

# Haplotype inference and applications

Olivier DELANEAU

Chaire de bioinformatique

Conservatoire des Arts et Métiers

[olivier.delaneau@gmail.com](mailto:olivier.delaneau@gmail.com)

I. Introduction

II. Combinatorial haplotype inference

III. Statistical haplotype inference

IV. Algorithmic tricks

V. Comparison of accuracy

VI. Application 1 : haplotype association tests

VII. Application 2 : genotype imputation

VIII. Application 3 : admixture

## I. Introduction

II. Combinatorial haplotype inference

III. Statistical haplotype inference

IV. Algorithmic tricks

V. Comparison of accuracy

VI. Application 1 : haplotype association tests

VII. Application 2 : genotype imputation

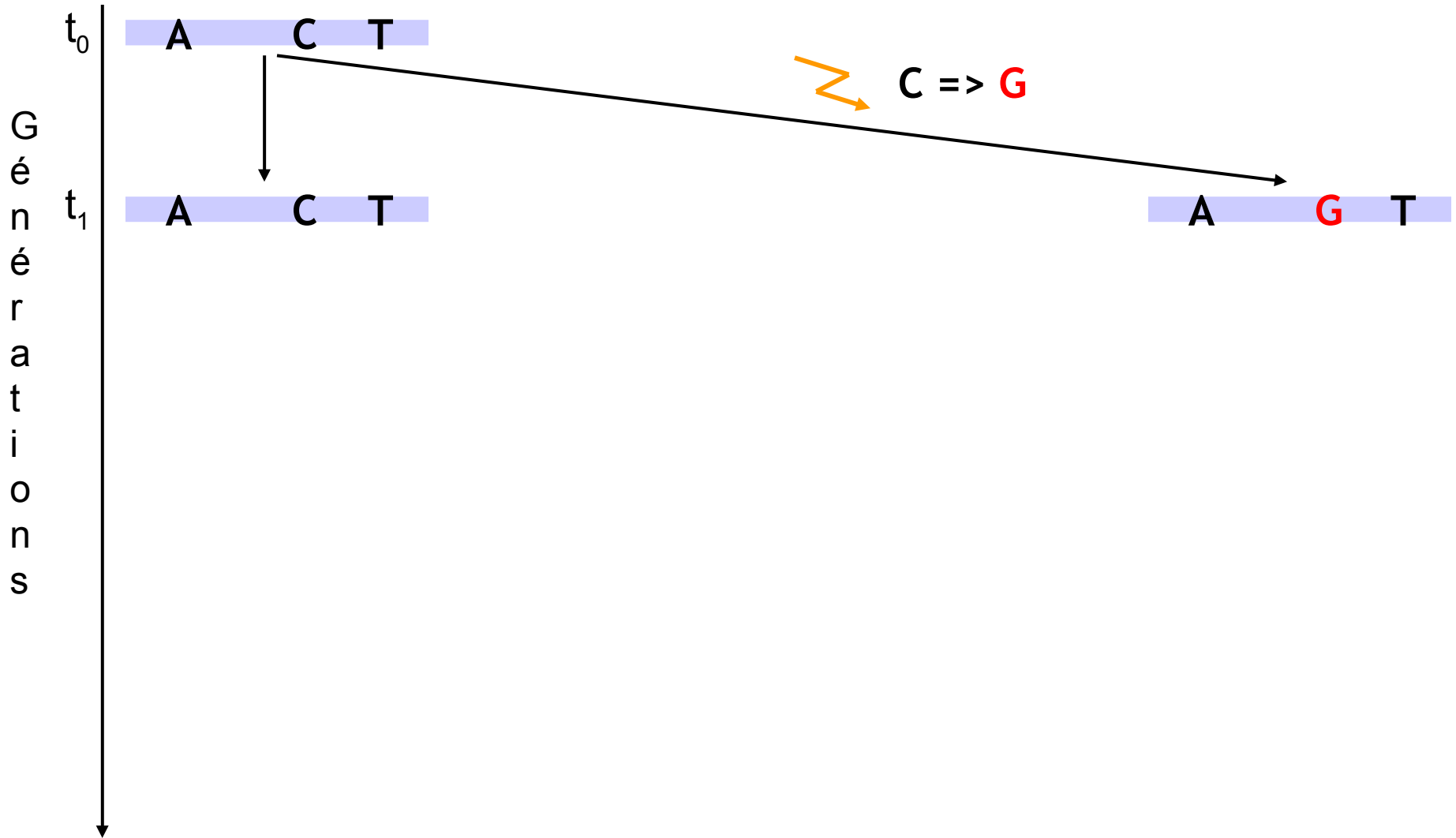
VIII. Application 3 : admixture



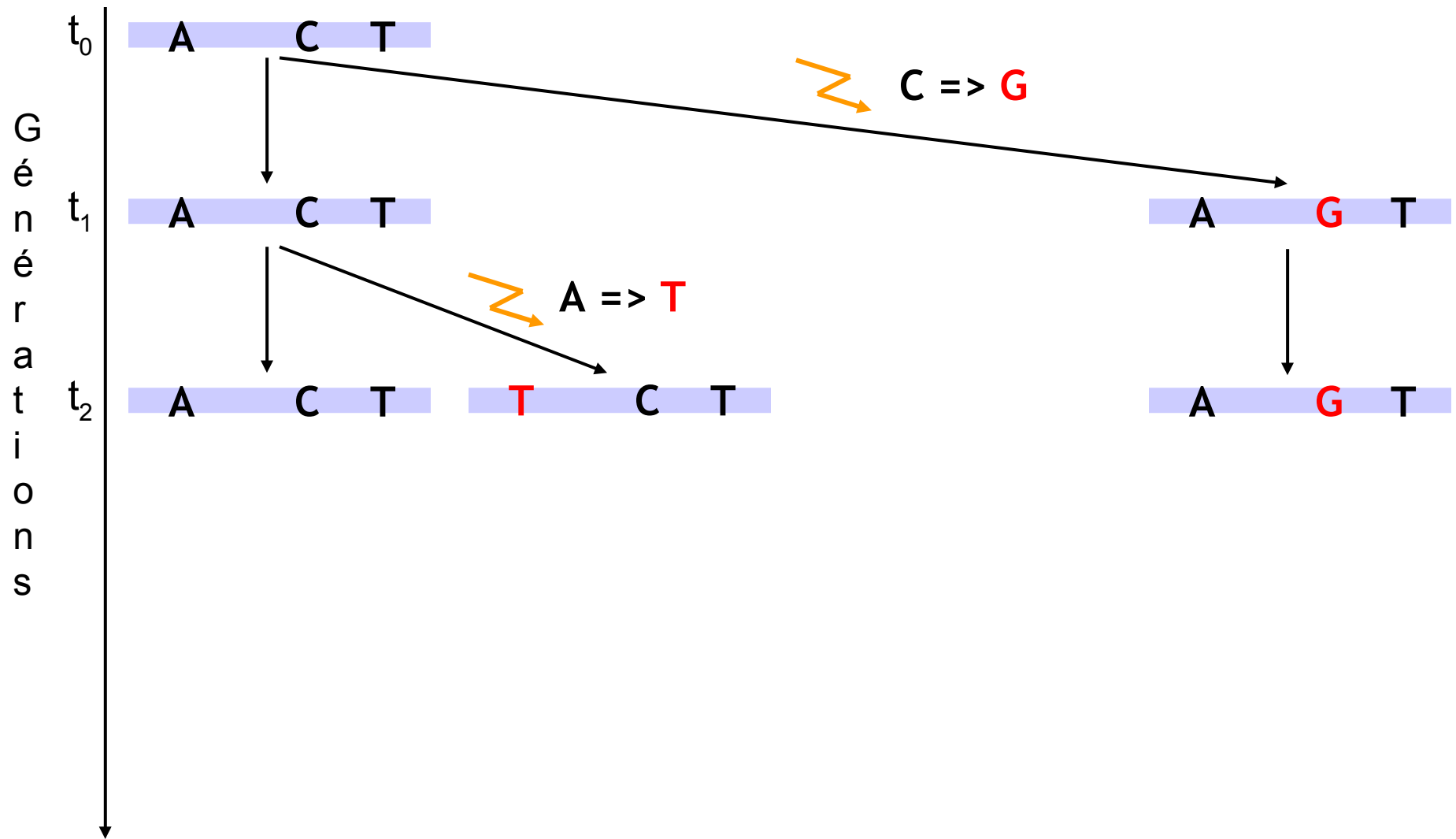
# HAPLOTYPES GENESIS

A C T

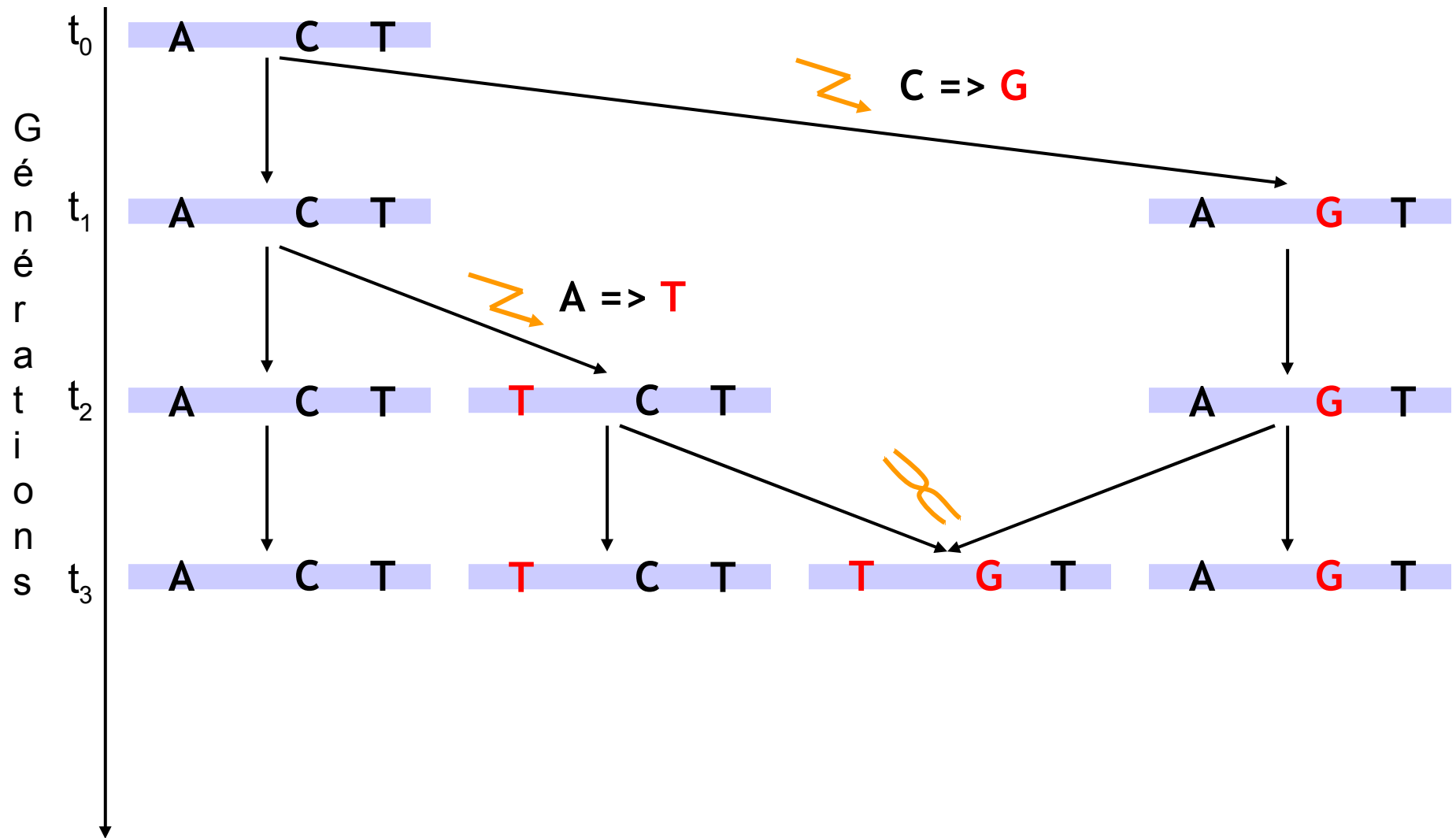
# HAPLOTYPES GENESIS



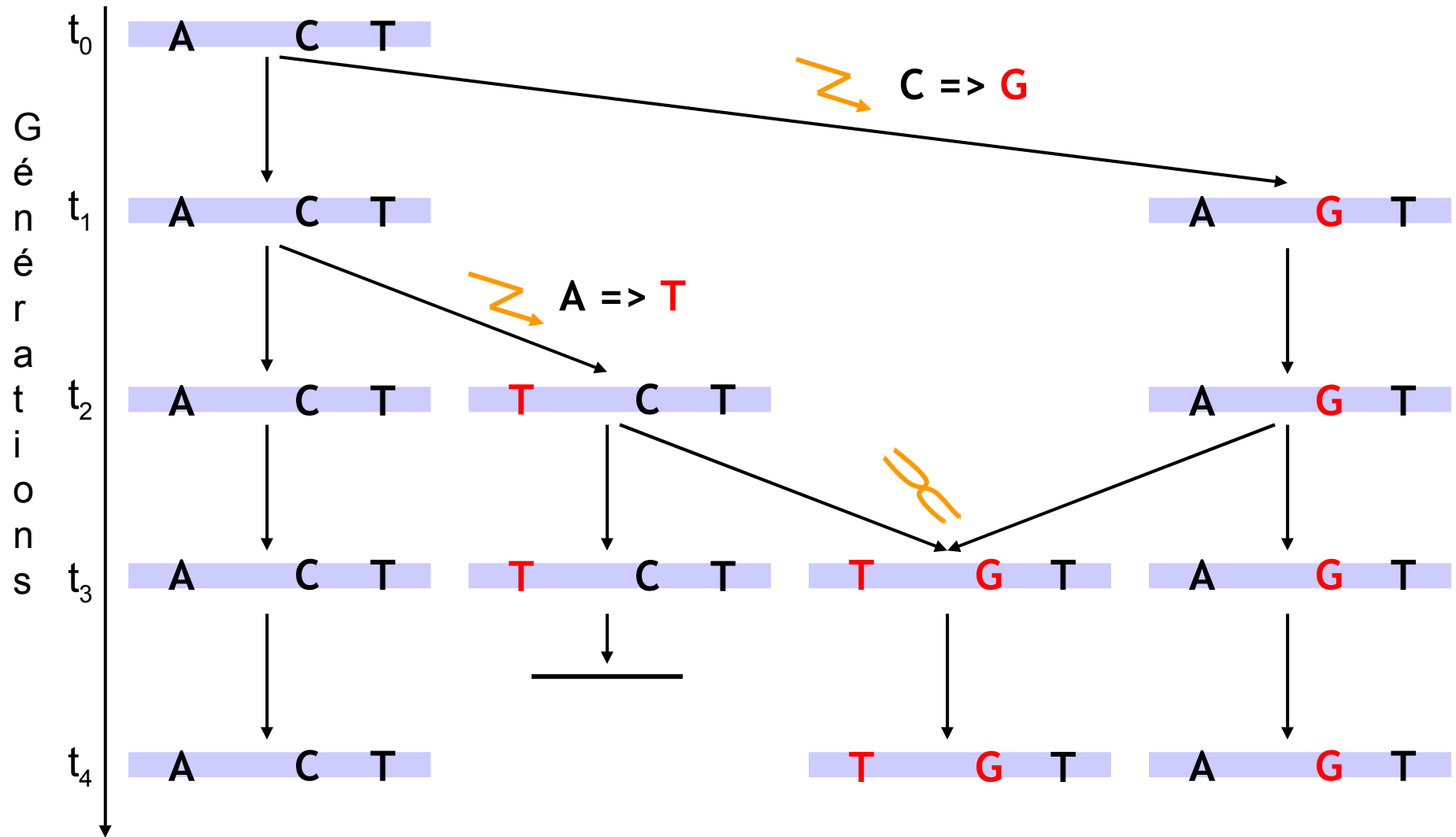
# HAPLOTYPES GENESIS



# HAPLOTYPES GENESIS



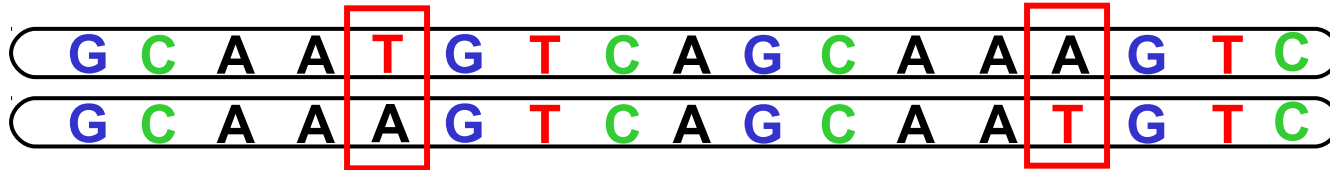
# HAPLOTYPES GENESIS





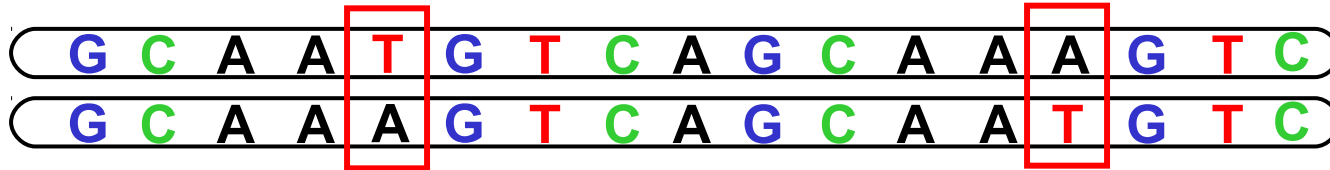
