# Estimating Local Ancestry in Admixed Populations (LAMP)

QIAN ZHANG
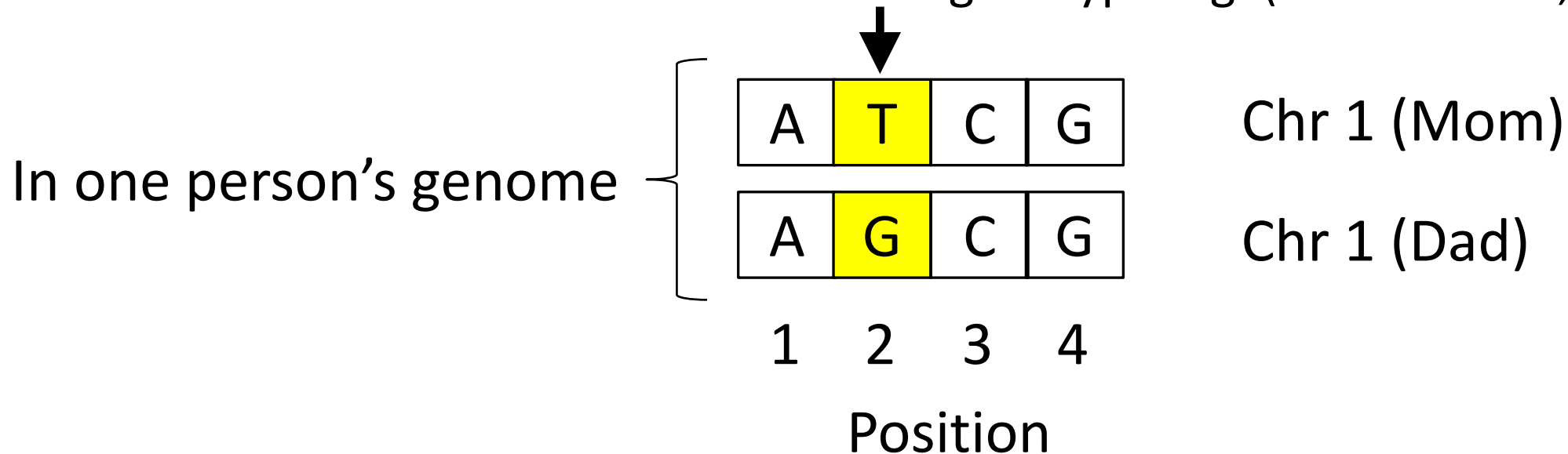
572 INTRO TALK

# Outline

- Human genome = 23 Chromosomes (Chr 1 – Chr 23)
- Each Chr comes as a pair (Mom, Dad)
- Each Chr has positions, where there are variants A, T, C, or G
- SNP: A position where people have different variants
- Most SNPs have 2 variants: Minor Allele, Major Allele
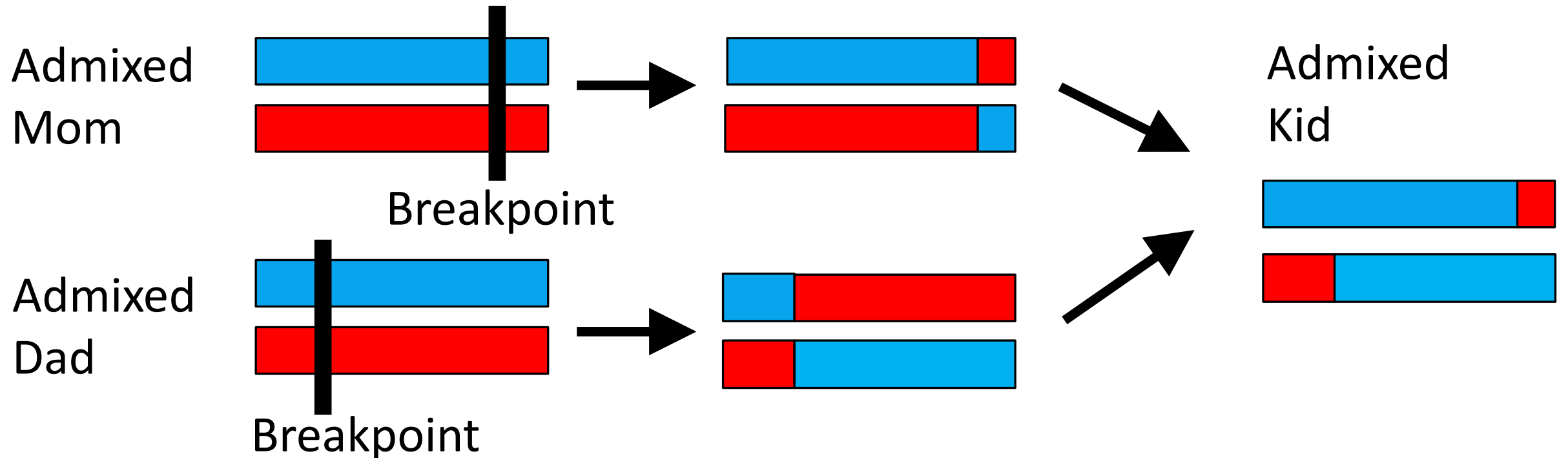- Genotype = Pair of alleles at a SNP

SNP with genotype e.g. (Minor Allele T, Major Allele G)

