

# Leveraging Local Ancestry with Admixture Mapping

Instructors: Andrey Ziyatdinov and Timothy Thornton

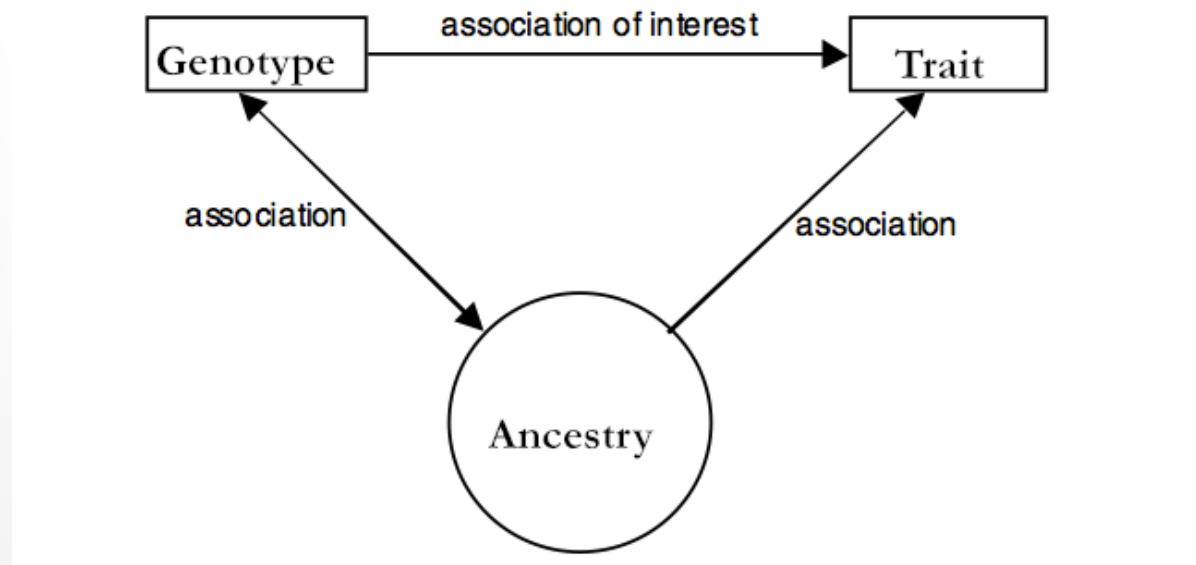
National Autonomous University of Mexico  
April 27, 2023

# Genetic Association Studies in Multi-Ethnic Populations

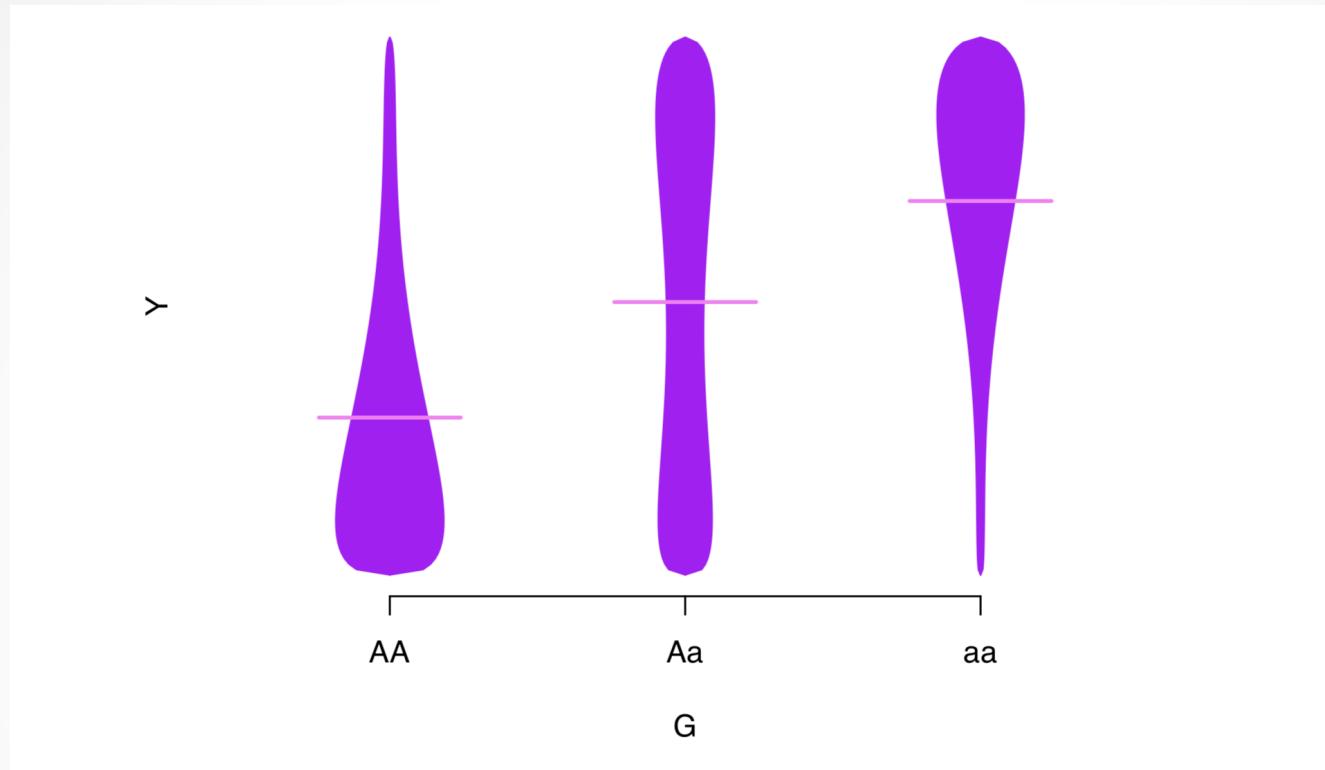
- There remain significant challenges with complex trait mapping in multi-ethnic populations
- Two well known challenges are:
  - **Heterogeneous genetic ancestry and environmental backgrounds among sampled individuals**
  - **Correlated genotype and phenotype data among relatives, known and/or cryptic, in the sample**