# PROBLEMATIC

Pair of chromosomes of an individual :

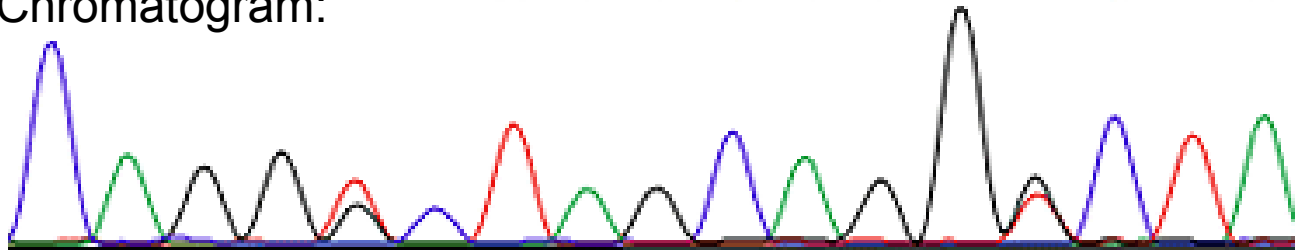


# PROBLEMATIC

Pair of chromosomes of an individual :

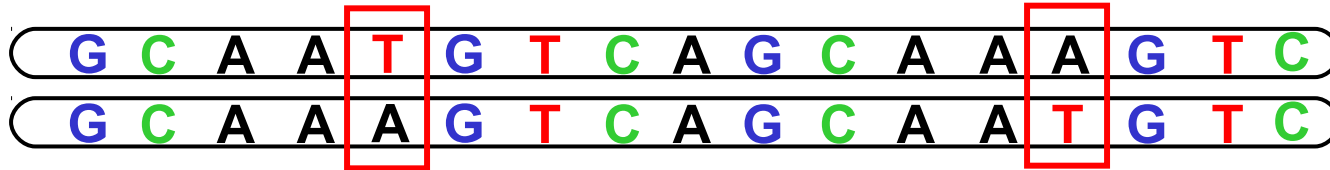


Chromatogram:

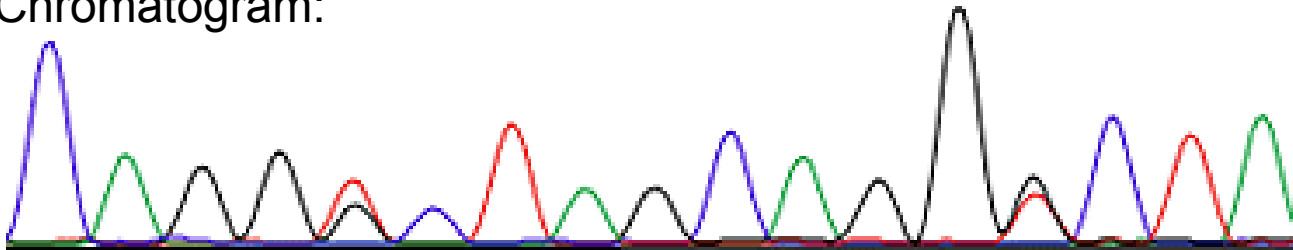


# PROBLEMATIC

Pair of chromosomes of an individual :



Chromatogram:



What is the true haplotype' pair ?

T A

A T

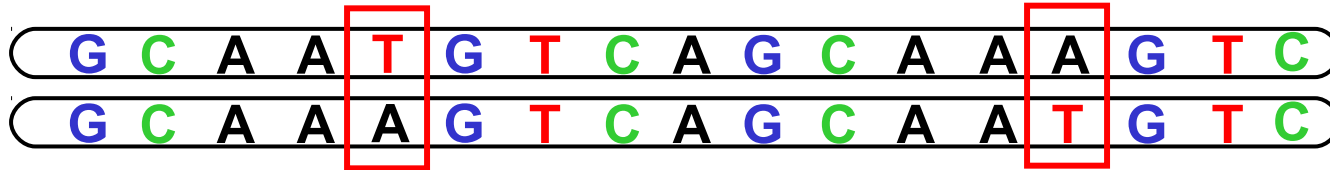
OR

T T

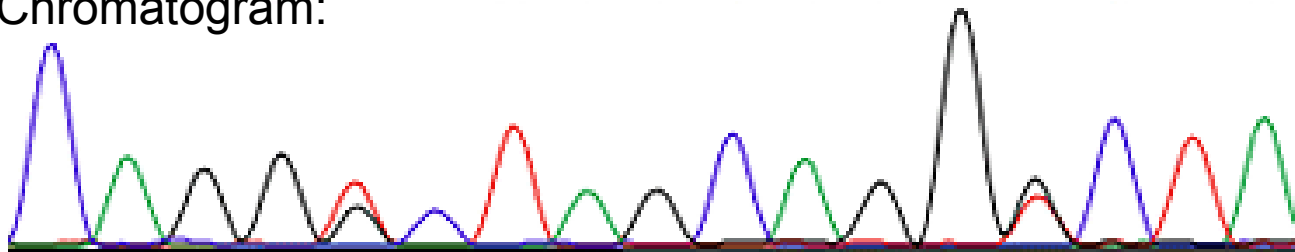
A A

# PROBLEMATIC

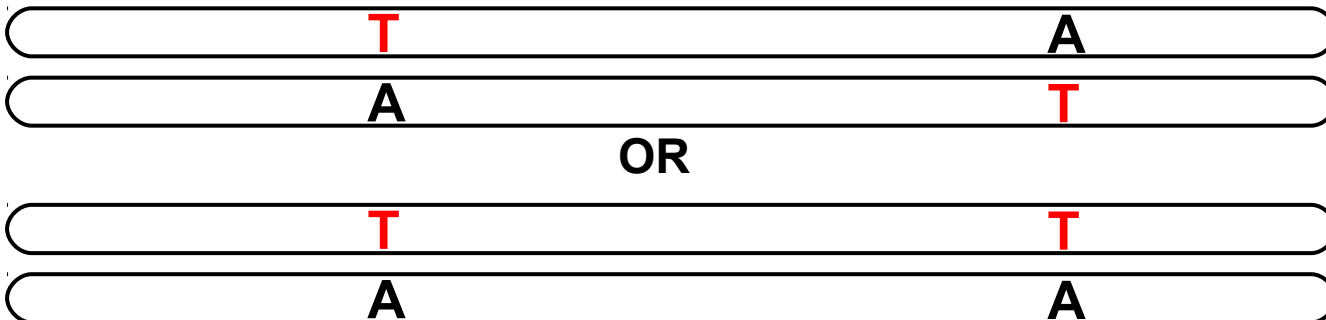
Pair of chromosomes of an individual :



Chromatogram:

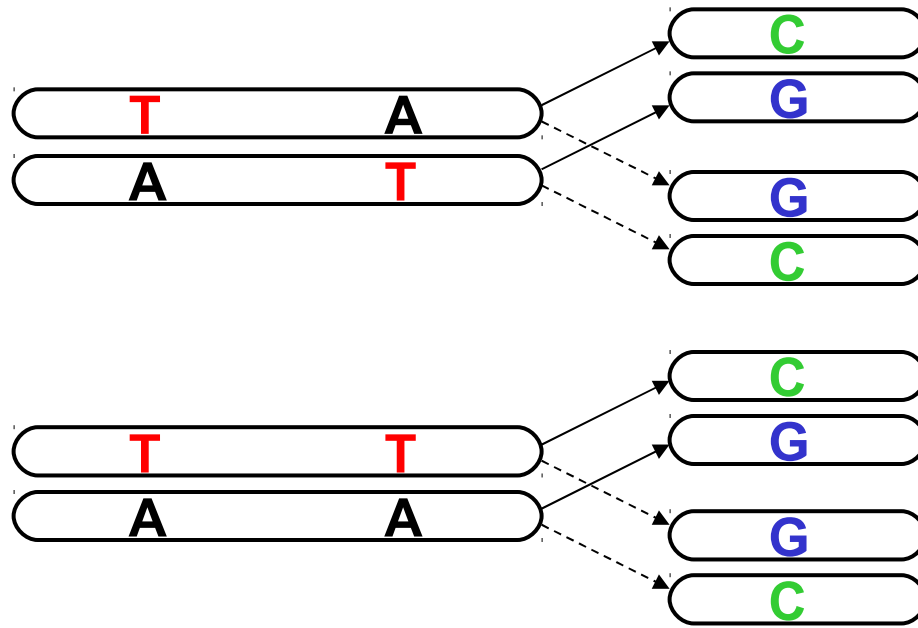


What is the true haplotype' pair ?

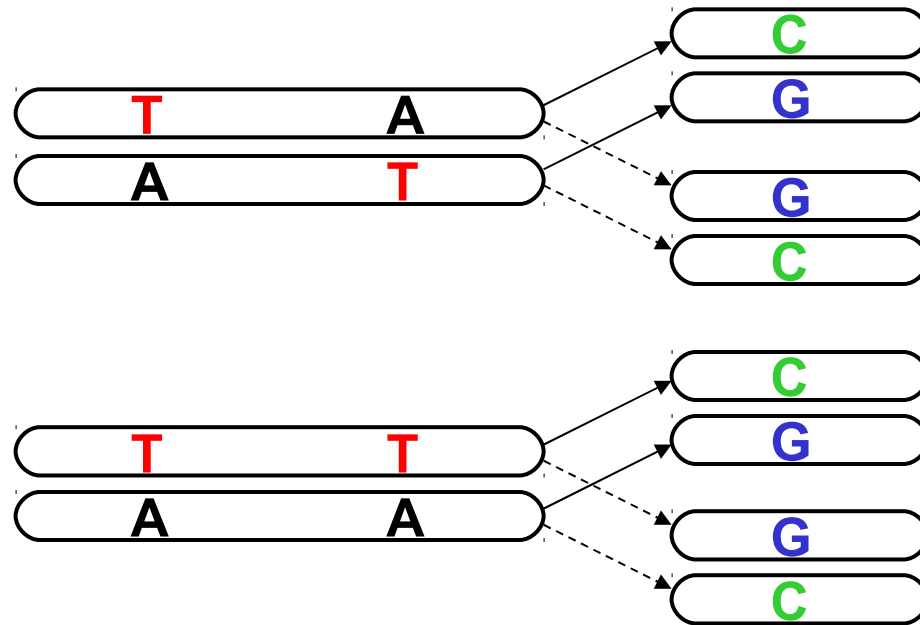


Haplotype inference

# NUMBER OF POSSIBLE HAPLOTYPES



## NUMBER OF POSSIBLE HAPLOTYPES



For a genotype with **S** heterozygous SNPs:

1.  $2^{S-1}$  possible haplotype' pairs.
2.  $2^S$  possible haplotypes.

# EXAMPLE AND NOTATIONS

Sample of genotypes **G**

**L** SNPs

*Genotype#1*

G G G G G T A T G A A A A T

G G G G G T A T G A A A A T

*Genotype#2*

G G T T G T A T G A A A A T

G G G G T T G G G A A A A T

.....

.....