In one person's genome

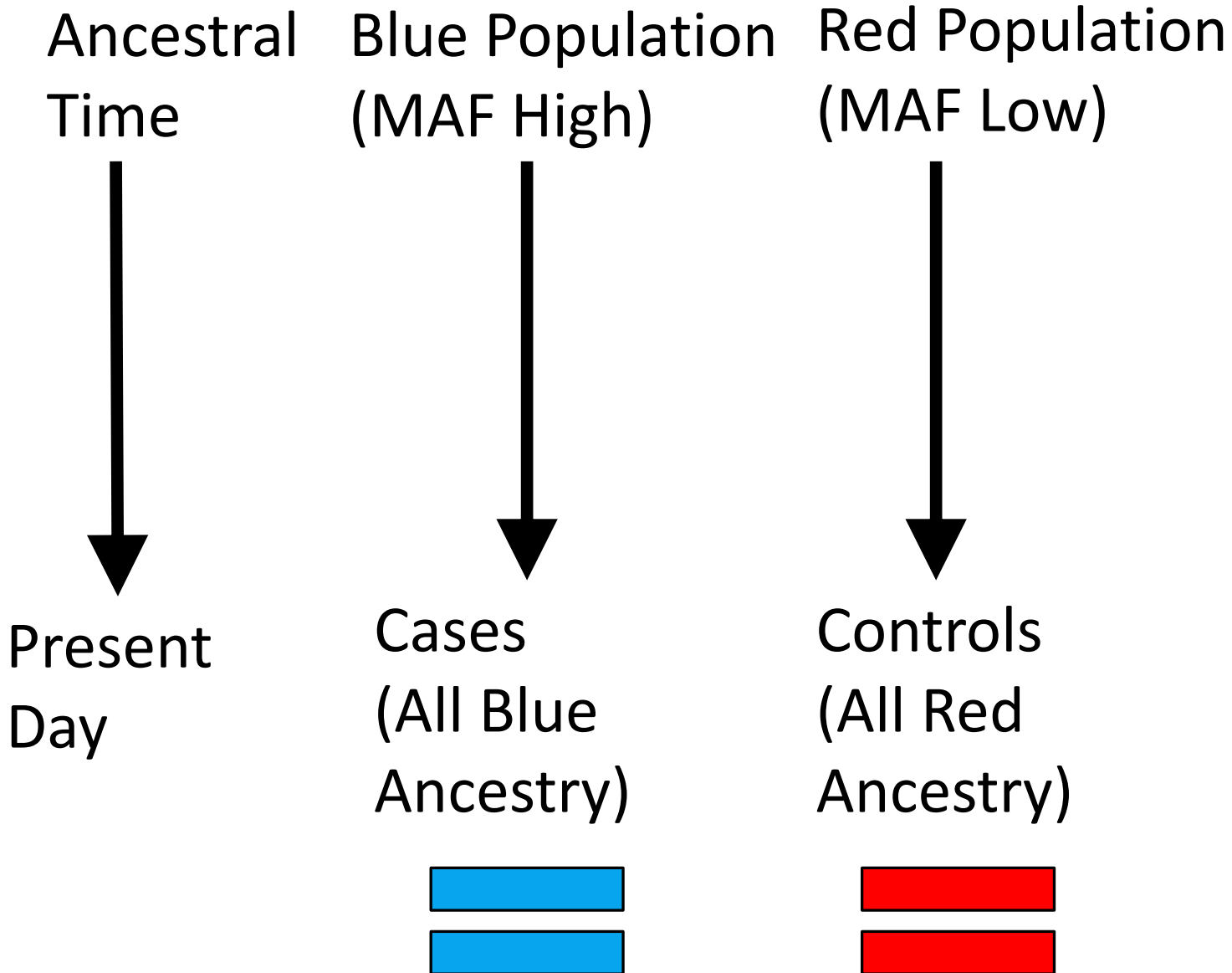| A | T | C | G | Chr 1 (Mom) |
| A | G | C | G | Chr 1 (Dad) |