# Confounding is a serious concern for genetic studies in multi-ethnic populations

- Ethnic groups (and subgroups) have often share distinct dietary habits and other lifestyle characteristics that result in traits of interest having **different distributions** that are correlated with genetic ancestry and/or ethnicity.



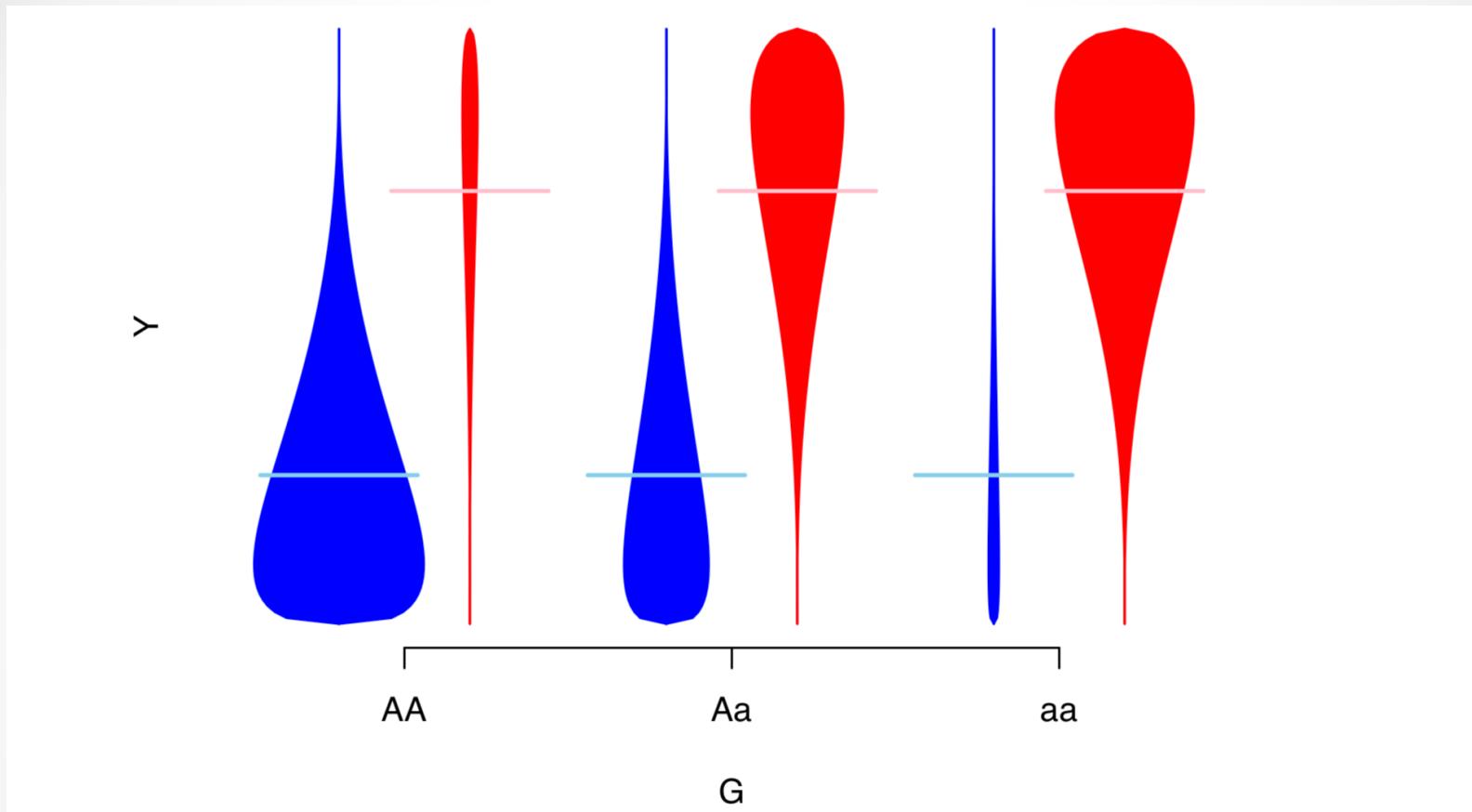
# Confounding: multi-ethnic studies



$$\mathbf{Y} = \beta_0 \mathbf{1} + \mathbf{g}_s \beta_1 + \epsilon$$

- $\mathbf{Y}$  is a vector of phenotypes
- $\mathbf{g}_s$  is an additive genotype count vector for a SNP  $s$ , where each entry corresponds to the number of reference alleles (A) an individual has, e.g., 0, 1, or 2;
- $\epsilon \sim \mathbf{N}(\mathbf{0}, \sigma_e^2 \mathbf{I})$

# Confounding: multi-ethnic studies



- The relationship between phenotype vector (**Y**) and genotype vector (**G**) looks much less interesting when broken down and assessed within ancestry groups;
- Well known that the genetic ancestry is a confounder in multi-ethnic studies that can lead to spurious associations

# PCs in Regression Models

- In practice, multi-ethnic cohort studies will not have discrete or a fixed number of ancestry groups.
- Eigenvectors (PCs) are often used as surrogates for ancestry (or population structure).
- To protect against spurious association due to genetic ancestry confounding, the top PCs are often included as **fixed effects** in regression models used for assessing genotype/phenotype associations;

$$E(\mathbf{Y}) = \beta_0 \mathbf{1} + \mathbf{g}_s \beta_1 + \gamma_1 PC_1 + \gamma_2 PC_2 + \gamma_3 PC_3 + \dots$$

