# SNV - disease association mapping
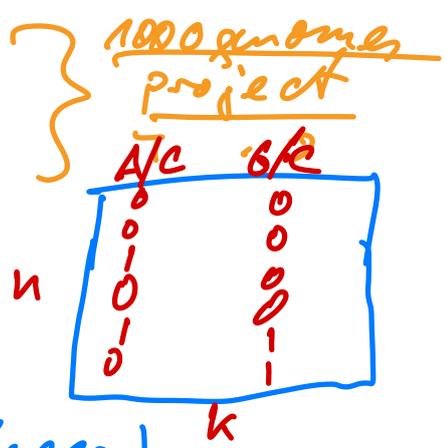
SNV : single nucleotide variant
(population genomics)
(VCF file format)



n=3
families/
individuals

$H$
$D$
$D$

} many individuals

H = healthy | type 1
D = disease | type 2

segregating sites

constant sites
(non-segregating)

→ 3 GB
81 MB

typical data :
n = thousands of individuals
k = millions of segregating sites

} 1000 genomes project
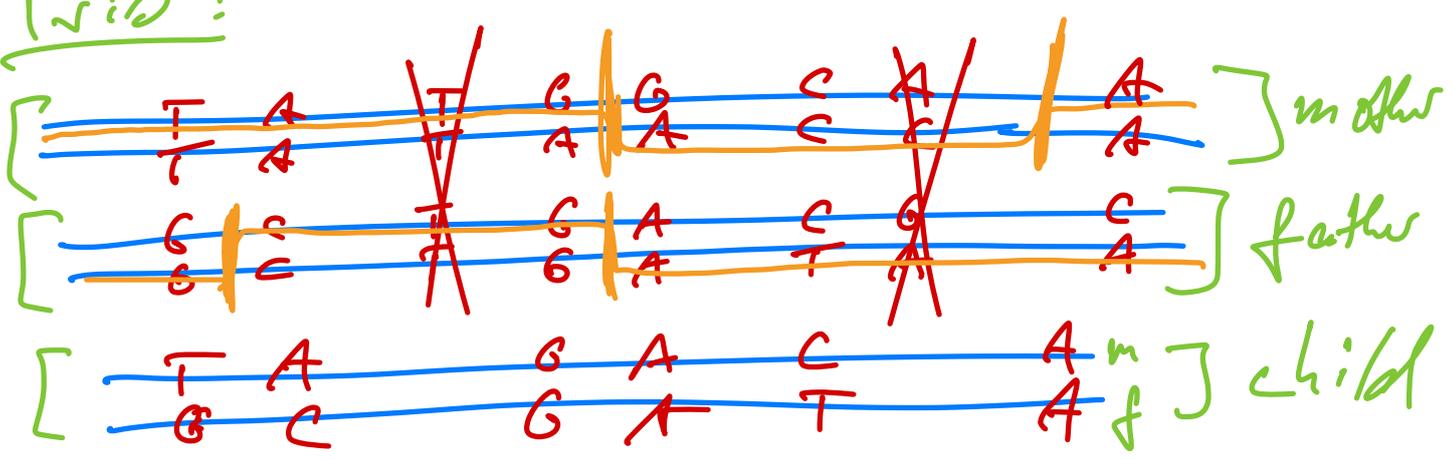
→ 4 × 10^11 bits
"big data"

(n = 2500 → 5000 haplotypes)
(k = 81.000.000 segregating sites)

A/C    G/C

| | | |
|---|---|---|
| 0 | | 0 |
| 0 | | 0 |
| 1 | | 0 |
| 0 | | 0 |
| 1 | | 1 |

n    k

Trio :



} mother
} father
] child

# correlated segregating sites

chromosome

## haplotype blocks

| | | | | | |
|---|---|---|---|---|---|
| . | . . | A A C A ① | - | - - | H |
| | | T T C A ② | | | H |
| | | A A C A | | | H |
| | | T A G G ③ | | | D |
| | | T A G G | | | D |
| | | A A C A | | | H |
| | | T A G G | | | D |

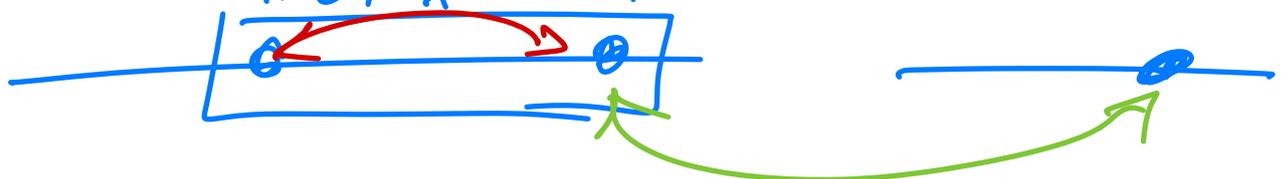need: ~~[ - - C G ]~~        ~~[ ]~~ does not exist
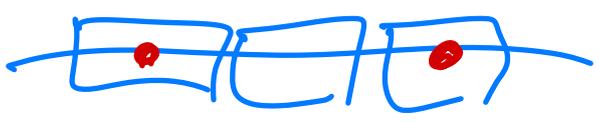
terms: linkage disequilibrium (LD)
"sites are in absolute linkage disequilib."
100% linked

"sites are in linkage equilibrium"
0% linked

# Problem:

Given an $n \times k$ haplotype matrix (binary) with $p$ "cases" and $q$ "healthy" individuals $(p + q = n)$, find segregating sites that correlate with the case/disease pattern.

# Classical solutions

apply statistical tests, site by site.
- Fisher's exact test
- chi-square test

$\Rightarrow$ slow, especially for multiple sites

# Another method (Blossoc):   [infinite sites assumption]

first, focus on a single haplotype block (no crossing-over) = recombination



perfect phylogeny

**Area 1** | **Area 2** | **Tree 1**

| 1 | A T C |
| 2 | A T C |
| 3 | A T T |
| 4 | C T T |
| 5 | C T C |
| 6 | C G T |
| 7 | C G T |

D D D  H H H H

① has perfect phylogeny $P$
② $P$ correlates with the genealogy of the disease

**Algorithm:**

1. Look at each SNV site from left to right.

2. From this extend to the right without violating the 4-gamete test.

3. For each such region, build a perfect phylogeny (which must exist).

4. Compare all these trees with the case-disease pattern.

**Tree 2**

H D  H D  H H D

↓ ↓ ↓ ↓ ↓ --

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