*Genotype#N-1*

G G G G T T G G G G A A A G

G G G G G T A T G A A A A T

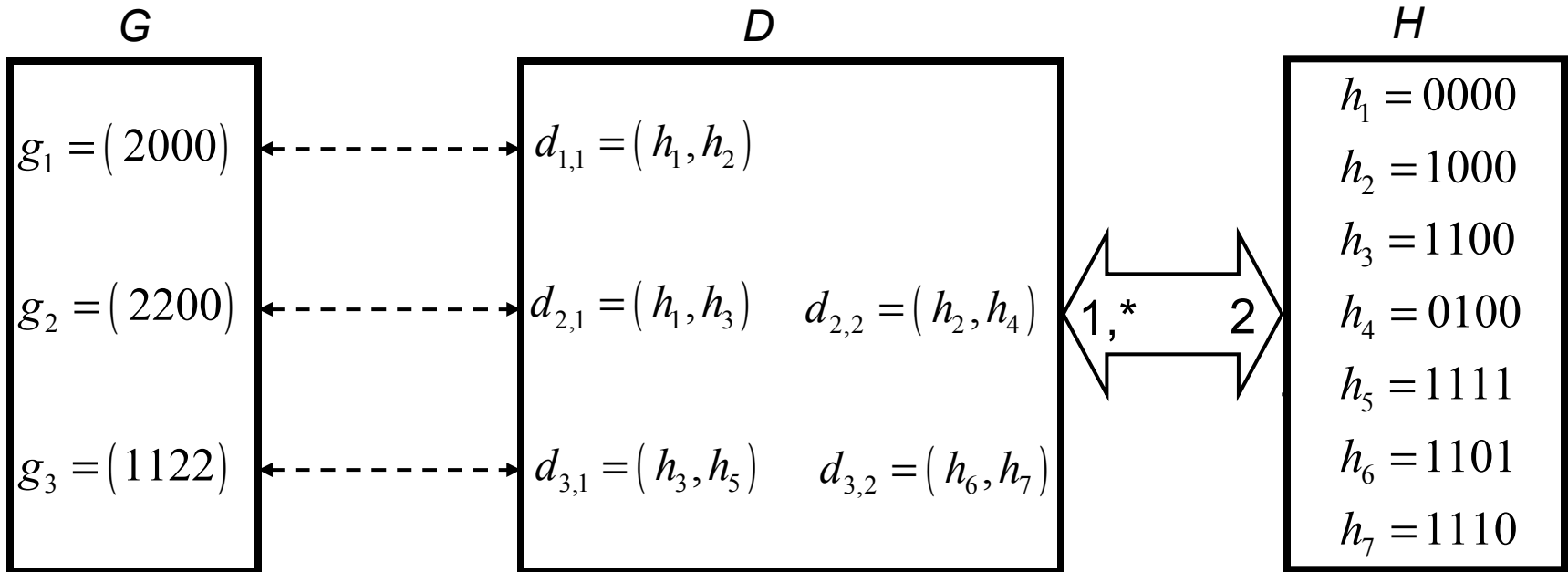
*Genotype#N*

G G G G G T A T G A A A A T

G C G G T T A G G A A A A T

**N** genotypes

# HAPLOTYPE SPACE



## Genotype :

- 0** : homozygous wild allele (0 / 0)
- 1** : homozygous mutant allele (1 / 1)
- 2** : heterozygous (0 / 1)

## Haplotype :

- 0** : wild allele
- 1** : mutant allele



I. Introduction

**II. Combinatorial haplotype inference**

III. Statistical haplotype inference

IV. Algorithmic tricks

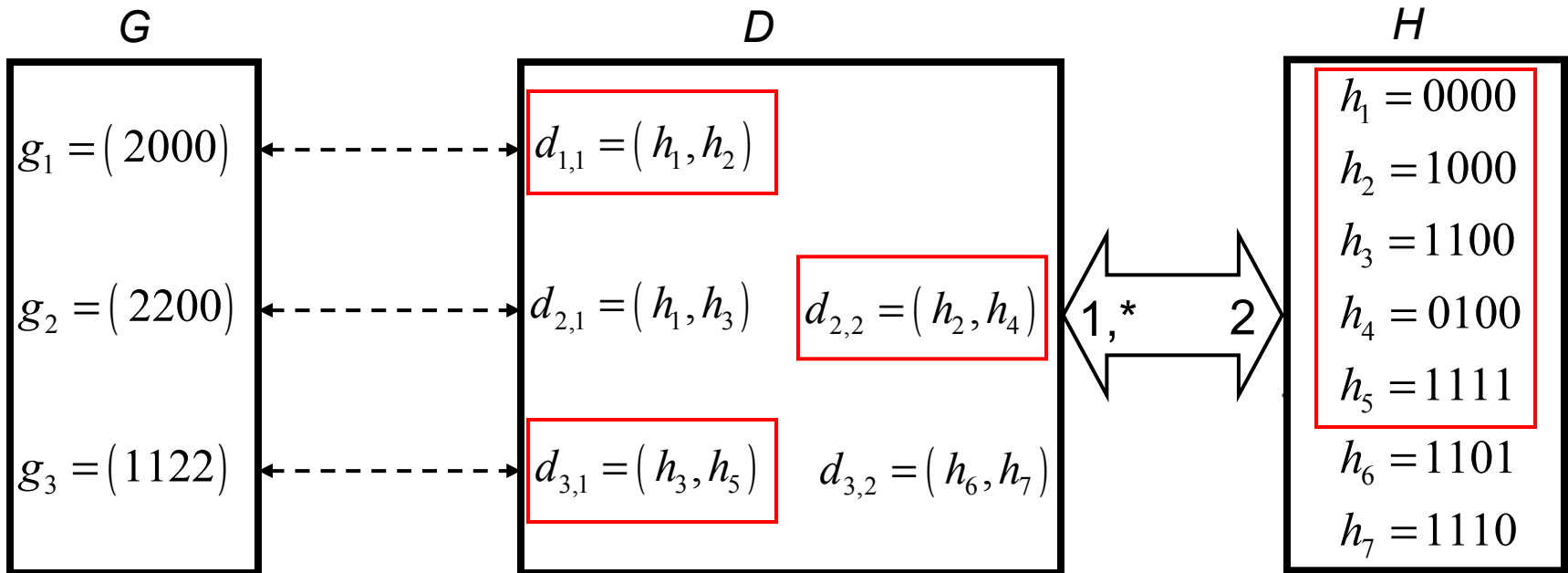
V. Comparison of accuracy

VI. Application 1 : haplotype association tests

VII. Application 2 : genotype imputation

VIII. Application 3 : admixture

# COMBINATORIAL APPROACH



$$D_1 = \{d_{1,1}, d_{2,2}, d_{3,1}\}$$

$$H_1 = \{h_1, h_2, h_3, h_4, h_5\}$$

## Idea :

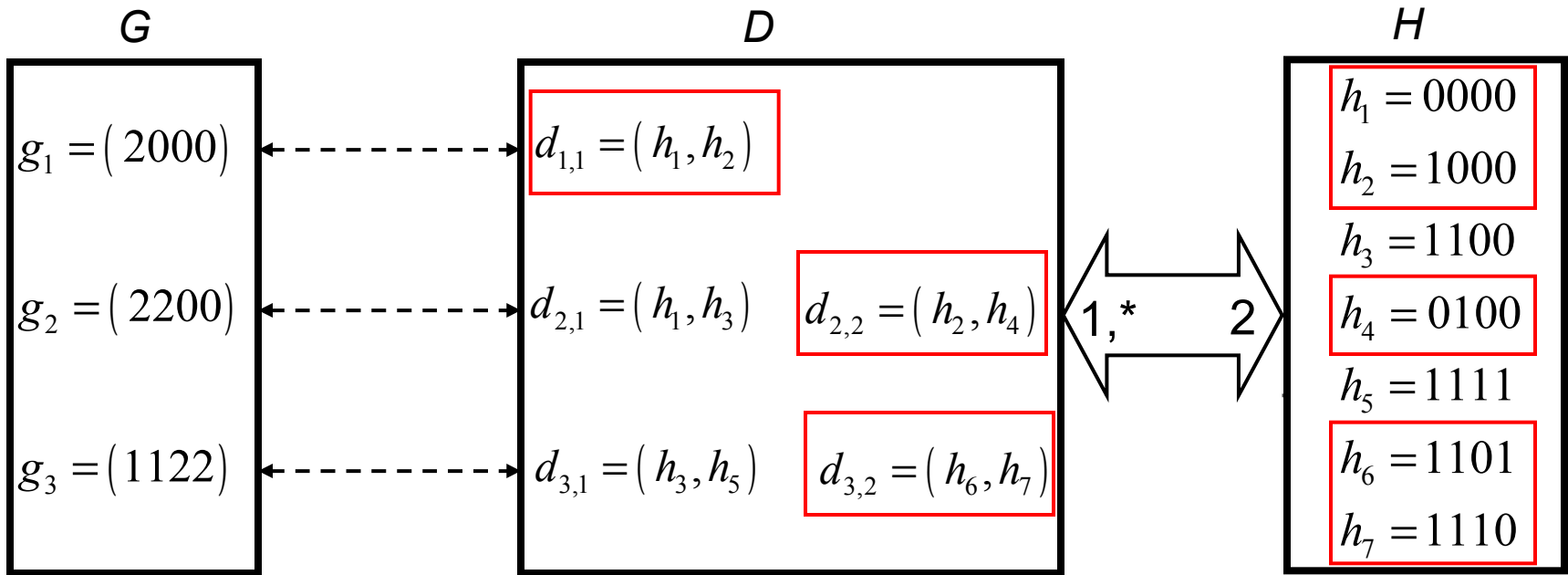
The individuals share a common evolutionary history

=>The individuals share common haplotypes

## Principle :

Minimize the number of distinct haplotypes required to solve all the individuals of the sample

# PARSIMONY

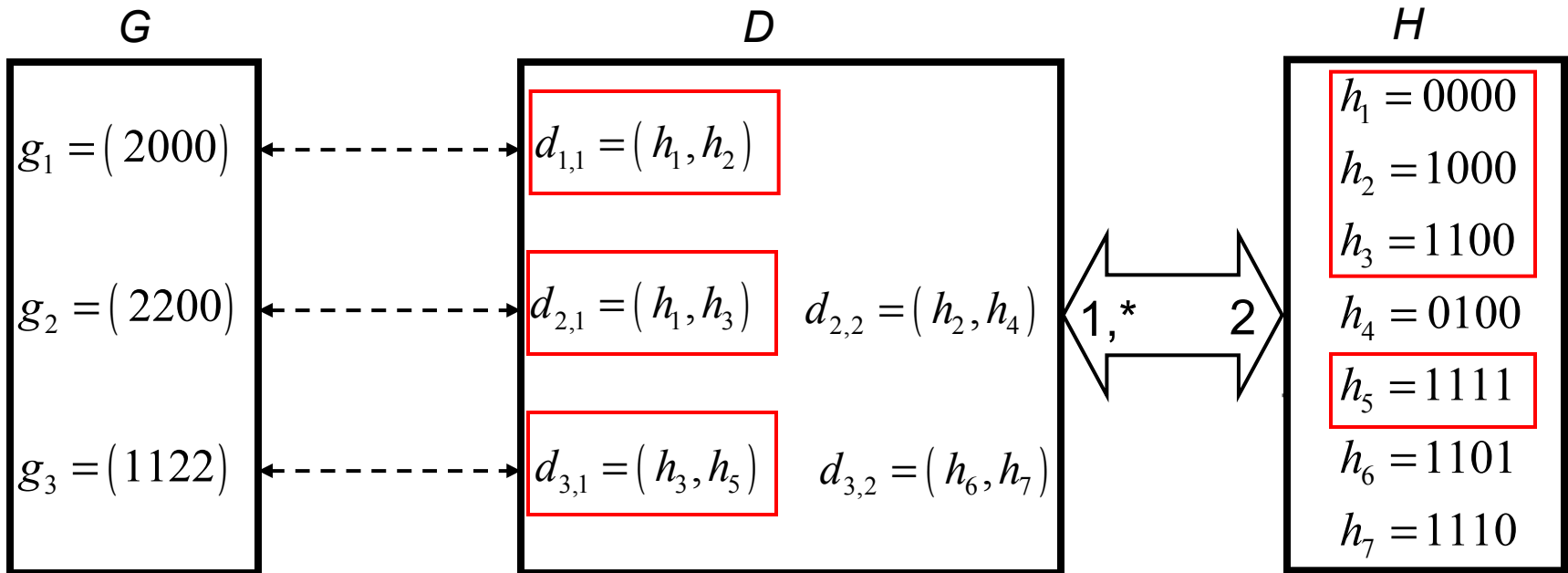


$$D_1 = \{d_{1,1}, d_{2,2}, d_{3,2}\}$$

$$H_1 = \{h_1, h_2, h_4, h_6, h_7\}$$

=> **5** haplotypes

# PARSIMONY



$$D_1 = \{d_{11}, d_{22}, d_{32}\}$$

$$H_1 = \{h_1, h_2, h_4, h_6, h_7\}$$

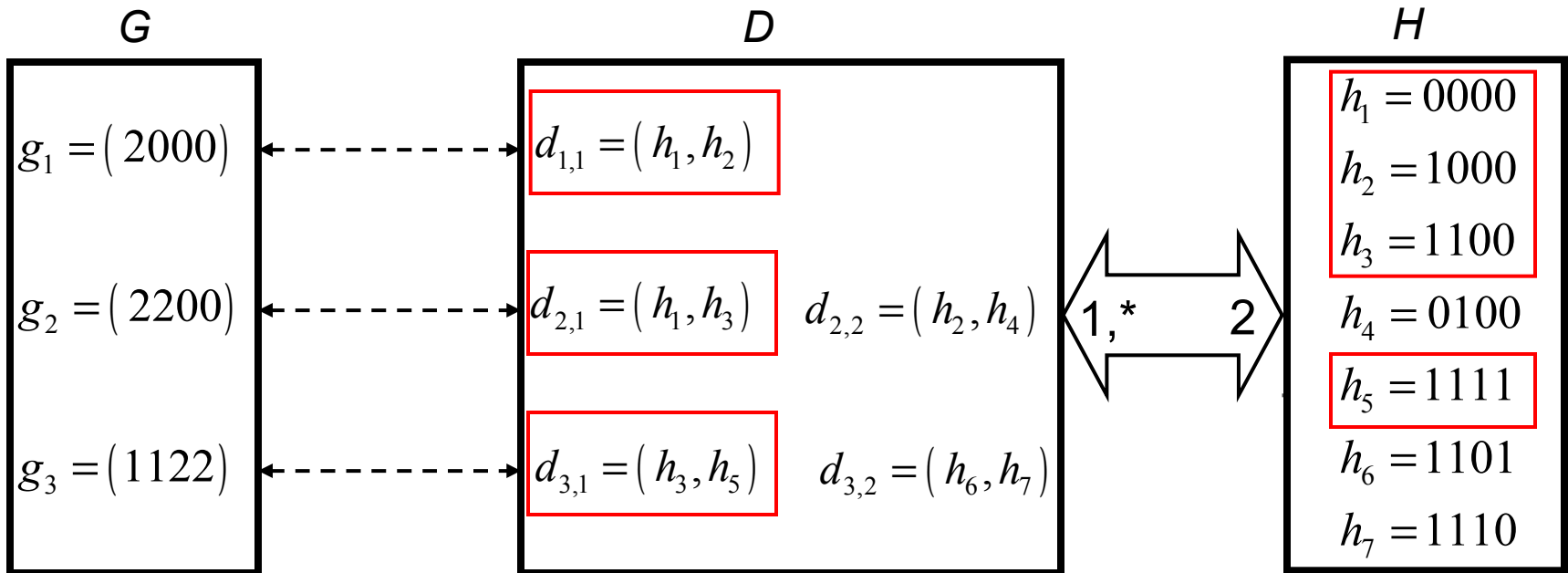
=> **5** haplotypes

$$D_2 = \{d_{11}, d_{21}, d_{31}\}$$

$$H_2 = \{h_1, h_2, h_3, h_5\}$$

=> **4** haplotypes

# PARSIMONY



$$D_1 = \{d_{11}, d_{22}, d_{32}\}$$

$$H_1 = \{h_1, h_2, h_4, h_6, h_7\}$$

=> **5** haplotypes

$$D_2 = \{d_{11}, d_{21}, d_{31}\}$$

$$H_2 = \{h_1, h_2, h_3, h_5\}$$

=> **4** haplotypes

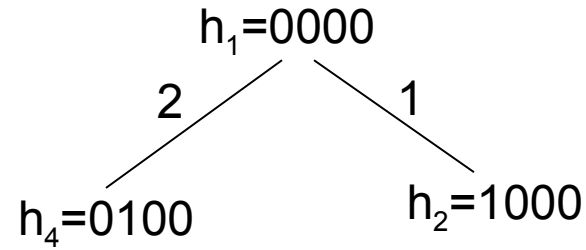
## Idea :