i.e., regression model adjusting for PC1, PC2, PC3 etc.  
(Logistic, Cox regression can be adjusted similarly)

- Among people with the same ancestry (i.e. the same PCs)  $\beta_1$  gives us the difference in mean phenotype, per 1-unit difference in  $\mathbf{g}_s$
- If the effect of  $\mathbf{g}_s$  differs by PCs,  $\beta_1$  provides a (sensible) average of these genetic effects

# Can genetic ancestry be leveraged?

- Genetic association mapping methods generally treat genetic ancestry as a potential confounder
- In multi-ethnic populations that have admixed ancestry derived from multiple continents, such as African Americans and Hispanic/Latino populations, the frequency of causal genetic variants may differ substantially in the underlying ancestral populations
- Admixed individuals with causal variants may be more likely to have inherited the variants from the ancestral population for which the variants are most prevalent
-

# Can genetic ancestry be leveraged?

- Mapping by admixture linkage disequilibrium, or **admixture mapping**, leverages heterogeneity in genetic ancestry by identifying loci that
  - (1) have unusual deviations in local ancestry, relative to what would be expected based on genome-wide ancestry

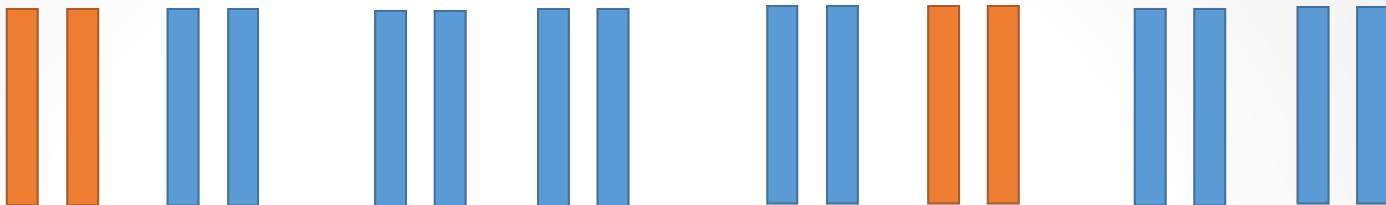
and

- (2) and that are significantly associated with a trait
- By leveraging genetic ancestry, regions of the genome harboring genetic variants that differ in frequency across the ancestral populations and that drive the observed phenotypic differences
- can be identified