1  2  3  4

Position

Genome-wide Association Study

- At each SNP, test H0:
  The minor allele is not associated with a disease-related outcome
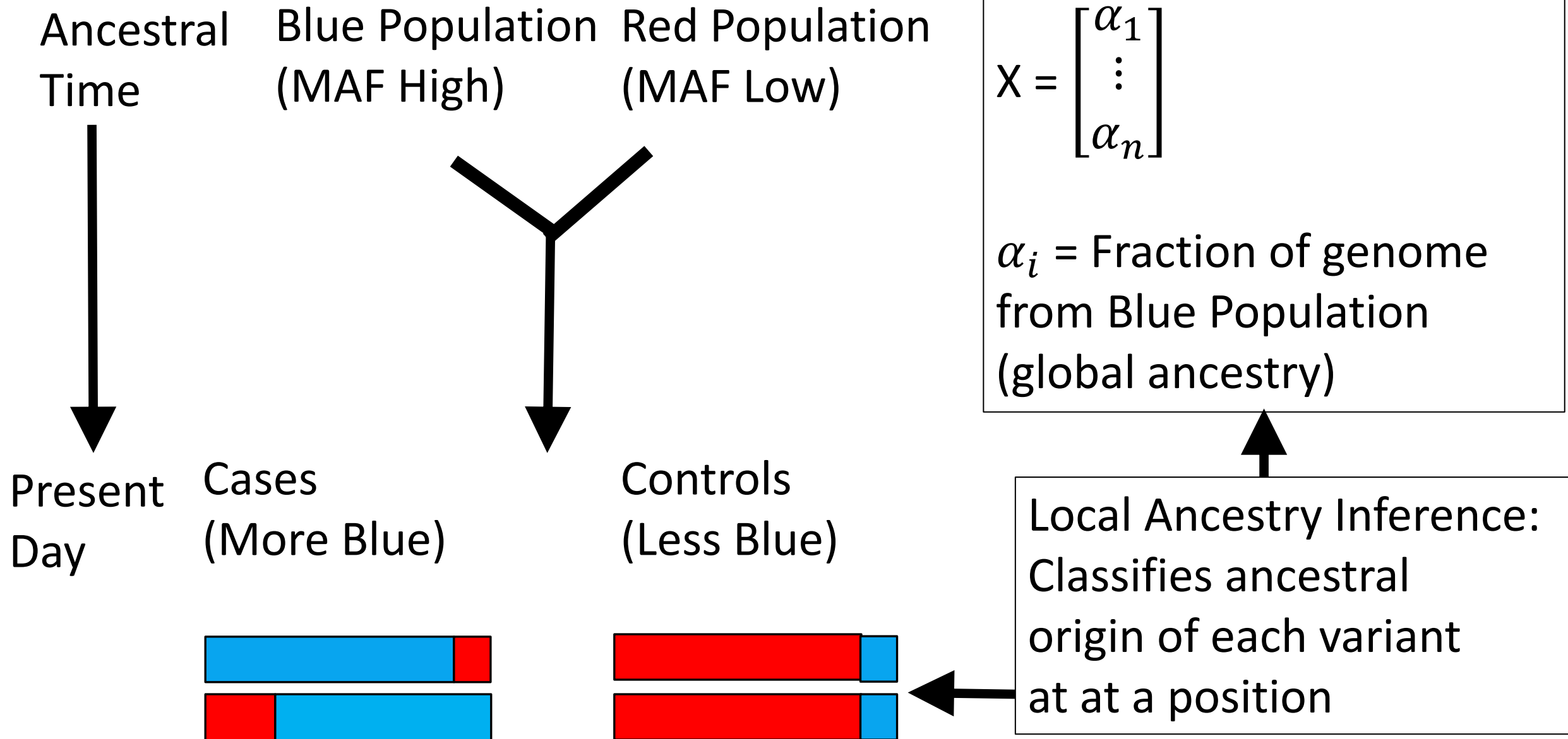  (e.g. case/control status, LDL level)

$$\text{logit}(1[\text{case}]) = B_0 + B_1[\text{\# copies of minor allele}] + B_2 X + \epsilon$$

- Is the minor allele more frequent in cases than in controls?
- X = Population history is a confounder because
Different Population history $\Rightarrow$ Different minor allele frequency (MAF)

# Mapping by Admixture Linkage Disequilibrium
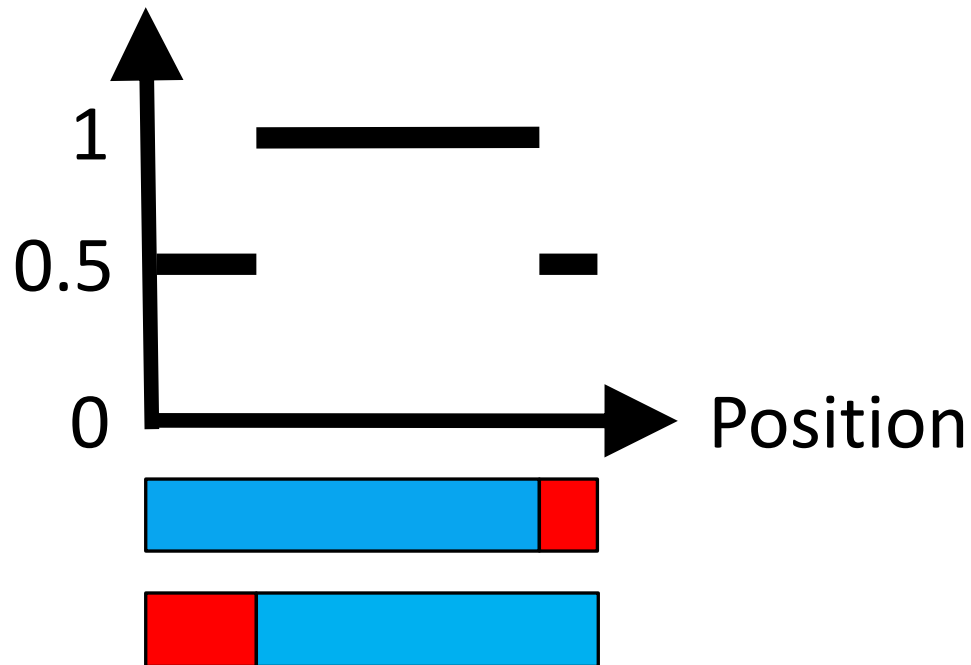- Disease-risk allele more frequent in Blue Population than Red Population
- Admixed cases