The haplotypes share a common evolutionary history which can be represented by a realistic phylogenetic tree

## Principle :

Find a set of haplotypes for  $G$  which fits a perfect phylogeny (each node is a haplotype, each branch is a mutation which appears only once)

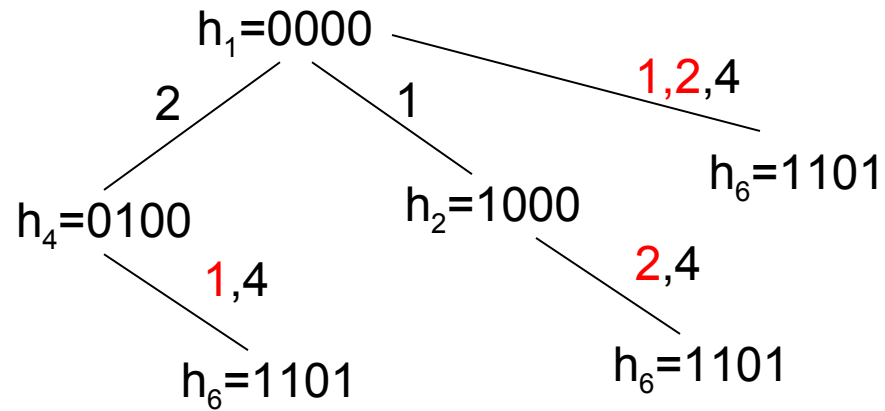
$$D_1 = \{d_{1,1}, d_{2,2}, d_{3,2}\} \quad H_1 = \{h_1=0000 \quad h_2=1000 \quad h_4=0100 \quad h_6=1101 \quad h_7=1110\}$$





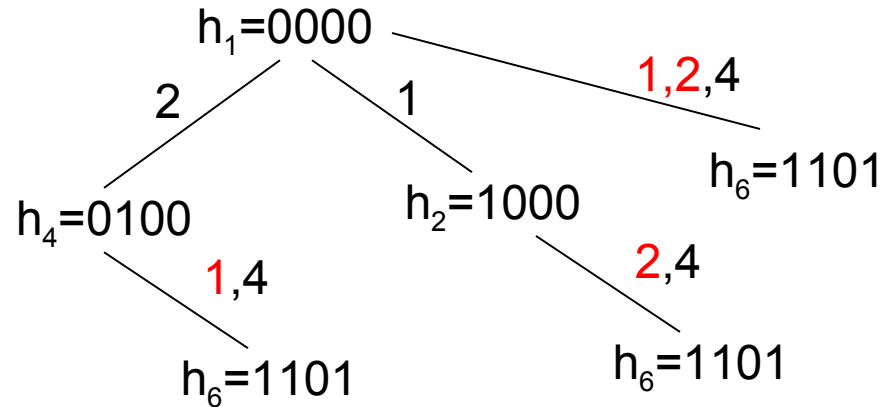
# PHYLOGENY

$$D_1 = \{d_{1,1}, d_{2,2}, d_{3,2}\} \quad H_1 = \{h_1=0000 \ h_2=1000 \ h_4=0100 \ h_6=1101 \ h_7=1110\}$$

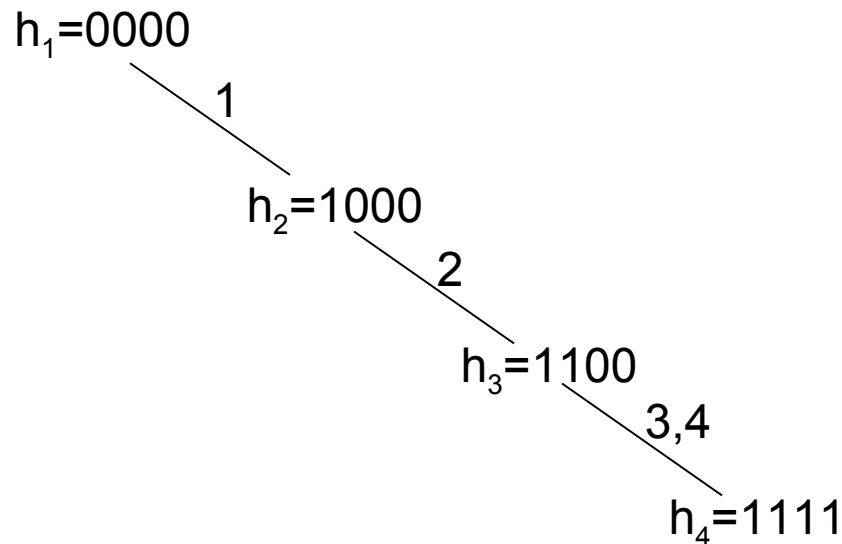


# PHYLOGENY

$$D_1 = \{d_{11}, d_{22}, d_{32}\} \quad H_1 = \{h_1=0000 \ h_2=1000 \ h_4=0100 \ h_6=1101 \ h_7=1110\}$$

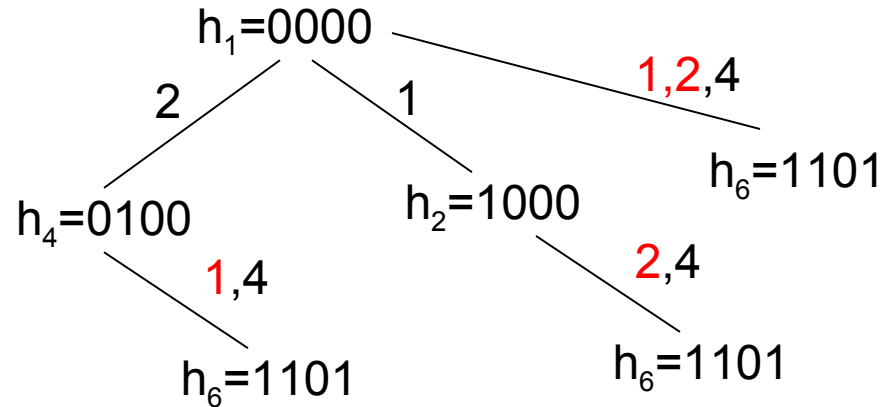


$$D_2 = \{d_{11}, d_{21}, d_{31}\} \quad H_2 = \{h_1=0000 \ h_2=1000 \ h_3=1100 \ h_5=1111\}$$

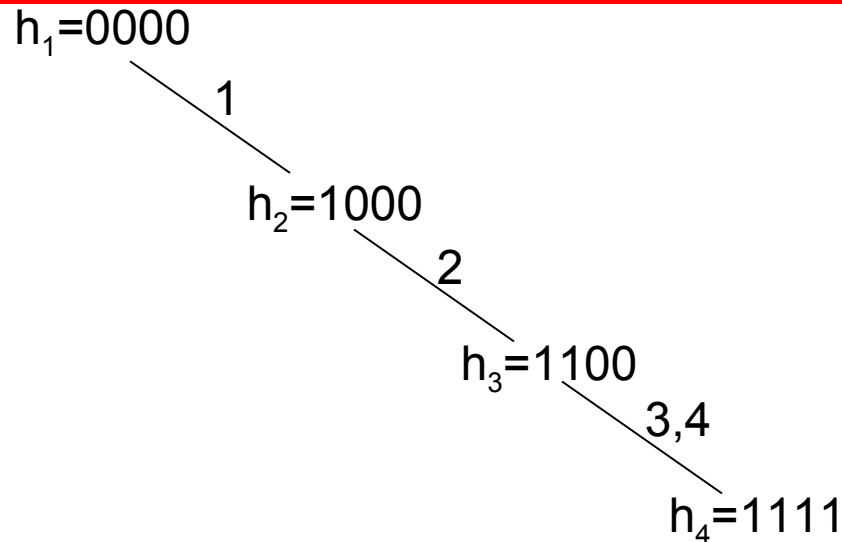


# PHYLOGENY

$$D_1 = \{d_{11}, d_{22}, d_{32}\} \quad H_1 = \{h_1=0000 \ h_2=1000 \ h_4=0100 \ h_6=1101 \ h_7=1110\}$$



$$D_2 = \{d_{11}, d_{21}, d_{31}\} \quad H_2 = \{h_1=0000 \ h_2=1000 \ h_3=1100 \ h_5=1111\}$$



I. Introduction

II. Combinatorial haplotype inference

**III. Statistical haplotype inference**

IV. Algorithmic tricks

V. Comparison of accuracy

VI. Application 1 : haplotype association tests

VII. Application 2 : genotype imputation

VIII. Application 3 : admixture

$H = \{ h_1, h_2, \dots, h_a, \dots, h_b, \dots, h_m \}$  with frequencies  $\{ f_1, f_2, \dots, f_a, \dots, f_b, \dots, f_m \}$

## I. Probability of a pair of haplotypes $(h_a, h_b)$ by Hardy Weinberg :

$$\Pr( g_i = (h_a, h_b) | H ) = 2 \times f_a \times f_b \quad \text{if } h_a \neq h_b$$

$$\Pr( g_i = (h_a, h_b) | H ) = f_a \times f_b \quad \text{if } h_a = h_b$$

II. Probability of a genotype by summing over all possible pairs of haplotypes :

$$\Pr( g_i | H ) = \sum_{(h_a, h_b)} \Pr( g_i = (h_a, h_b) | H )$$

II. Probability of a genotype by summing over all possible pairs of haplotypes :

$$\Pr( g_i | H ) = \sum_{(h_a, h_b)} \Pr( g_i = (h_a, h_b) | H )$$

III. Likelihood of the entire sample **G** by multiplying the N genotype probabilities :

$$\Pr( G | H ) = \prod_{i=1}^{i=N} \Pr( g_i | H )$$

II. Probability of a genotype by summing over all possible pairs of haplotypes :

$$\Pr( g_i | H ) = \sum_{(h_a, h_b)} \Pr( g_i = (h_a, h_b) | H )$$

III. Likelihood of the entire sample **G** by multiplying the N genotype probabilities :

$$\Pr( G | H ) = \prod_{i=1}^{i=N} \Pr( g_i | H )$$

IV. Maximum likelihood estimate  $\hat{H}$  by Expectation – Maximization (EM) :

$$\hat{H} = \max_H \Pr(G | H)$$



$D = \{ d_1, \dots, d_i, \dots, d_N \}$  where  $d_i$  is a random variable defined on all possible pairs of haplotypes of  $g_i$ .

Construct a Monte Carlo Markov chain with states  $D^0, D^1, \dots, D^t, \dots$  which converges to the true reconstruction of haplotypes of  $G$ .

Gibbs sampler :

**0.** Start with a random reconstruction  $D^0$

**1.** Iterate a large number of times the following step to go from state  $D^t$  to state  $D^{t+1}$  :

Construct  $D^{t+1}$  from  $D^t$  by sampling a new pair of haplotypes for a randomly chosen genotype  $g_i$  given all the  $2N - 2$  haplotypes previously sampled  $D^t_i$

$D = \{ d_1, \dots, d_i, \dots, d_N \}$  where  $d_i$  is a random variable defined on all possible pairs of haplotypes of  $g_i$ .

Construct a Monte Carlo Markov chain with states  $D^0, D^1, \dots, D^t, \dots$  which converges to the true reconstruction of haplotypes of  $G$ .

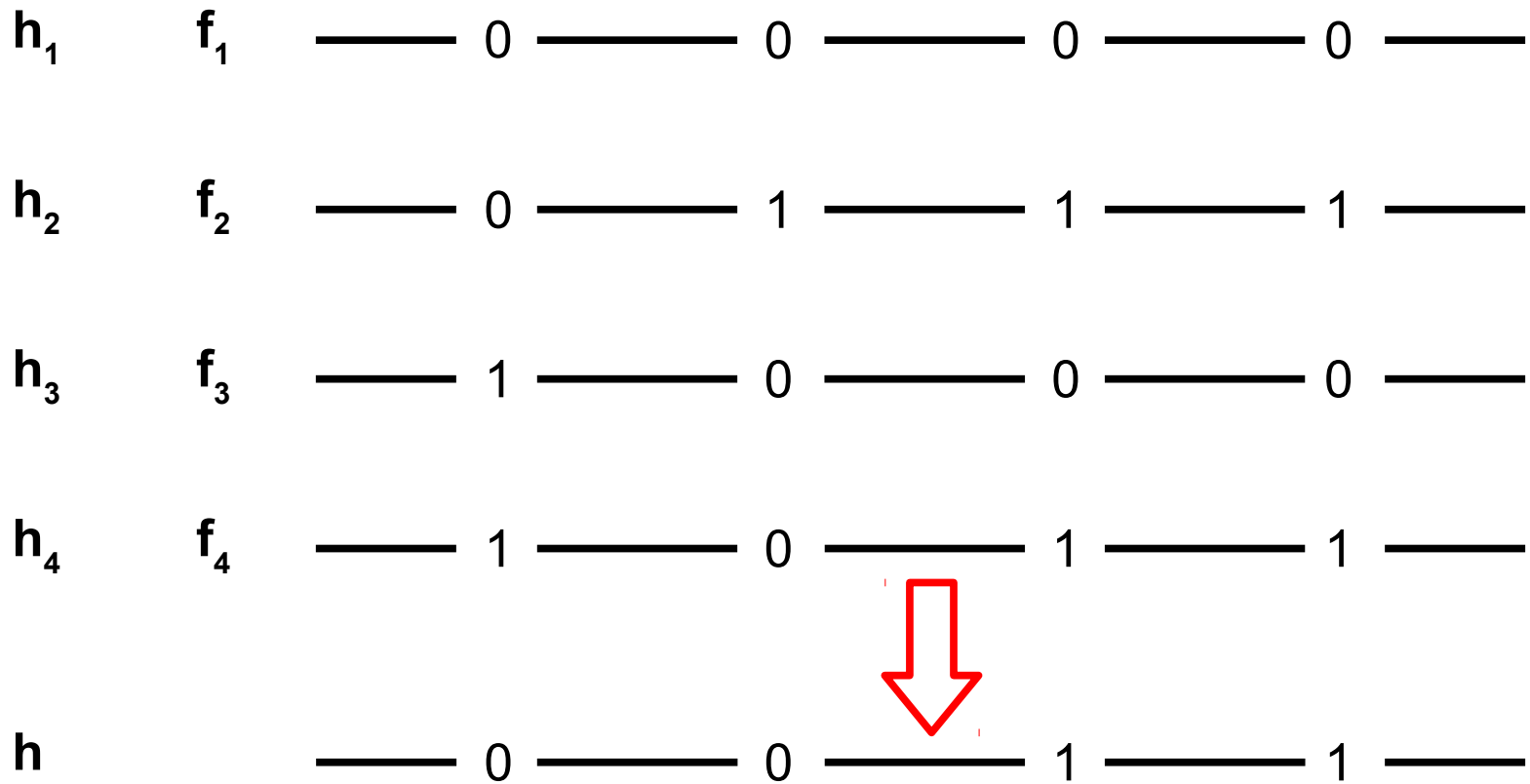
Gibbs sampler :