# Admixed Populations

Generation 0:  
2 Ancestral  
Populations  
(orange, blue)



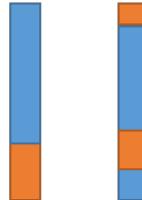
Generation 1



Generation 2



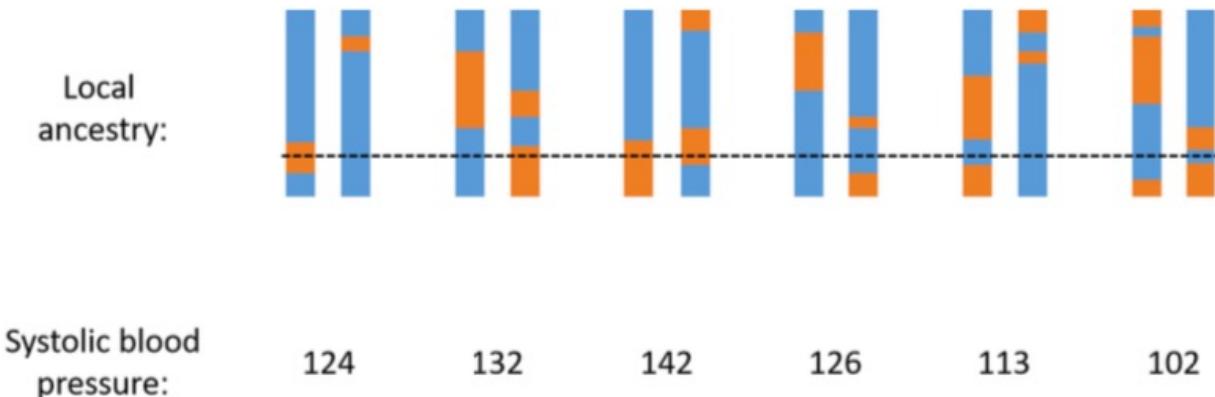
Generation 3



# Admixture Mapping

**Goal:** Identify genetic variants associated with trait of interest

**How?** Look for associations between the trait and local ancestry



- Previously proposed regression approach for admixture with two ancestral populations:

$$\mathbf{Y} = \beta_0 \mathbf{1} + \mathbf{a}_{\ell p} \beta_{\ell p} + \mathbf{V} \gamma + \epsilon$$

- $\mathbf{a}_{\ell p}$  is a local ancestry vector with entries consisting of the number of alleles (0,1,2) inherited from ancestral population  $p$  at locus  $\ell$  for each individual
- $\mathbf{V}$  is a vector of covariates for genome-wide average ancestry (e.g., admixture proportions or PCs)

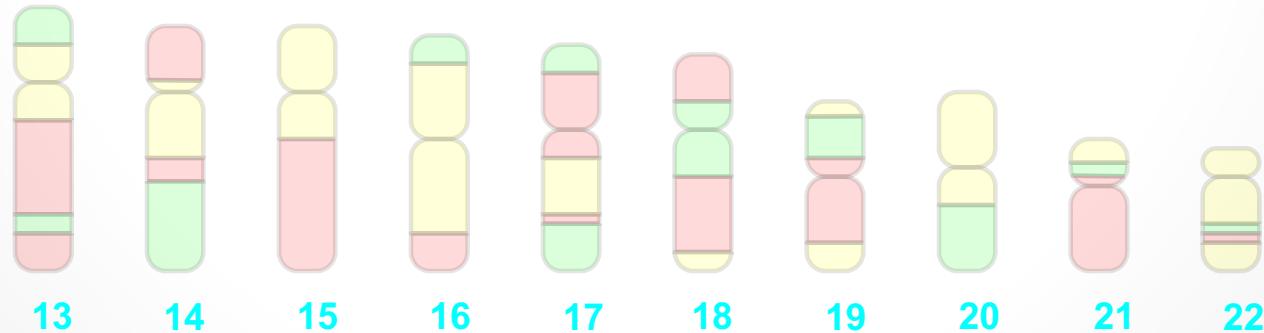
# Admixture Mapping with Multiple Ancestral Populations

- Proposed admixture mapping model:

$$\mathbf{Y} = \beta_0 \mathbf{1} + \mathbf{a}_{\ell p} \beta_{\ell p} + \mathbf{a}_{\ell p'} \beta_{\ell p'} + \mathbf{V} \gamma + \mathbf{r} + \epsilon$$

