j \geq c_{ij}, \quad & \forall i \in U, \forall j \in V \\ 1 - v_i - v_j \leq c_{ij}, \quad & \forall i \in U, \forall j \in V \\ \epsilon \times \sum_{i \in U} \sum_{j \in V} (1 - A_{i,j})c_{ij} \geq \sum_{i \in U} \sum_{j \in V} A_{i,j}c_{ij} \\ v_i \in \{0, 1\}, c_{ij} \in \{0, 1\} \quad & \forall i \in U \cup V \end{aligned}$$

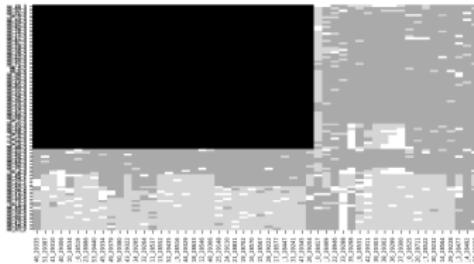
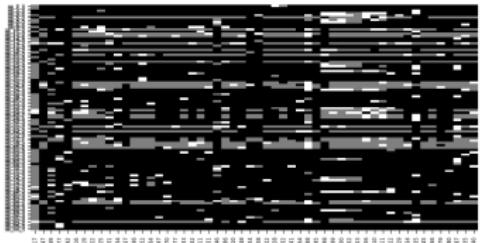
Weighted minimal vertex cover problem, tolerating an  $\epsilon$  ratio of 1 in the submatrix of 0s

- ▶ Run through Gurobi Solver

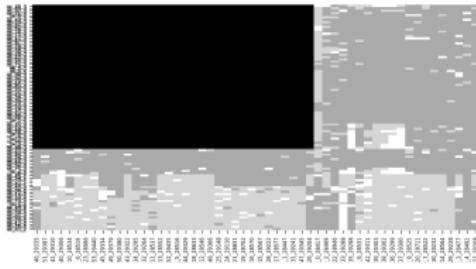
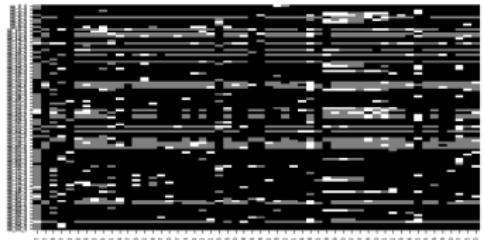
## Result on a mix of 3 *E. coli*



# Result on a mix of 3 *E. coli*

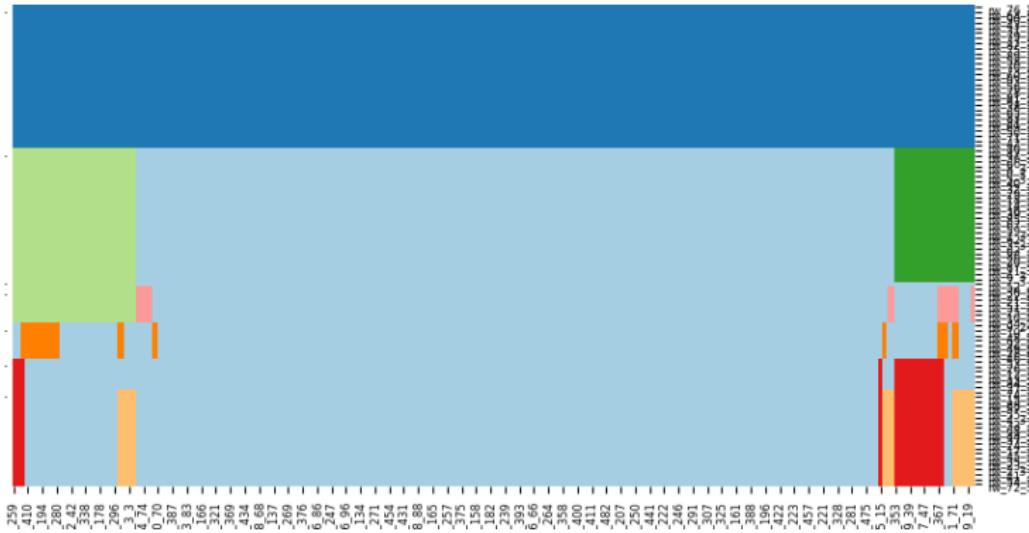


## Result on a mix of 3 *E. coli*



- ▶ Successfully separated two haplotypes from the reference
- ▶ Some false negative but no false positive

# Separating all the haplotypes



- ▶ Iteratively apply the method on remaining matrix
- ▶ The difficult part is knowing when to stop separating

# Conclusion

- ▶ This is the first formulation of the haplotyping problem as a linear programming problem
- ▶ Not yet competitive with existing heuristics

## Acknowledgments

- ▶ Tam K.M. Truong (implementation), Rumen Andonov (supervision) and Dominique Lavenier (original idea)