**0.** Start with a random reconstruction  $D^0$

**1.** Iterate a large number of times the following step to go from state  $D^t$  to state  $D^{t+1}$  :

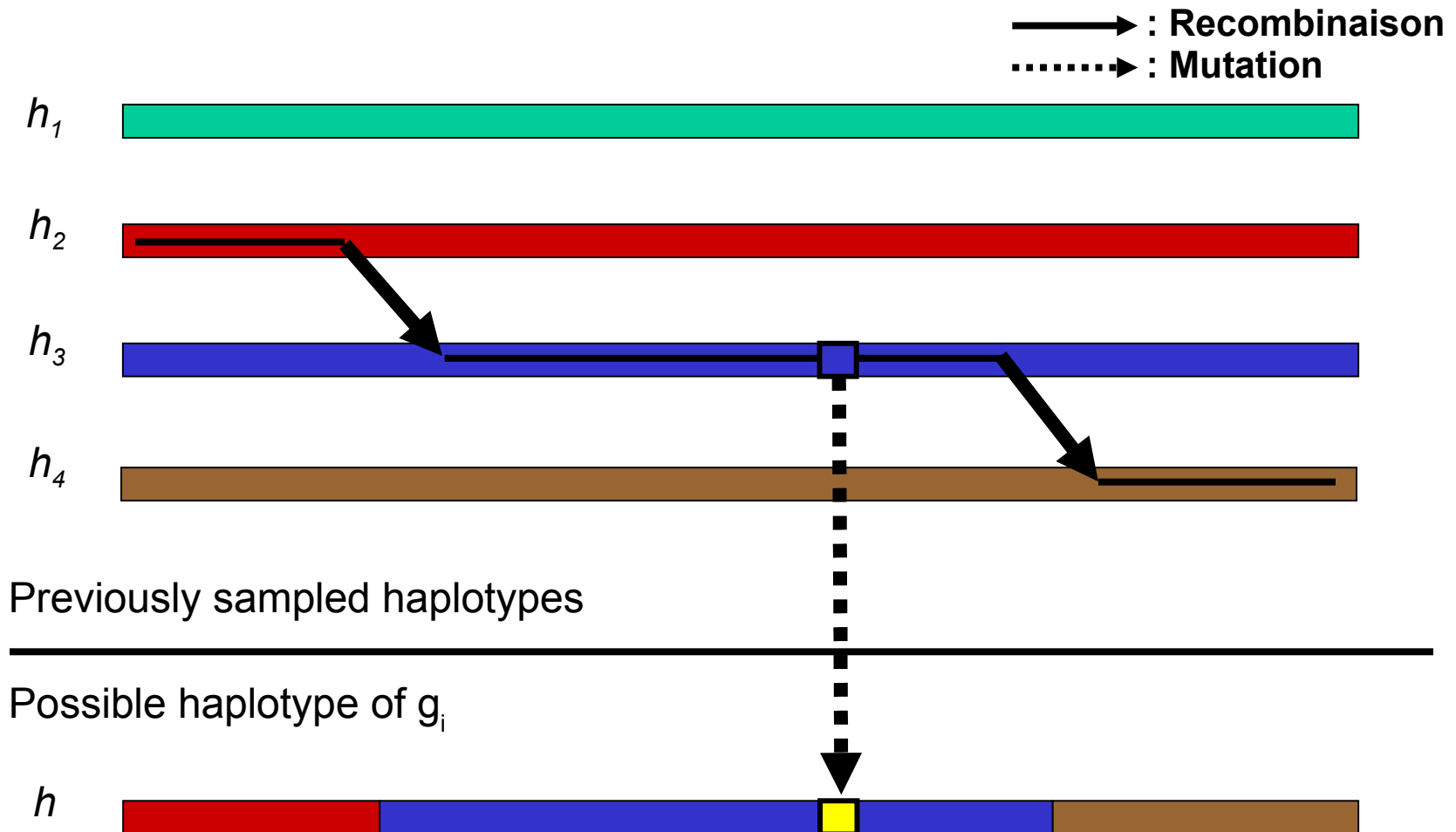
Construct  $D^{t+1}$  from  $D^t$  by sampling a new pair of haplotypes for a randomly chosen genotype  $g_i$  given all the  $2N - 2$  haplotypes previously sampled  $D_i^t$

# HIDDEN MARKOV MODEL OF HAPLOTYPES



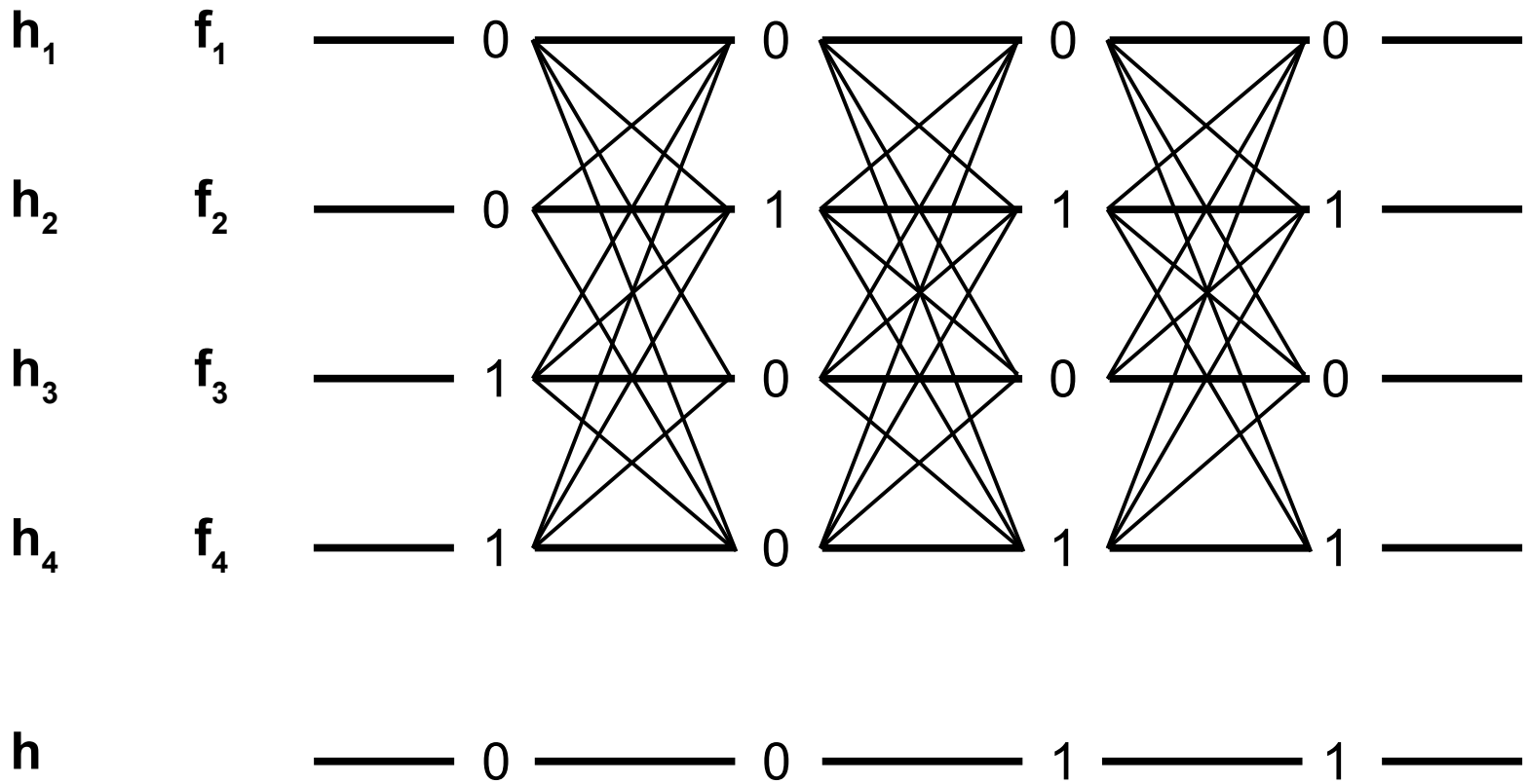
1. How likely  $(h_1, h_2, h_3, h_4)$  have produced  $h$  ?