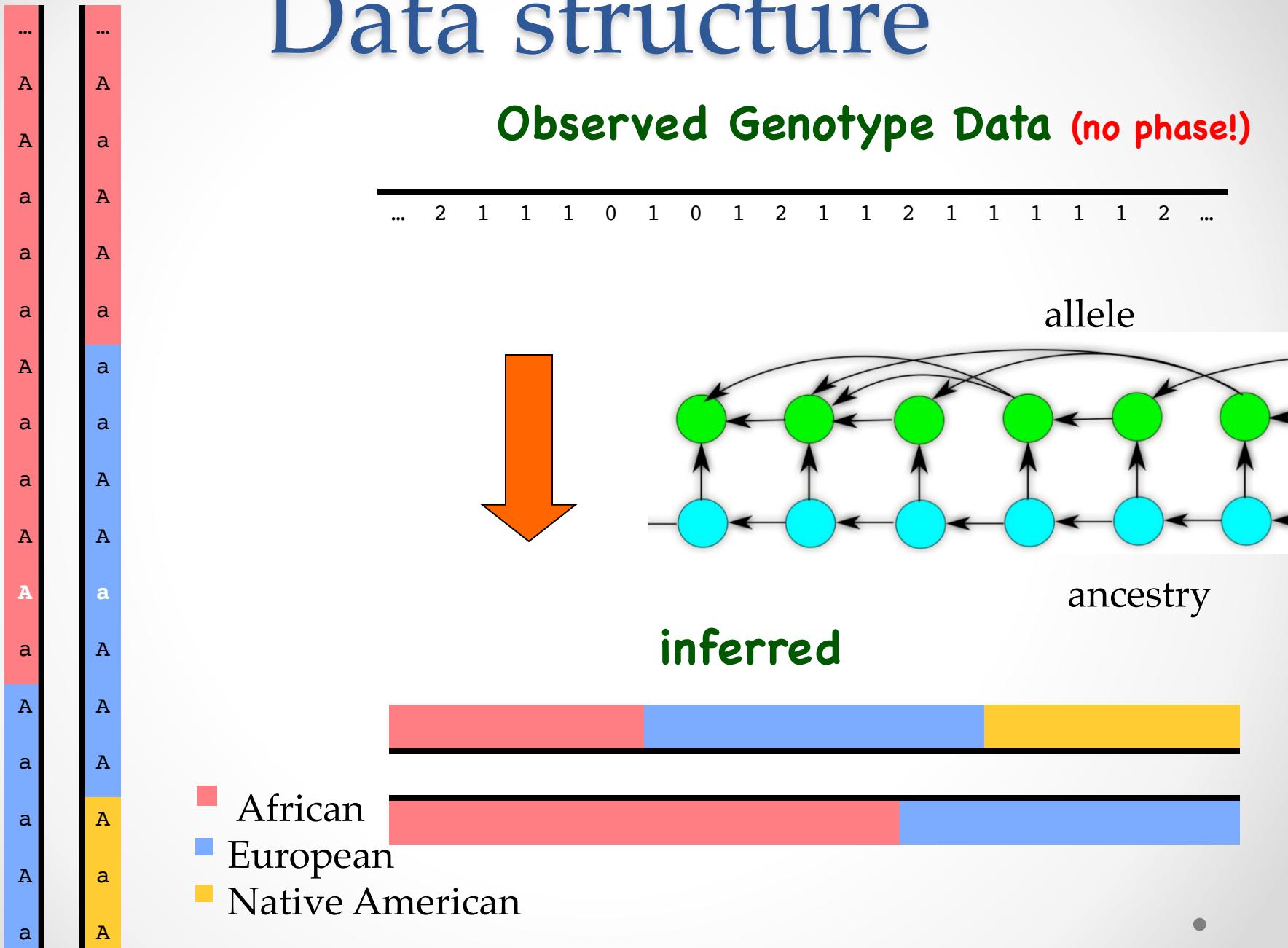
- To test for an association with local ancestry at a locus  $\ell$ , we propose using a score test to test the null hypothesis

$$H_0 : \beta_{\ell p} = \beta_{\ell p'} = 0$$



# Data structure

# Observed Genotype Data (no phase!)



# RFMix Method for Local Ancestry Inference

## ARTICLE

---

### RFMix: A Discriminative Modeling Approach for Rapid and Robust Local-Ancestry Inference

Brian K. Maples,<sup>1,2</sup> Simon Gravel,<sup>1,3</sup> Eimear E. Kenny,<sup>1,4,5,6,7,8</sup> and Carlos D. Bustamante<sup>1,8,\*</sup>

- RFMix (Maples et al.; AJHG 2013) continues to be one of the most widely used software programs as it can accurately infer local ancestry for admixed samples derived from multiple ancestral populations