Fraction Blue Ancestry

1

0.5

0

Position

Average step functions over cases

Fraction Blue Ancestry

Position

Disease risk allele here
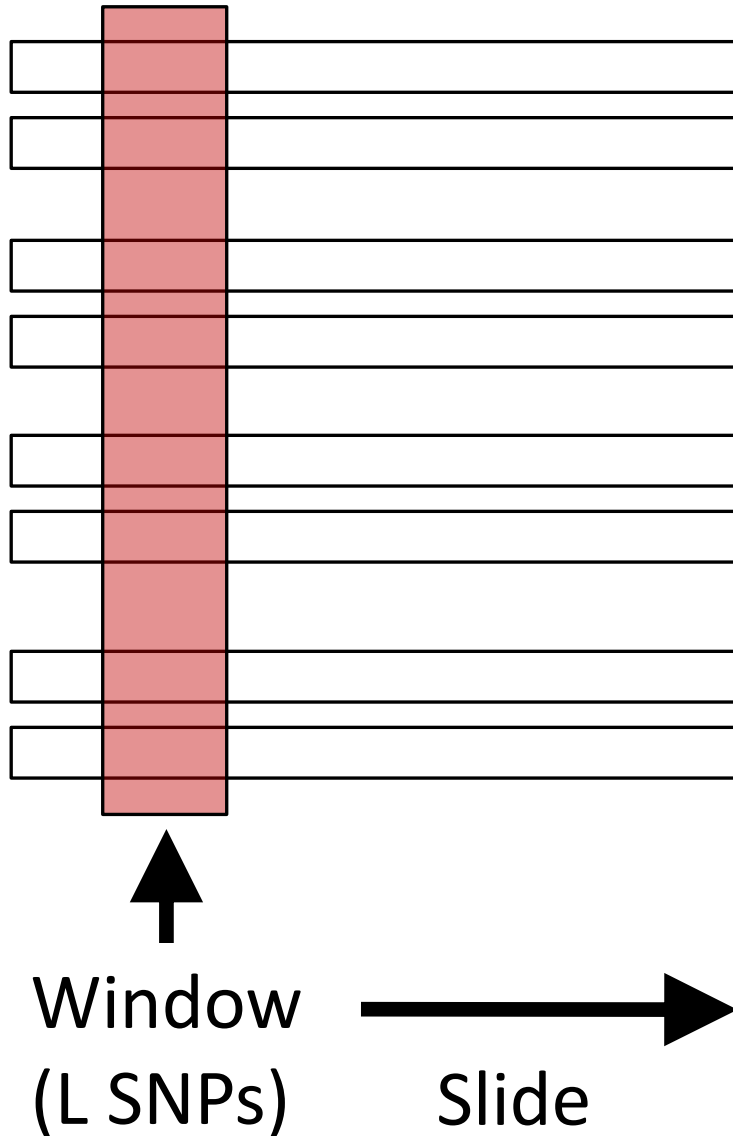
Summary of Uses of Local Ancestry Inference
(the classification of variants to ancestral populations):


1) Calculate global ancestry → Control for population history in genome-wide association studies


2) Perform mapping by admixture linkage disequilibrium

# Outline

1. Scientific and Statistical Motivation
2. **How LAMP Works**
3. LAMP vs Other Methods
4. What's Next?

# LAMP: Overview



Window (L SNPs)    Slide

0) Remove 1 of each pair of correlated SNPs
→ Work with smaller set of independent SNPs
1) Consider a window of L SNPs:
Assume L so small that nobody has a breakpoint in the window, so in the window, a Chr is all from one ancestry
2) MAXVAR: Get initial ancestry estimates in the window

$$\theta(i) = (\theta_1(i), \theta_2(i)) \in \{1, \ldots, K\} \text{ x } \{1, \ldots, K\}$$

From Mom    From Dad    Number of Ancestral Populations

for individuals i = 1, ..., m
3) Iterative clustering algorithm: Get final ancestry estimates $\theta(i)$ in the window
4) Slide window over to get a new window overlapping the old one. Repeat 1-3.
5) Majority vote (over windows) to call a SNP's ancestry

# LAMP: Overview



SNP from Blue Population

0) Remove 1 of each pair of correlated SNPs
→ Work with smaller set of independent SNPs
1)  Consider a window of L SNPs:
Assume L so small that nobody has a breakpoint in the window, so in the window, a Chr is all from one ancestry
2) MAXVAR: Get initial ancestry estimates in the window

$$\theta(i) = (\theta_1(i), \theta_2(i)) \in \{1, \ldots, K\} \text{ x } \{1, \ldots, K\}$$