# HIDDEN MARKOV MODEL OF HAPLOTYPES

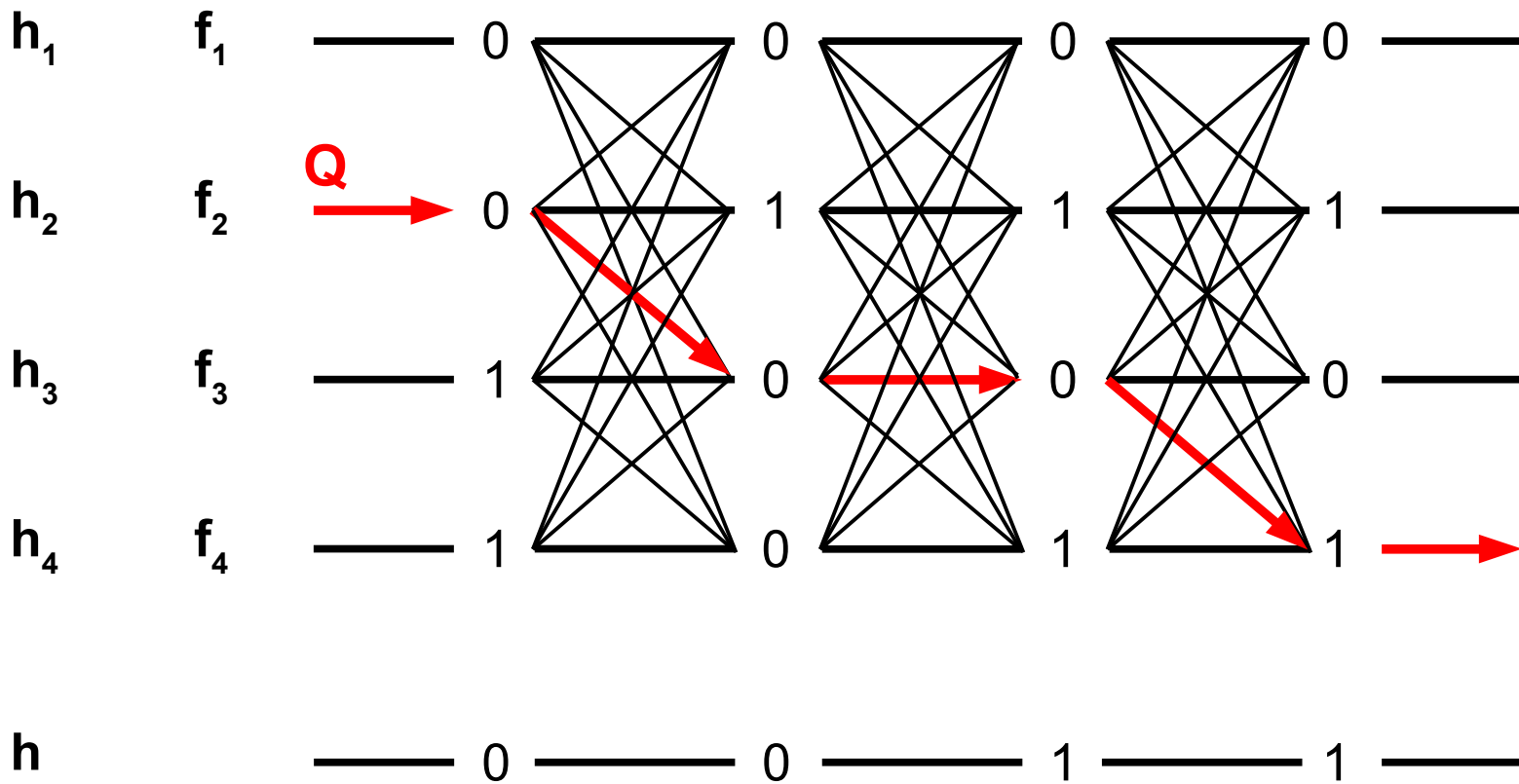


Haplotype  $h$  is an imperfect mosaic of  $(h_1, h_2, h_3, h_4)$

# HIDDEN MARKOV MODEL OF HAPLOTYPES

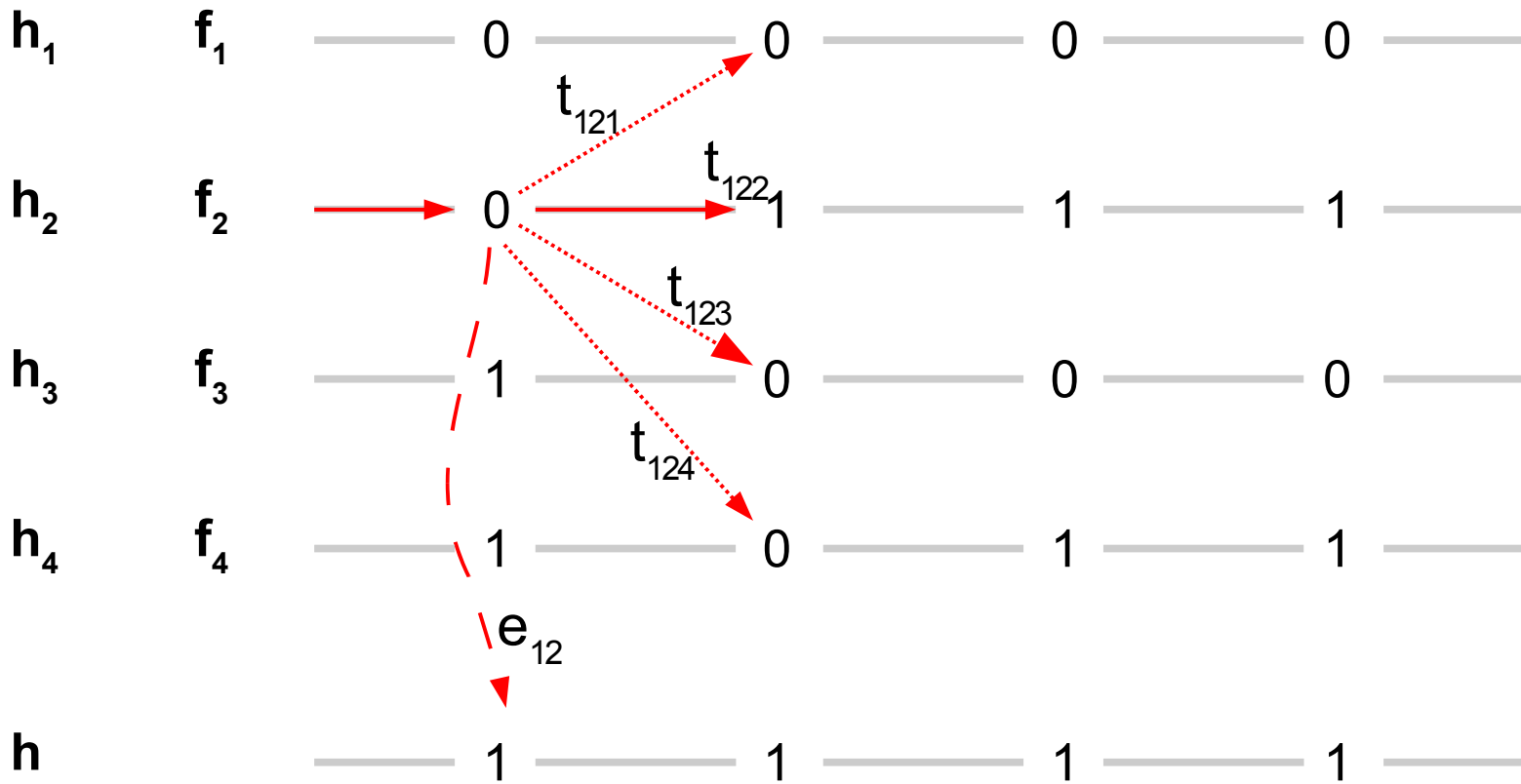


# HIDDEN MARKOV MODEL OF HAPLOTYPES



$Q$  is a path in the graph =  $Q$  is a mosaic of  $h_2, h_3$  and  $h_4$   
Haplotype  $h$  is an imperfect copy of  $Q$

# HIDDEN MARKOV MODEL OF HAPLOTYPES



$K = 4$  haplotypes  
 $L = 4$  SNPs

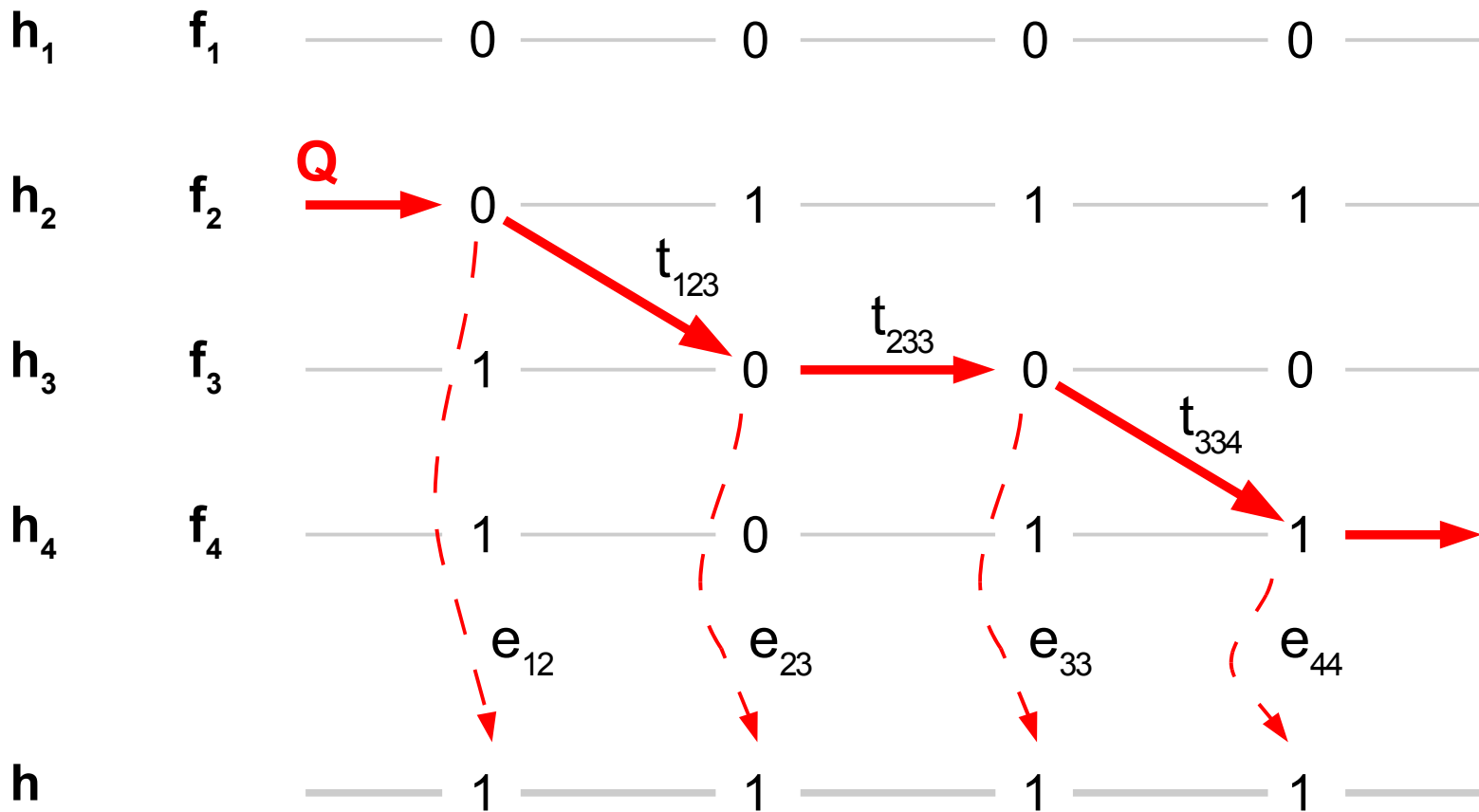
$$\text{HMM } \lambda = \begin{cases} F : \text{initial probabilities } f_k \\ T : \text{transition probabilities } t_{jkl} = \Pr(X_{j+1} = h_l | X_j = h_k) \\ E : \text{emission probabilities } e_{jk} = \Pr(o_j | X_j = h_k) \end{cases}$$

## 1. How to define the parameters ?

- Prior on the transitions : a function of the genetic distance between SNPs (a large genetic distance implies greater chance of recombination).
- Prior on the emissions : a function of the mutation rate



# HIDDEN MARKOV MODEL OF HAPLOTYPES



$$\Pr(Q | \lambda) = f_2 \times t_{123} \times t_{233} \times t_{334}$$

$$\Pr(h | Q, \lambda) = e_{12} \times e_{23} \times e_{33} \times e_{44}$$

$$\Pr(h, Q | \lambda) = \Pr(h | Q, \lambda) \times \Pr(Q | \lambda)$$

2. How likely  $(\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3, \mathbf{h}_4)$  can produce  $\mathbf{h}$  ?

$$\begin{aligned}\Pr(\mathbf{h} \mid \lambda) &= \sum_Q \Pr(\mathbf{h}, Q \mid \lambda) \\ &= \sum_Q \Pr(\mathbf{h} \mid Q, \lambda) \times \Pr(Q \mid \lambda)\end{aligned}$$

=> Naïve approach  $O(K^L)$

=> Forward-backward algorithm  $O(LK^2)$

I. Introduction

II. Combinatorial haplotype inference

III. Statistical haplotype inference

**IV. Algorithmic tricks**

V. Comparison of accuracy

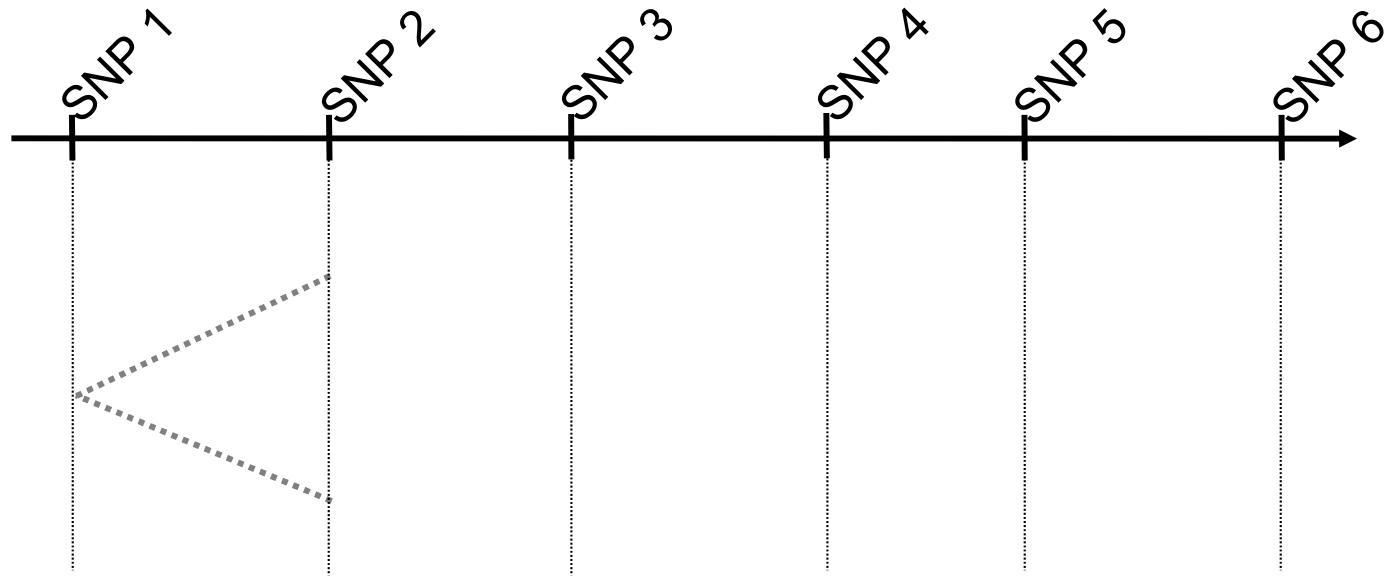
VI. Application 1 : haplotype association tests