# Application to HCHS/SOL

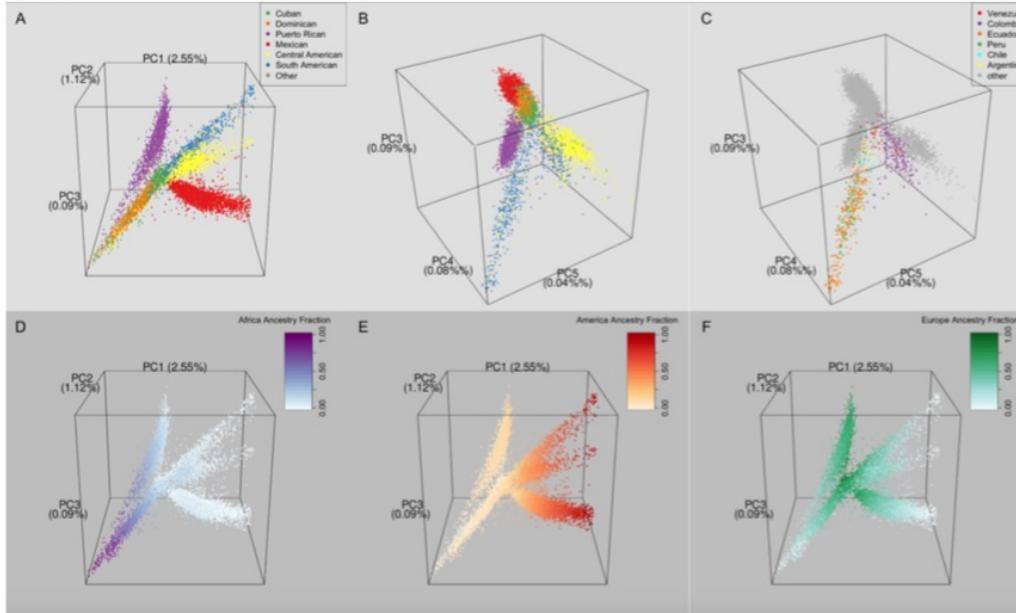
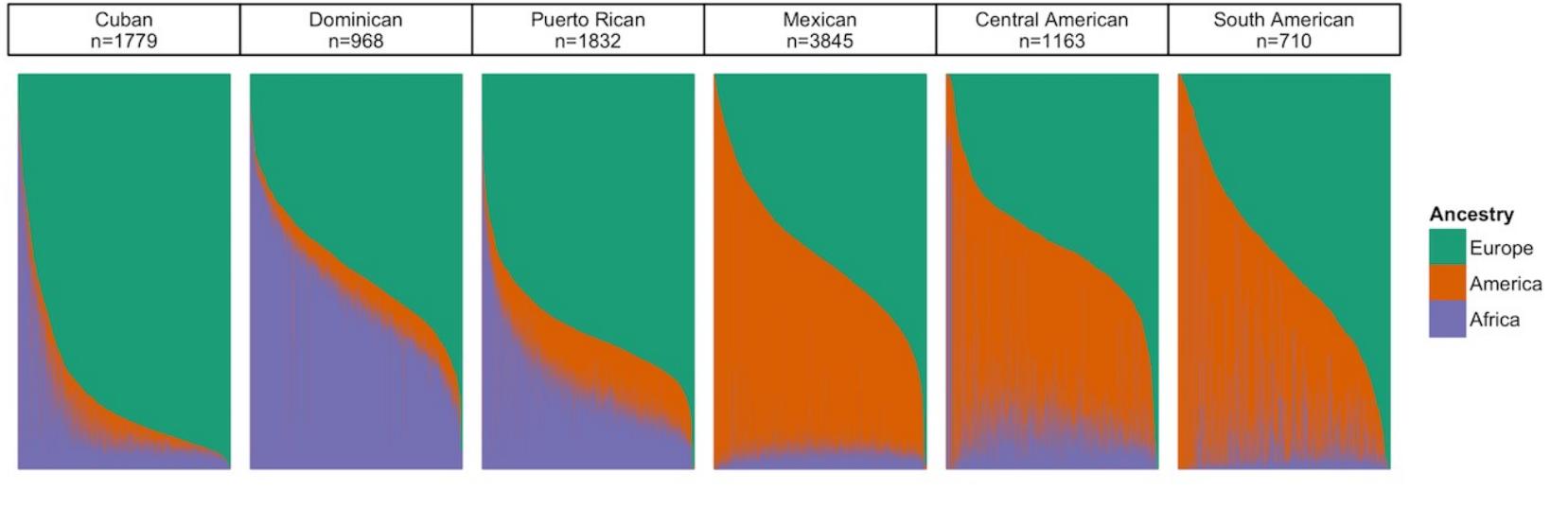
Hispanic Community Health Study/Study of  
Latinos