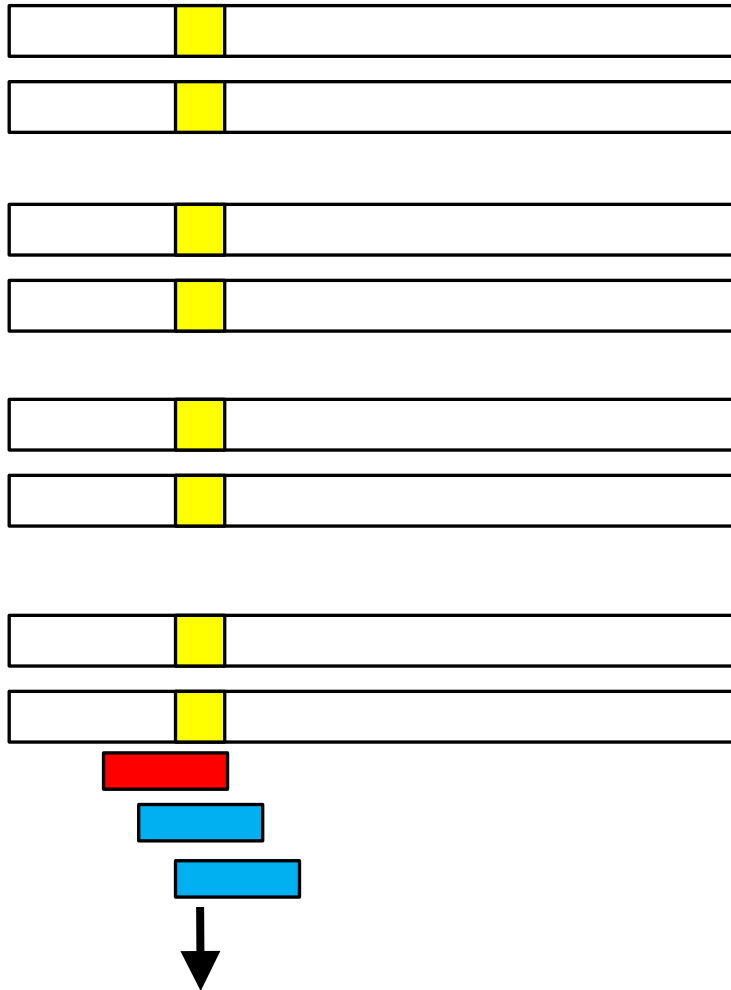↑ From Mom    ↑ From Dad    ↑ Number of Ancestral Populations

for individuals i = 1, …, m
3) Iterative clustering algorithm: Get final ancestry estimates $\theta(i)$ in the window
4) Slide window over to get a new window overlapping the old one. Repeat 1-3.
5) Majority vote (over windows) to call a SNP's ancestry

2) MAXVAR: Get initial ancestry estimates in the window

Define a similarity score between individuals i1 and i2

$$S(i_1, i_2) = \frac{\sum_{j=1}^{n}(G_{i1,j} - u_j)(G_{i2,j} - u_j)}{\sigma_j^2}$$

which is like a sample correlation($G_{i1}$, $G_{i2}$).

$G_{i,j}$ is the Genotype (# of minor alleles) at position j in individual i.

$u_j = \sum_{i=1}^{m}\frac{G_{ij}}{m}$ is the Genotype at position j averaged over m individuals i.

$\sigma_j^2 = \sum_{i=1}^{m}\frac{(G_{ij}-u_j)^2}{m}$ is the variance in Genotype

# 2) MAXVAR: Get initial ancestry estimates in the window

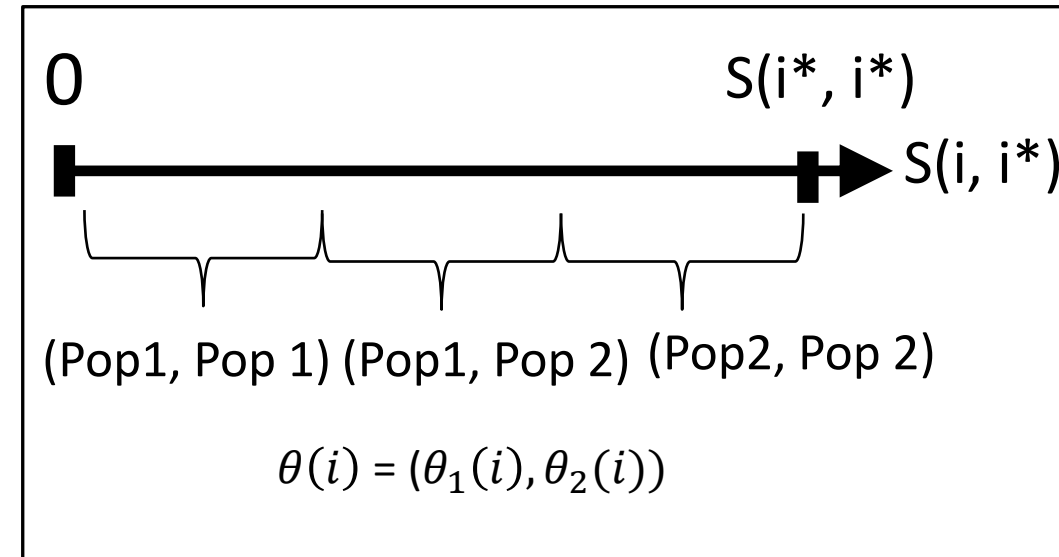A) For each individual i, $Var(i) = \sum_{\{i':i'\neq i\}} S(i,i')^2$
measures how genetically similar i is to everyone else (the other m - 1 people).

