# Haplotype Inference

22.1.2021

**Idea:** diploid

A  G/T  chr 7

C  chr 7

**2 haplotypes (phased)**

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

C        T

**1 genotype**

(unphased)

↑
segregating site
= heterozygous allele

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

0
0
1
1

Assumption: haplotypes are binary
↪ encoded by 0/1

# 3 approaches:

1) population based haplotyping
2) genetic haplotyping
3) molecular haplotyping

## ad 1)

6 individuals

**1** $\underline{11000}$ / $\underline{01100}$ → 21200

**2** $\underline{11000}$ / $\underline{11000}$ → 11000

**3** $\underline{11100}$ / $\underline{11010}$ → 11220

**4** $\underline{11000}$ / $\underline{11100}$ → 11200

**5**

**6**

haplotypes 0/1-seq.

genotypes

$$\frac{0/1/2}{both}$$

easy: haplotypes → genotype

hard: genotype → haplotypes !

approach:

1) look for individuals with 0 or 1 heterozygous sites
→ here haplotyping is trivial

1)

indiv. 2 $\longrightarrow$ 11000

indiv. 4 $\longrightarrow$

~~11000~~

11100

2) 11000 + indiv.3 : 11220

$\longrightarrow$ 11110

11100 + ind.1 : 21200

$\longrightarrow$ 01000

OR

11000 + ind.1 = 21200

$\longrightarrow$ 01100

Problem :

we can get different solutions

ad 2) Like above, but with pedigree information.

Example 1:

parent 1: 201
parent 2: 022
child: 221 → [101] [011]