- Study Subjects
  - ~13,000 participants who self-identified as Hispanic/Latino
  - randomly selected households from four US cities
- Local ancestry estimation
  - performed at 236,456 SNPs with RFMix



# HCHS/SOL Ancestry

A



ARTICLE

## Genetic Diversity and Association Studies in US Hispanic/Latino Populations: Applications in the Hispanic Community Health Study/Study of Latinos

Matthew P. Conomos,<sup>1,14,\*</sup> Cecelia A. Laurie,<sup>1,14</sup> Adrienne M. Stilp,<sup>1,14</sup> Stephanie M. Gogarten,<sup>1,14</sup> Caitlin P. McHugh,<sup>1</sup> Sarah C. Nelson,<sup>1</sup> Tamar Sofer,<sup>1</sup> Lindsay Fernández-Rhodes,<sup>2</sup> Anne E. Justice,<sup>2</sup> Mariaelisa Graff,<sup>2</sup> Kristin L. Young,<sup>2</sup> Amanda A. Seyerle,<sup>2</sup> Christy L. Avery,<sup>2</sup> Kent D. Taylor,<sup>3</sup> Jerome I. Rotter,<sup>3</sup> Gregory A. Talavera,<sup>4</sup> Martha L. Daviglus,<sup>5</sup> Sylvia Wassertheil-Smoller,<sup>6</sup> Neil Schneiderman,<sup>7</sup> Gerardo Heiss,<sup>2</sup> Robert C. Kaplan,<sup>6</sup> Nora Franceschini,<sup>2</sup> Alex P. Reiner,<sup>8</sup> John R. Shaffer,<sup>9</sup> R. Graham Barr,<sup>10</sup> Kathleen F. Kerr,<sup>11</sup> Sharon R. Browning,<sup>11</sup> Brian L. Browning,<sup>11</sup> Bruce S. Weir,<sup>1</sup> M. Larissa Avilés-Santa,<sup>12</sup> George J. Papanicolaou,<sup>12</sup> Thomas Lumley,<sup>13</sup> Adam A. Szpiro,<sup>1</sup> Kari E. North,<sup>2</sup> Ken Rice,<sup>1</sup> Timothy A. Thornton,<sup>1</sup> and Cathy C. Laurie<sup>1,\*</sup>

# HCHS/SOL: Admixture Mapping of Kidney Traits

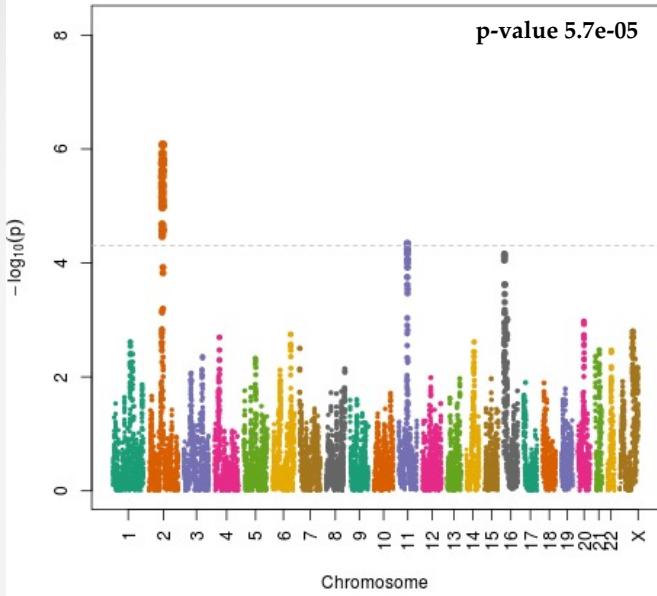
## **Genome-Wide Admixture Mapping of Estimated Glomerular Filtration Rate and Chronic Kidney Disease Identifies European and African Ancestry-of-Origin Loci in Hispanic and Latino Individuals in the United States**

Andrea R.V.R. Horimoto ,<sup>1</sup> Diane Xue,<sup>2</sup> Jianwen Cai ,<sup>3</sup> James P. Lash ,<sup>4</sup>  
Martha L. Daviglus,<sup>5</sup> Nora Franceschini,<sup>6</sup> and Timothy A. Thornton<sup>1,7</sup>

- Horimoto et al. (JASN; 2022)
- Analysis model
  - linear mixed model for both association and admixture mapping
  - adjust for PCs1-5, genetic background group, sex, age, center, smoking status
- ◦ admixture mapping significance threshold: 6e-5 •