VII. Application 2 : genotype imputation

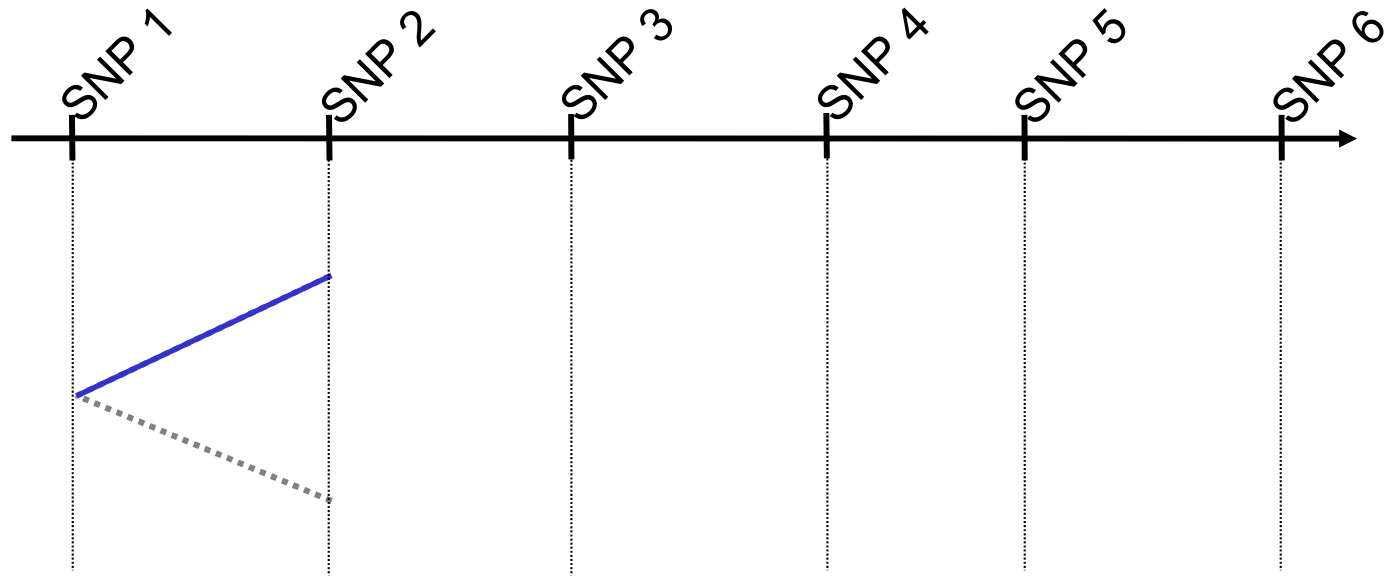
VIII. Application 3 : admixture



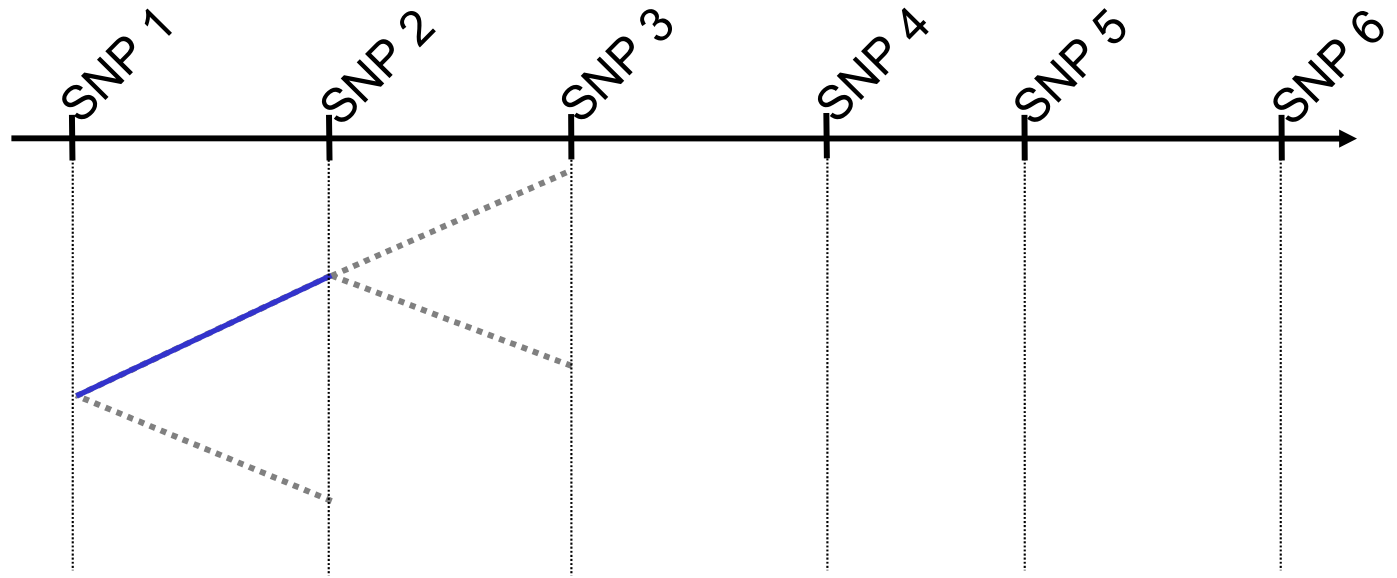
# ITERATIVE SNP INCLUSION



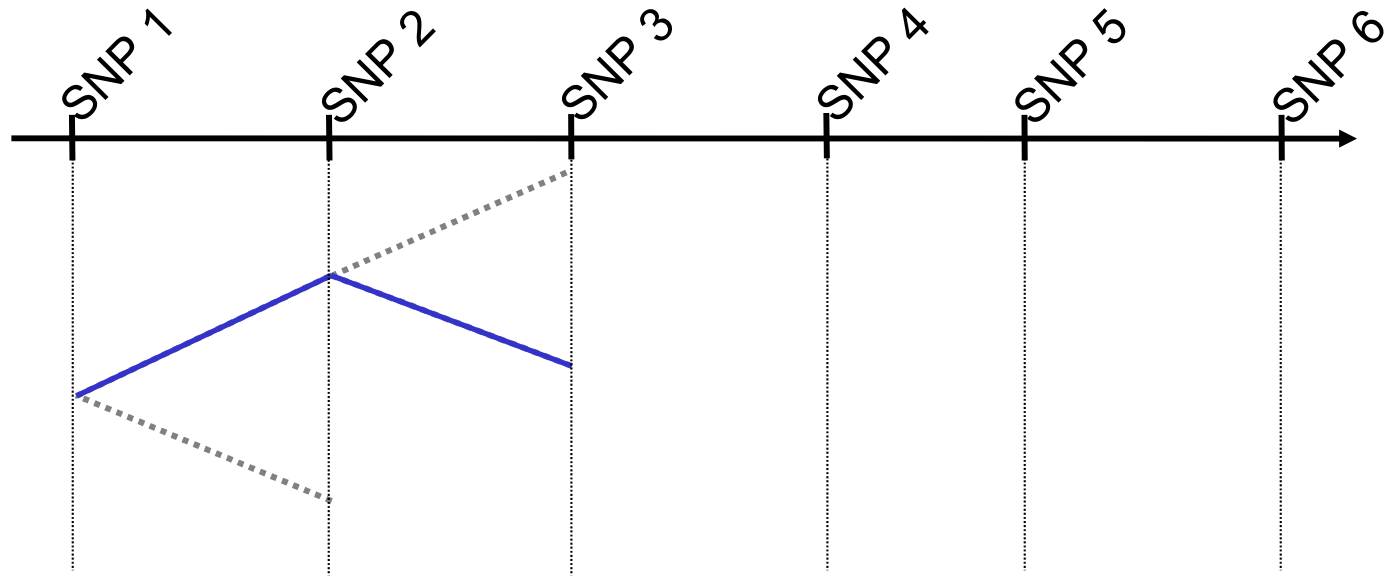
# ITERATIVE SNP INCLUSION



# ITERATIVE SNP INCLUSION

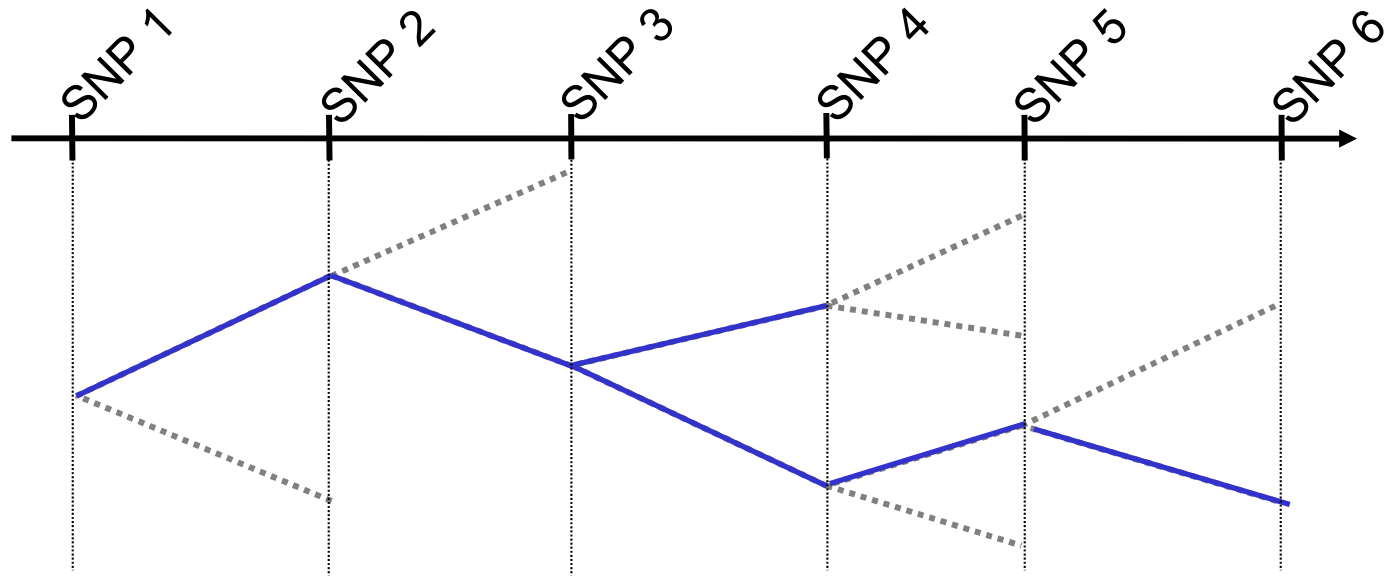


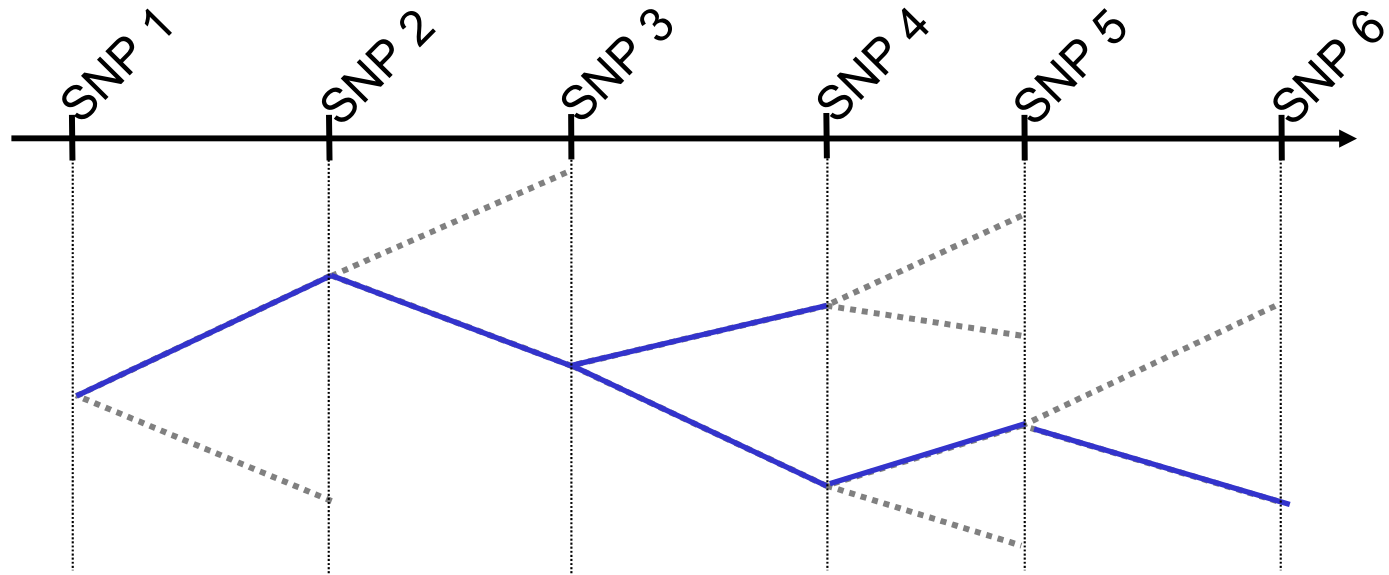
# ITERATIVE SNP INCLUSION





# ITERATIVE SNP INCLUSION





Implementation of the forward – backward algorithm with this tree strategy

*Delaneau et al, 2007*

*Delaneau et al, 2008*

I. Introduction

II. Combinatorial haplotype inference

III. Statistical haplotype inference

IV. Algorithmic tricks

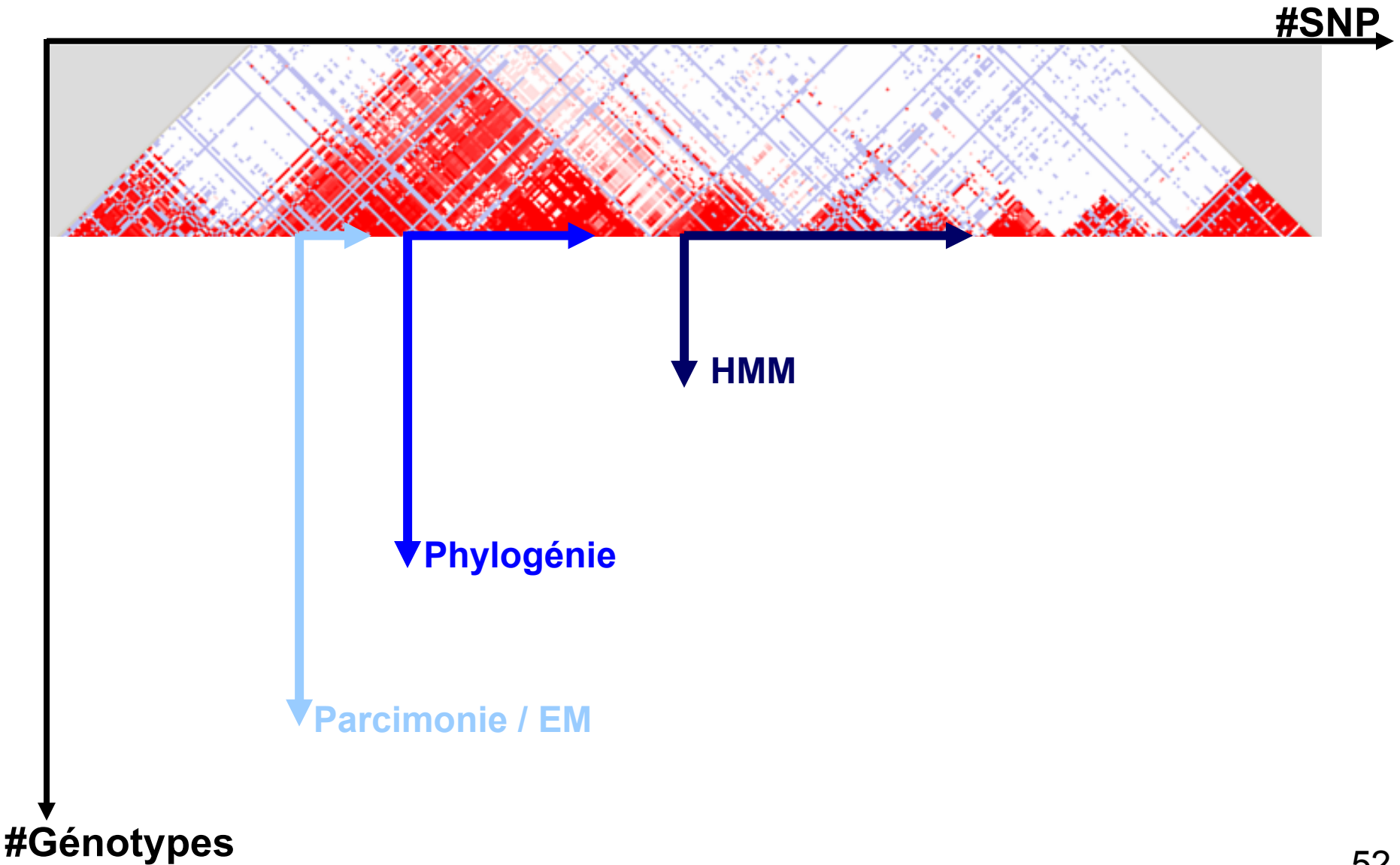
**V. Comparison of accuracy**

VI. Application 1 : haplotype association tests

VII. Application 2 : genotype imputation

VIII. Application 3 : admixture

# ACCURACY OF HAPLOTYPE INFERENCE



HapMap : 300 x (60 genotypes / 50 SNP)