green: inferred haplotypes

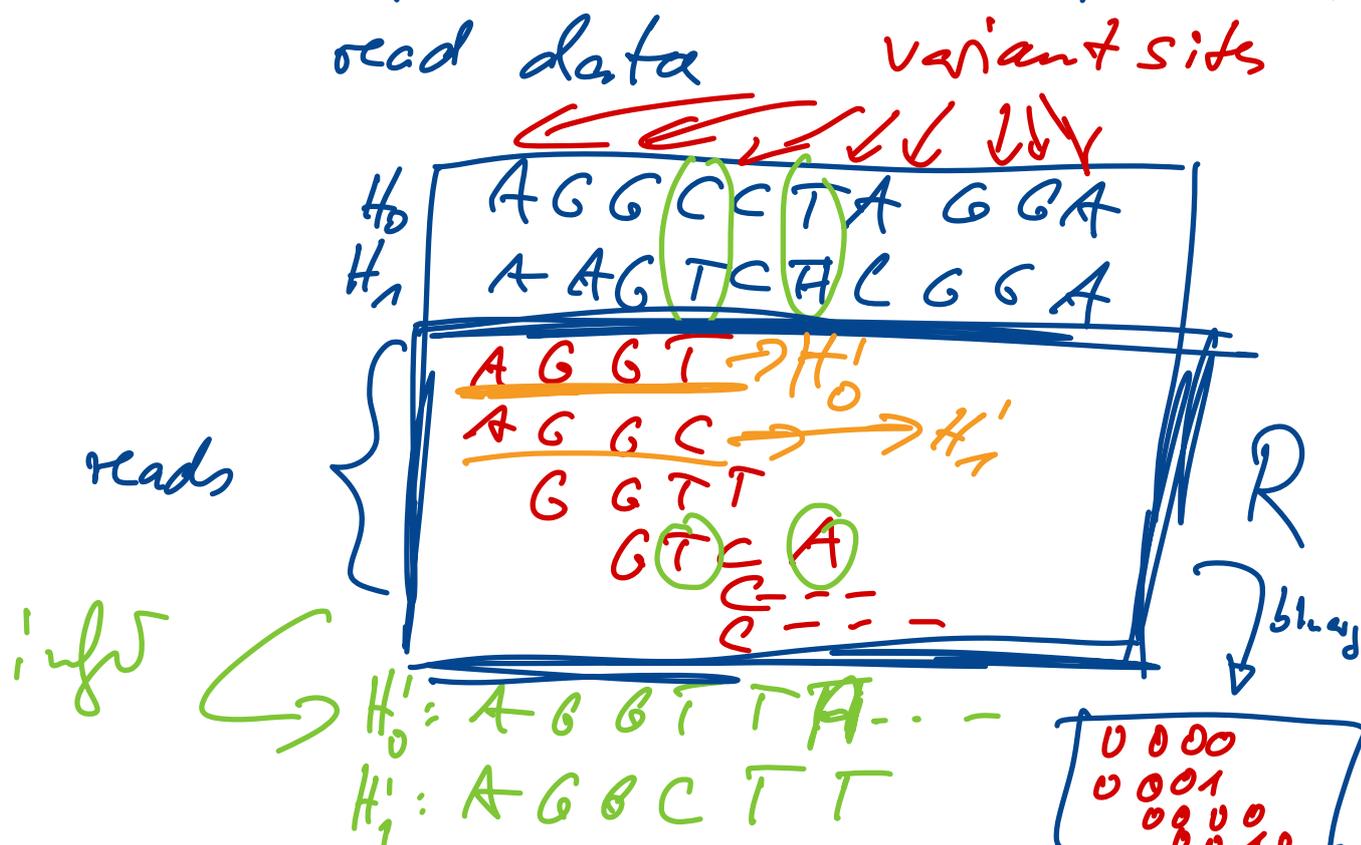→ parent 1: 101, 001
parent 2: 011, 000

✓

Example 2: parent 1: 222
parent 2: 222
child: 222

↳ nothing can be inferred

<u>ad 3)</u> Haplotype Inference from Sequencing
read data     <span style="color:red">variant sites</span>

$H_0$: A G G C C T A G G A
$H_1$: A A G T C A C G G A

reads {
A G G T $\rightarrow H_0'$
A G G C $\rightarrow H_1'$
G G T
G T C A
C - - -
C - - -
} R

blueq

info $\Rightarrow$
$H_0'$: A G G T T A - - -
$H_1'$: A G G C T T

$$\begin{matrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{matrix}$$

<span style="color:red">$c(1,1)=0$   $c(1,2)=0$</span>

<u>Haplotype Assembly Problem</u> (a.k.a. Minimum Error Correction Problem, MEC)

<u>Input</u>: $n \times m$ <u>fragment matrix</u> $R$ (usually binary)

<u>Def</u>: <u>conflict</u> of two <u>fragments</u> if they are unequal at some site

<u>Def</u>: Let $\underline{d(A,B)}$ be the number of corrections to transform fragment matrix $A$ into fragment matrix $B$

<u>Problem</u>: Given a fragment matrix $R$, find a fragment matrix $R'$ that has a conflict-free <u>bipartition</u>, minimizing $d(R, R')$.

# Result: MEC is NP-complete.
## (reduction from MAX CUT.)

# Solution 1:

(ILP)

$n \left\{ \boxed{R} \right._m$

An Integer Linear Program solving the Haplotype Assembly Problem:

$$\boxed{\text{Minimize} \sum_{j=1}^{j=n} z(j)}$$

$z(j)$ = mismatch counter for read $j$

Such that:

For each *Read* $j$ from 1 to $n$, and each position $k$ from 1 to $m$:

$z_0$ : like $z$, assuming that read $j$ is assigned to $H_0$

$z_1$ : like $z$, assuming that read $j$ is assigned to $H_1$

$z(j,k)$ : mismatch indicator for read $j$ in site $k$ in assigned haplotype

$$\boxed{\begin{aligned} z_0(j,k) &\geq c(j,k) - H_0'(k) \\ z_0(j,k) &\geq H_0'(k) - c(j,k) \\ z_1(j,k) &\geq c(j,k) - H_1'(k) \\ z_1(j,k) &\geq H_1'(k) - c(j,k) \\ z_0(j,k) - A(j) - z(j,k) &\leq 0 \\ z_1(j,k) + A(j) - z(j,k) &\leq 1 \end{aligned}}$$

$c(j,k)$ : input $R$

For each *Read* $j$:

$$\boxed{z(j) = \sum_{k=1}^{k=m} z(j,k)}$$

All variables are binary

$$A(j) = \begin{cases} 0 & \text{if read } j \text{ is in } H_0' \\ 1 & \text{if read } j \text{ is in } H_1' \end{cases}$$
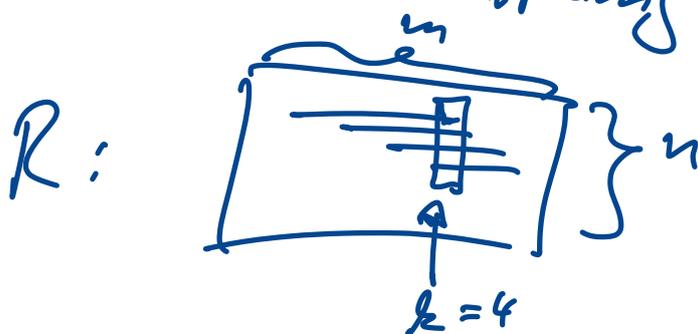
$H_0'$, $H_1'$ the two haplotypes

# Solution 2: WhatsHap (Patterson et al. 2015)

## dynamic programming solution along the columns of $R$.

## analysis: $O(2^{k-1} \cdot m)$ time

where $k$ is the maximum coverage at any site.



$R:$   $k = 4$

$\rightarrow$ FPT algorithm