B) Find the person i* who is most genetically similar to everyone else:
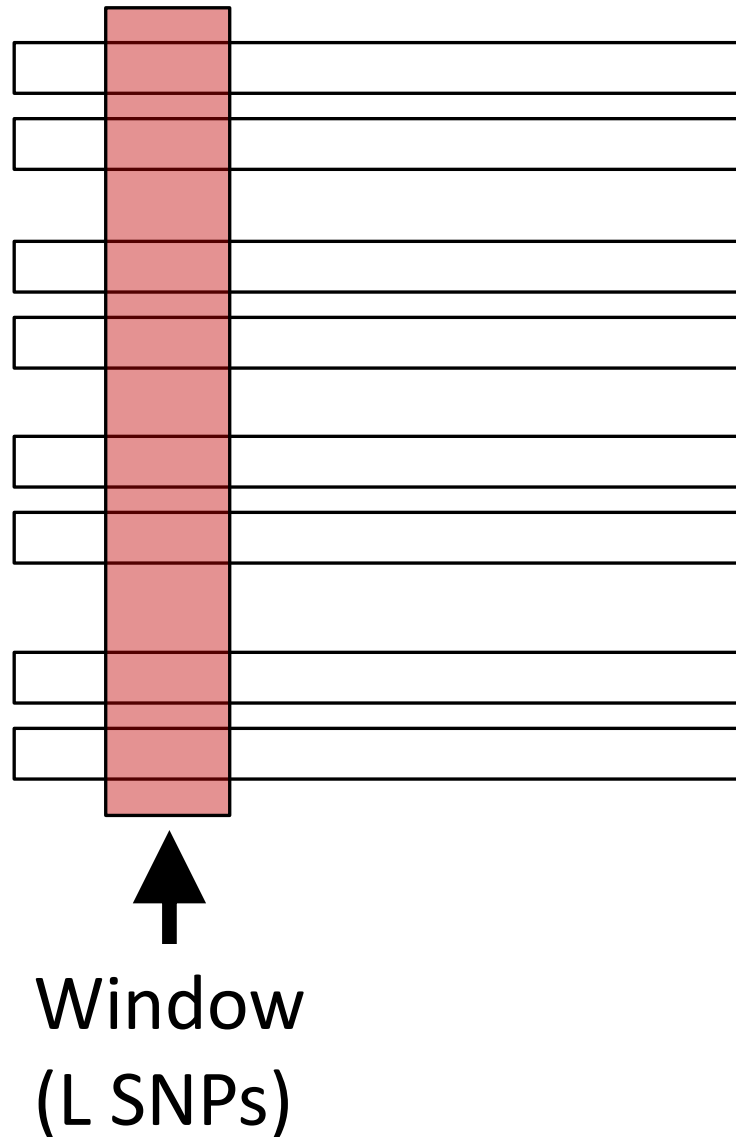i* = $\mathrm{argmax_i} \left( Var(i) \right)$.

C) Assign initial ancestry estimates of $\theta(i)$ based on how similar i is to i*:
- Assume
  - K = 2 ancestral populations, Pop 1 and Pop 2
  - i* has ancestry (Pop 2, Pop 2)
  - Know $\alpha$: fraction of variants from Pop 1
- For each individual i, find how similar i is to i*: S(i, i*)
  - Lowest $(1-\alpha)^2 n$ scores $S(i, i^*)$ → (Pop 1, Pop 1)
  - Highest $\alpha^2 n$ scores $S(i, i^*)$ → (Pop 2, Pop 2)
  - Everyone else → (Pop 1, Pop 2)



0        S(i*, i*)

S(i, i*)

(Pop1, Pop 1) (Pop1, Pop 2) (Pop2, Pop 2)

$\theta(i) = (\theta_1(i), \theta_2(i))$

# 3) Iterative clustering algorithm: Get final ancestry estimates $\theta(i)$



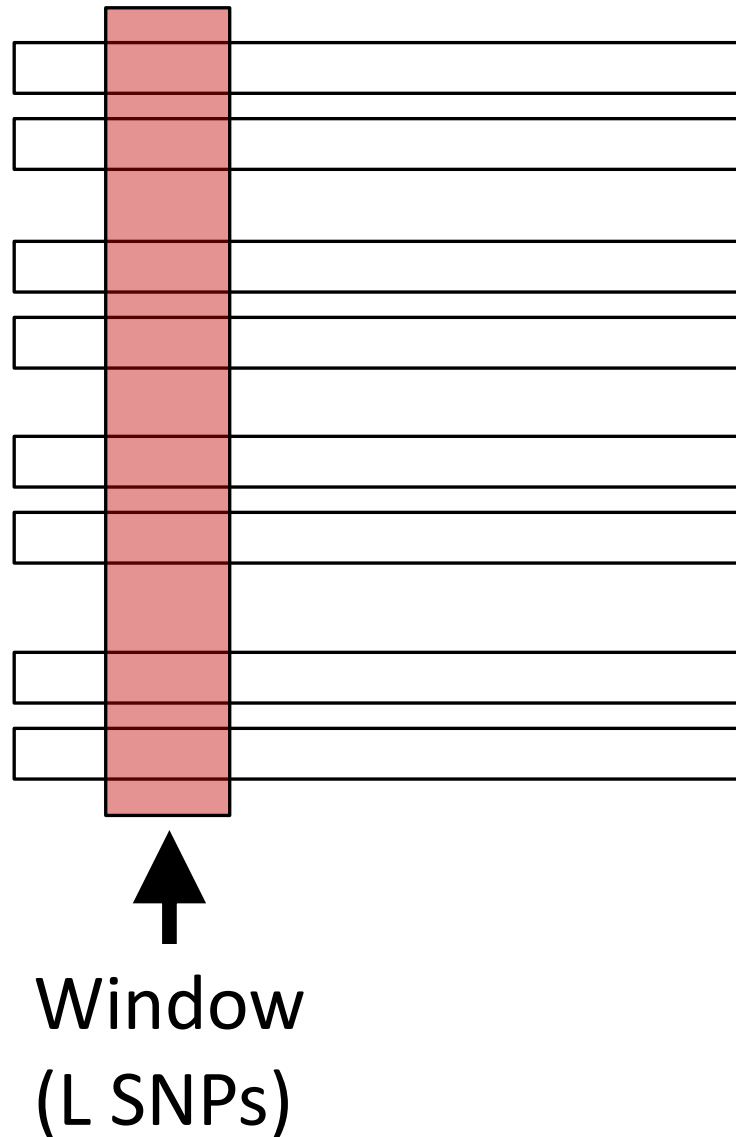Goal: For each individual i, finally assign to each Chr in the window an ancestry
$$\theta(i) = (\theta_1(i), \theta_2(i)) \in \{1, \ldots, K\} \text{ x } \{1, \ldots, K\}$$