	Error rate (%)	Running time (sec.)
ShapeIT	6,3	50
Phase v2.1	6,3	1215
EM	12,2	10

Studied genomic region



# HAPLOTYPE INFERENCE STRATEGY

Studied genomic region

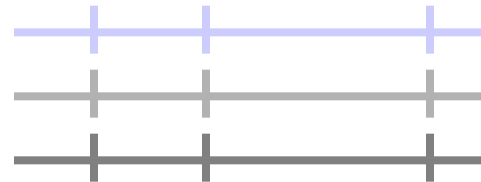


Local haplotype inference

1. Extraction of SNP



2. Haplotype inference



# HAPLOTYPE INFERENCE STRATEGY

Studied genomic region

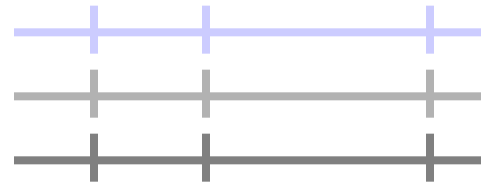


Local haplotype inference

1. Extraction of SNP

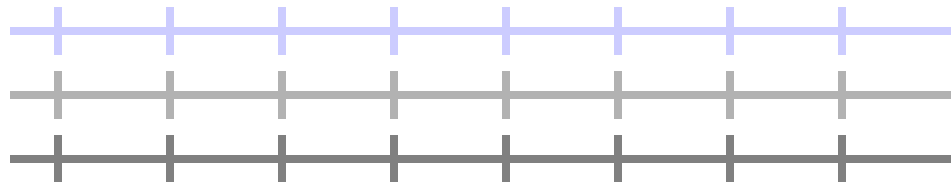


2. Haplotype inference

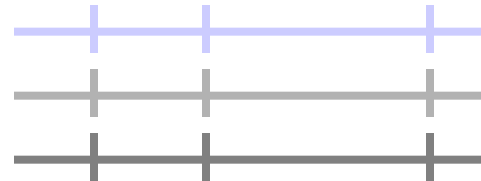


Global haplotype inference

1. Haplotype inference



2. Extraction of sub-haplotypes





# HAPLOTYPE INFERENCE STRATEGY

Studied genomic region

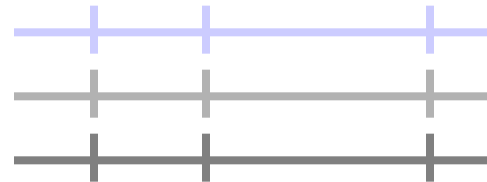


Local haplotype inference

1. Extraction of SNP

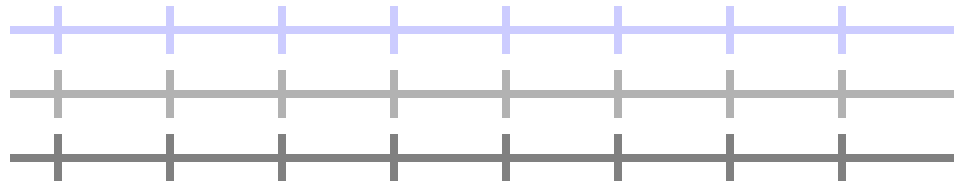


2. Haplotype inference

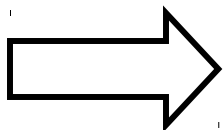
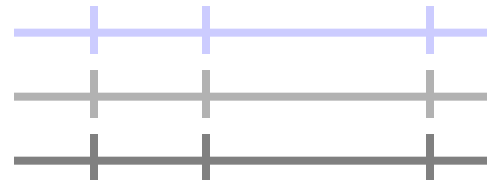


Global haplotype inference

1. Haplotype inference



2. Extraction of sub-haplotypes



Reduce from 1% to 6% the error rate

I. Introduction

II. Combinatorial haplotype inference

III. Statistical haplotype inference

IV. Algorithmic tricks

V. Comparison of accuracy

**VI. Application 1 : haplotype association tests**

VII. Application 2 : genotype imputation

VIII. Application 3 : admixture

The simplest approach :

1. Estimate haplotype frequencies in cases and control separately.
2. Find the most likely haplotype pair for each individual.
3. Fill in a contingency table to compare haplotypes in the cases versus the controls.

The simplest approach :

- ~~1. Estimate haplotype frequencies in cases and control separately.~~
2. Find the most likely haplotype pair for each individual.
3. Fill in a contingency table to compare haplotypes in the cases versus the controls.

Cases		Controls
$g_1$ : A/A A/A		$g_5$ : A/A T/T
$g_2$ : A/T A/T	Versus	$g_6$ : A/T A/T
$g_3$ : A/T A/T		$g_7$ : A/T A/T
$g_4$ : A/T A/T		$g_8$ : A/T A/T

Cases		Controls
$g_1$ : AA/AA		$g_5$ : AT/AT
$g_2$ : AA/TT		$g_6$ : AT/TA
$g_3$ : AA/TT		$g_7$ : AT/TA
$g_4$ : AA/TT		$g_8$ : AT/TA

**! Do not infer haplotypes in two samples separately !**

The simplest approach :

1. Estimate haplotype frequencies in cases and control separately.
- ~~2. Find the most likely haplotype pair for each individual.~~
3. Fill in a contingency table to compare haplotypes in the cases versus the controls.

g :     AA/TT with probability of 0.7  
          AT/AT with probability of 0.3

g :     AA/TT

**! Does not take into account uncertainty !**

Solution 1: Likelihood ratio test + label permutations

$$L_A = (Pr(G_{cases} | H))$$

$$L_B = (Pr(G_{controls} | H))$$

$$L_C = (Pr(G_{cases+controls} | H))$$

$$2 \ln \left( \frac{L_A \times L_B}{L_C} \right) \sim \text{chi}^2$$

## Solution 2 : Regression model

genotype	phenotype	hap pair 1	prob 1	hap pair 2	prob 2
g1	p1	h1/h1	1.0		
g2	p2	h1/h2	0.7	h3/h4	0.3
g3	p3	h3/h4	1.0		

# HAPLOTYPE ASSOCIATION TESTS

## Solution 2 : Regression model

genotype	phenotype	hap pair 1	prob 1	hap pair 2	prob 2
g1	p1	h1/h1	1.0		
g2	p2	h1/h2	0.7	h3/h4	0.3
g3	p3	h3/h4	1.0		

$$\begin{pmatrix} p1 \\ p2 \\ p3 \end{pmatrix} = \begin{matrix} & \mu & h1 & h2 & h3 & h4 \\ \begin{pmatrix} 1.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 1.0 & 0.35 & 0.35 & 0.15 & 0.15 \\ 1.0 & 0.0 & 0.0 & 0.5 & 0.5 \end{pmatrix} & \cdot & \begin{pmatrix} \mu \\ \beta1 \\ \beta2 \\ \beta3 \\ \beta4 \end{pmatrix} \end{matrix}$$



## Solution 2 : Regression model

genotype	phenotype	hap pair 1	prob 1	hap pair 2	prob 2
g1	p1	h1/h1	1.0		
g2	p2	h1/h2	0.7	h3/h4	0.3
g3	p3	h3/h4	1.0		

$$\begin{pmatrix} p1 \\ p2 \\ p3 \end{pmatrix} = \begin{pmatrix} \mu & h1 & h2 & h3 & h4 \\ 1.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 1.0 & 0.35 & 0.35 & 0.15 & 0.15 \\ 1.0 & 0.0 & 0.0 & 0.5 & 0.5 \end{pmatrix} \cdot \begin{pmatrix} \mu \\ \beta1 \\ \beta2 \\ \beta3 \\ \beta4 \end{pmatrix}$$

Testing simultaneously all haplotypes

$$H_0 : \beta_1=0 \ \beta_2=0 \ \beta_3=0 \ \beta_4=0 \quad \text{VS} \quad H_1 : \beta_1 \neq 0 \ \beta_2 \neq 0 \ \beta_3 \neq 0 \ \beta_4 \neq 0$$

Testing only haplotype h1

$$H_0 : \beta_1=0 \quad \text{VS} \quad H_1 : \beta_1 \neq 0$$

### Strategies :

1. Biological meaning : promoters, exons, etc...
2. Most significant regions obtained by single SNP analysis
3. Systematic screening by sliding window

I. Introduction

II. Combinatorial haplotype inference

III. Statistical haplotype inference

IV. Algorithmic tricks

V. Comparison of accuracy

VI. Application 1 : haplotype association tests

**VII. Application 2 : genotype imputation**

VIII. Application 3 : admixture



HapMap haplotypes :

G	A	G	A	T	C	T	C	C	T	T	C	T	T	C	T	G	T	G	C
G	A	G	A	T	C	T	C	C	C	G	A	C	C	T	C	A	T	G	G
C	A	A	G	C	T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
G	A	A	G	C	T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
G	A	G	A	C	T	C	T	C	C	G	A	C	C	T	T	A	T	G	C
G	G	G	A	T	C	T	C	C	C	G	A	C	C	T	C	A	T	G	G

Genotype :

.	.	.	A	.	.	.	.	.	.	.	C	.	.	.	.	A	.	.	.
.	.	.	G	.	.	.	.	.	.	.	A	.	.	.	.	A	.	.	.



HapMap haplotypes :

G	A	G	A	T	C	T	C	C	T	T	C	T	T	C	T	G	T	G	C
G	A	G	A	T	C	T	C	C	C	G	A	C	C	T	C	A	T	G	G
C	A	A	G	C	T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
G	A	A	G	C	T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
G	A	G	A	C	T	C	T	C	C	G	A	C	C	T	T	A	T	G	C
G	G	G	A	T	C	T	C	C	C	G	A	C	C	T	C	A	T	G	G

Genotype :

.	.	.	A	.	.	.	.	.	.	.	C	.	.	.	.	A	.	.	.
.	.	.	G	.	.	.	.	.	.	.	A	.	.	.	.	A	.	.	.

Haplotypes :

.	.	.	A	.	.	.	.	.	.	.	A	.	.	.	.	A	.	.	.
.	.	.	G	.	.	.	.	.	.	.	C	.	.	.	.	A	.	.	.



HapMap haplotypes :

G	A	G	A	T	C	T	C	C	T	T	C	T	T	C	T	G	T	G	C
G	A	G	A	T	C	T	C	C	C	G	A	C	C	T	C	A	T	G	G
C	A	A	G	C	T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
G	A	A	G	C	T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
G	A	G	A	C	T	C	T	C	C	G	A	C	C	T	T	A	T	G	C
G	G	G	A	T	C	T	C	C	C	G	A	C	C	T	C	A	T	G	G

Genotype :

.	.	.	A	.	.	.	.	.	.	.	C	.	.	.	.	A	.	.	.
.	.	.	G	.	.	.	.	.	.	.	A	.	.	.	.	A	.	.	.

Haplotypes :

.	.	.	A	.	.	.	.	.	.	.	A	.	.	.	.	A	.	.	.
.	.	.	G	.	.	.	.	.	.	.	C	.	.	.	.	A	.	.	.

Imputation :

g	a	g	A	t	c	t	c	c	c	g	A	c	c	t	c	A	t	g	g
g	a	a	G	c	t	c	t	t	t	t	C	t	t	t	c	A	t	g	g



HapMap haplotypes :

```

G A G A T C T C C T T C T T C T G T G C
G A G A T C T C C C G A C C T C A T G G
C A A G C T C T T T T C T T C T G T G C
G A A G C T C T T T T C T T C T G T G C
G A G A C T C T C C G A C C T T A T G C
G G G A T C T C C C G A C C T C A T G G

```

Genotype :

```

. . . A . . . . . C . . . . A . . .
. . . G . . . . . A . . . . A . . .

```

Haplotypes :

```

. . . A . . . . . A . . . . A . . .
. . . G . . . . . C . . . . A . . .

```

Imputation :

```

g a g A t c t c c c g A c c t c A t g g
g a a G c t c t t t t C t t t c A t g g

```

In practice, we have a probability distribution :

c/c	c/t	t/t
0.01	0.18	0.81

I. Introduction

II. Combinatorial haplotype inference

III. Statistical haplotype inference

IV. Algorithmic tricks

V. Comparison of accuracy

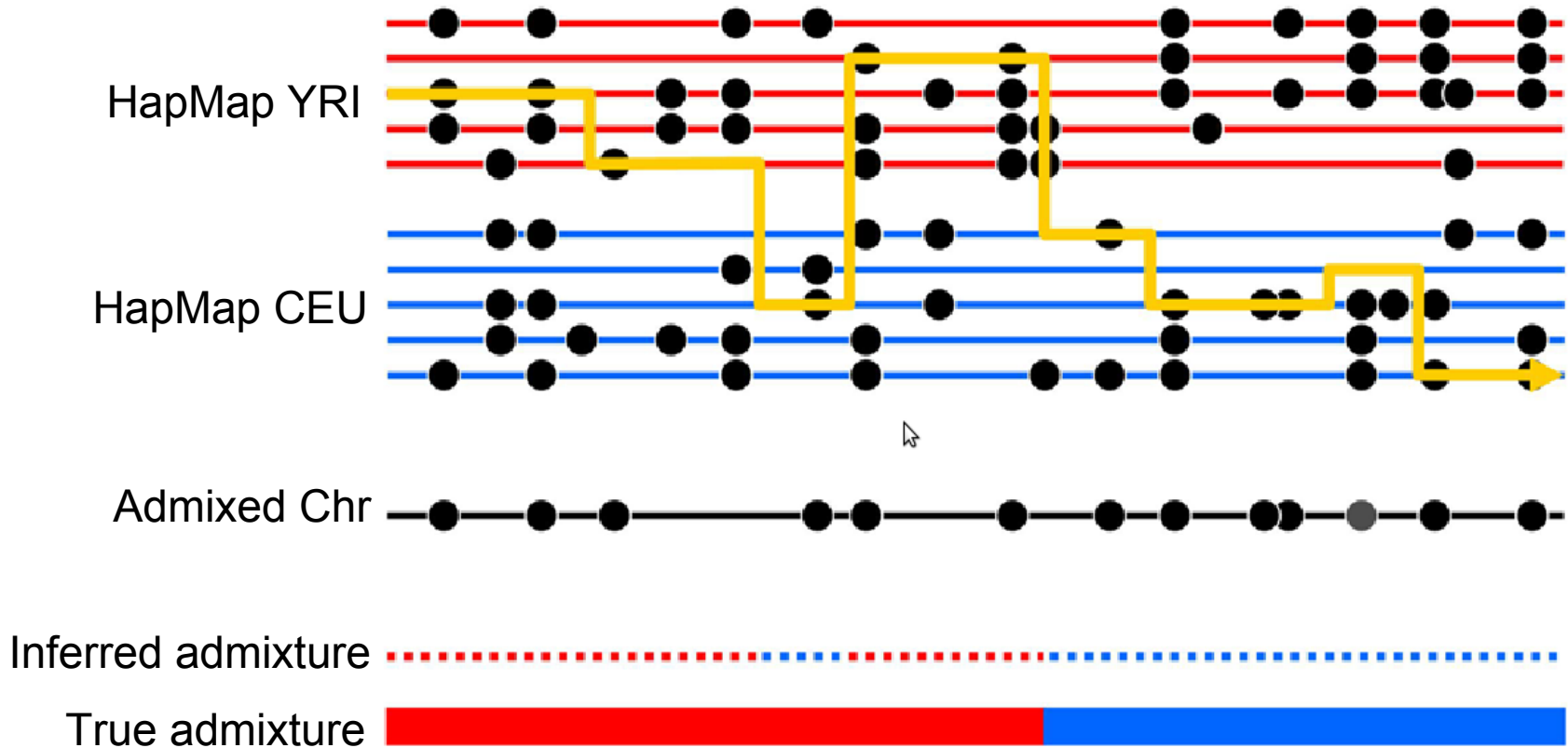
VI. Application 1 : haplotype association tests

VII. Application 2 : genotype imputation

**VIII. Application 3 : admixture**



# ADMIXTURE



1. Most accurate haplotype inference methods rely on hidden Markov model.
2. Algorithmic improvements still needed to fit the models on large datasets.
3. Haplotype association tests remain an interesting approach in association studies.
4. Several useful applications rely on haplotype inference: genotype imputation, admixture inference, etc ...

## Chaire de Bioinformatique du CNAM

Jean-François Zagury

Cédric Coulonges

Sigrid Le Clerc

Sophie Limou

Lieng Taing

Hervé Do

Taoufik Labib

Christiane Morel

Matthieu Montes

Rojo Ratsimandresy

Les nouveaux collaborateurs



## CNRS

Christine Sinoquet