From Mom    From Dad

Number of Ancestral Populations

Window (L SNPs)

# 3) Iterative clustering algorithm: Get final ancestry estimates $\theta(i)$



If know:

- MAFs $\vec{f_1}, \dots, \vec{f_K}$ of K ancestral populations
- Genotypes $G_1, \dots, G_m$ of m individuals

Then estimate individual i's ancestries
$\theta(i) = (\theta_1(i), \theta_2(i)) \in \{1, \dots, K\} \times \{1, \dots, K\}$ in the window as
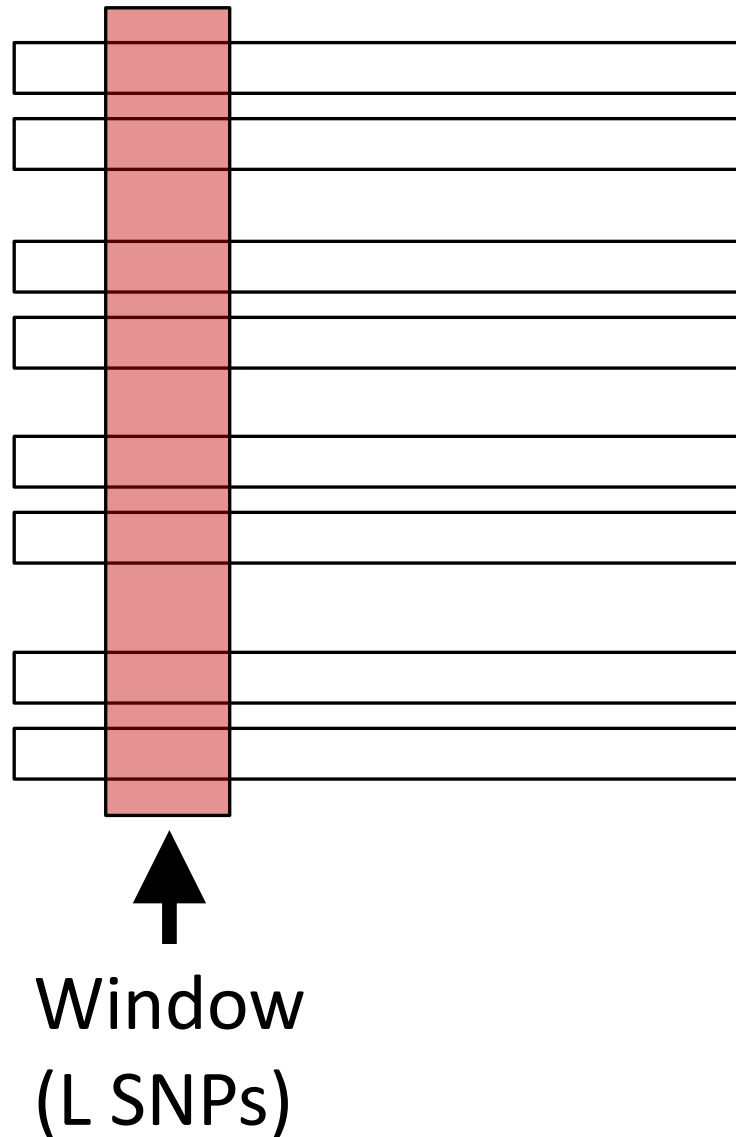
From Mom    From Dad    Number of Ancestral Populations

$$\text{argmax}_{\theta(\text{i})} \, P(\theta(i)) \, | \vec{f_1}, \dots, \vec{f_K}, G_i) \qquad (1)$$

Window (L SNPs)

# 3) Iterative clustering algorithm: Get final ancestry estimates $\theta(i)$



If know:

- Individual i's ancestries

$$\theta(i) = (\theta_1(i), \theta_2(i)) \in \{1, \ldots, K\} \times \{1, \ldots, K\}$$

From Mom  From Dad  Number of Ancestral Populations

- Genotypes $G_1, \ldots, G_m$ of m individuals

Then estimate ancestral MAFs $\vec{f_1}, \ldots, \vec{f_K}$ as

$$\text{argmax}_{\vec{f_1}, \ldots, \vec{f_K}} \prod_{i=1}^{m} P(G_i | \vec{f_1}, \ldots, \vec{f_K}, \theta(i)) \qquad (2)$$

Start $\theta(i)$ from MAXVAR

Window
(L SNPs)

# 3) Iterative clustering algorithm: Get final ancestry estimates $\theta(i)$



Iterate equations (1) and (2) to finally estimate individual i's ancestries

$$\theta(i) = (\theta_1(i), \theta_2(i)) \in \{1, \dots, K\} \text{ x } \{1, \dots, K\}$$

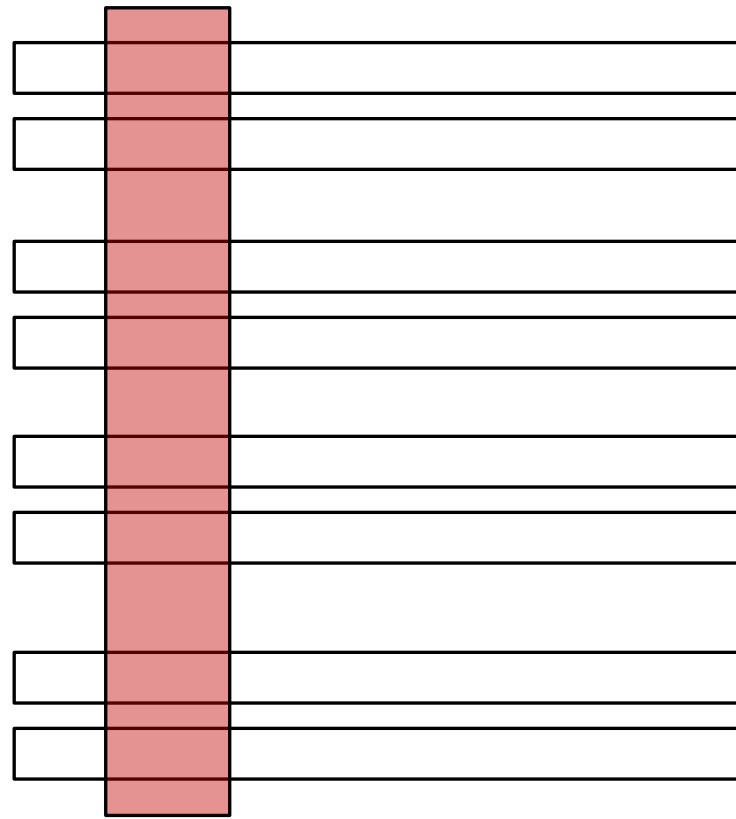From Mom    From Dad    Number of Ancestral Populations

Window
(L SNPs)

# Outline

Methods that detect global (whole-genome ancestry proportion) or local ancestry (ancestry at particular positions):

- STRUCTURE (> 11,000 citations)
- SABER
- **LAMP (LAMP-ANC)**
- HAPMIX, RFMIX, FAST STRUCTURE, and others!

Methods differ in:
- Estimating global vs local ancestry
- Can handle = 2 or >= 2 ancestral populations
- Input parameters to program
(SABER, LAMP-ANC assume ancestral MAFs known)
- Speed, Accuracy
- Whether they crash when given a lot of data (# SNPs, # people)

Paper in a Nutshell
- Simulate admixed genetic data from ancestral populations (African + European, European + Japanese, Japanese + Chinese)
- Apply STRUCTURE, SABER, LAMP (LAMP-ANC) to simulated data
- LAMP is faster, more accurate than STRUCTURE and SABER except with Japanese + Chinese
- LAMP-ANC is faster, more accurate than STRUCTURE and SABER in all cases
    - Accuracy measured as fraction of variants with correctly classified ancestry
        - Measure makes sense for SABER and LAMP (local ancestry)
        - Measure does not make sense for STRUCTURE (global ancestry)

**Table 1. A Summary of the Comparison between LAMP, LAMP-ANC, SABER, and STRUCTURE**

| Dataset | Distance | LAMP | LAMP-ANC | SABER | STRUCTURE |
|---|---|---|---|---|---|
| YRI-CEU | 0.055 | 0.94 | 0.95 | 0.87 | 0.84 |
| CEU-JPT | 0.036 | 0.87 | 0.93 | 0.82 | 0.47 |
| JPT-CHB | 0.0045 | 0.48 | 0.72 | 0.68 | 0.40 |
| Time (s) | | 394 | 246 | 7681 | $2.57 \times 10^5$ |
| Number of SNPs | | 38,864 | 38,864 | 4000 | 4000 |

The accuracy across all positions on chromosome 1 is shown for the three admixed populations. The distance between the admixing population (measured by the mean squared distance between the allele frequency vectors) is also shown, indicating the difficulty in separating alleles from the populations. The time taken to run each of the methods is shown. LAMP and LAMP-

Maybe why LAMP is faster than STRUCTURE:
- STRUCTURE uses MCMC (computationally intensive) to generate draws from a posterior distribution
- LAMP does not use MCMC

Expect SABER to be more accurate than LAMP:
- SABER extends a hidden Markov model to account for linkage disequilibrium (correlation between closely positioned SNPs)
- LAMP removes 1 of each pair of highly correlated SNPs, uses reduced set of independent SNPs

Maybe why LAMP is more accurate than SABER:
- All the SNPs have only a little more info than the set of independent SNPs
- SABER estimates more parameters, so algorithm might not converge to global optimum

# Outline

1. Scientific and Statistical Motivation

2. LAMP (LAMP-ANC) and Others

3. How LAMP Works

4. **What's Next?**

# What's Next?

- Choice of window length

- Why MAXVAR is a reasonable thing to do

- Simulate data and run 3 methods to compare accuracy


OLGA (Omics in Latinos Genetic Analysis) at UW

will use local ancestry inference.