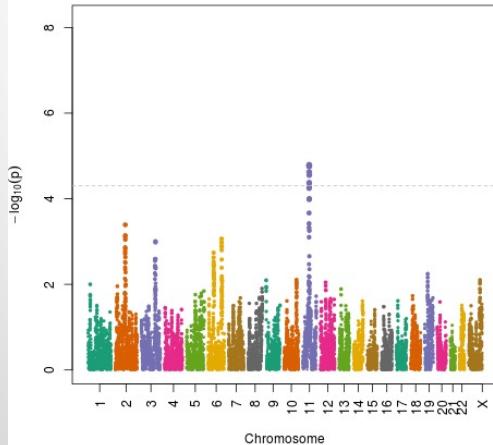
# HCHS/SOL: Albumin-to-creatinine ratio Admixture Mapping

## Joint Test

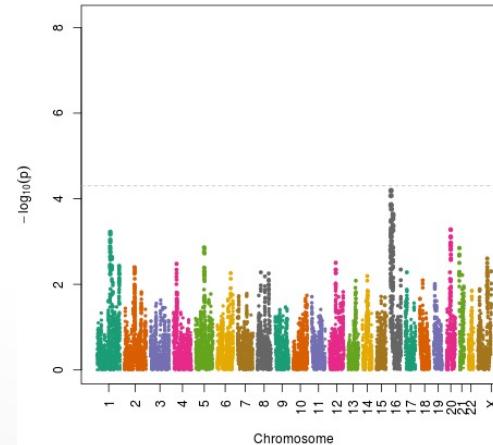


chr	# SNPs	pos	p-value	gene
2	1 (rs2139376)	112143413	8.4e-07	MIR4435-2HG (2q13)
11	13	69533129 - 69579070	4.5e-05	intergenic region
16	8	4964076 - 4999392	7.0e-05	PPL (16p13.3)

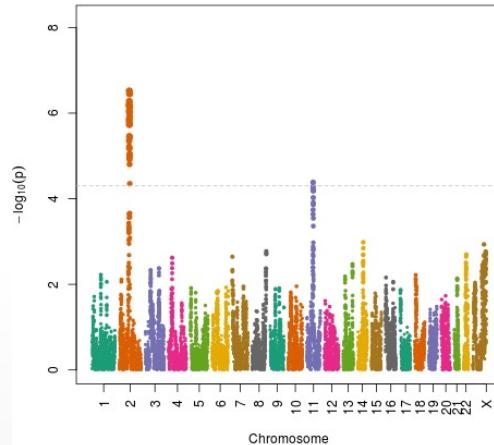
EUR



AFR



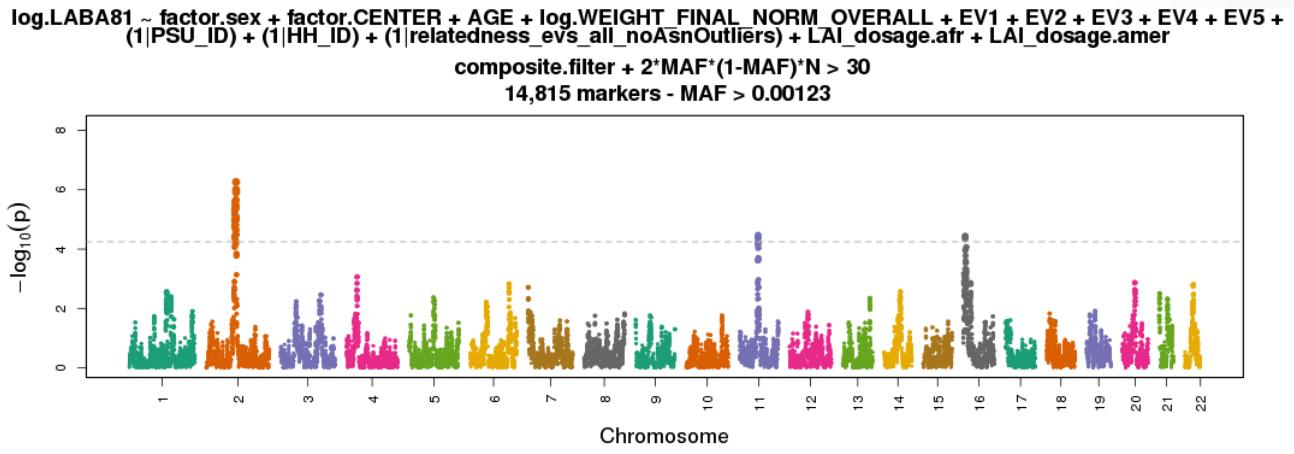
NAM



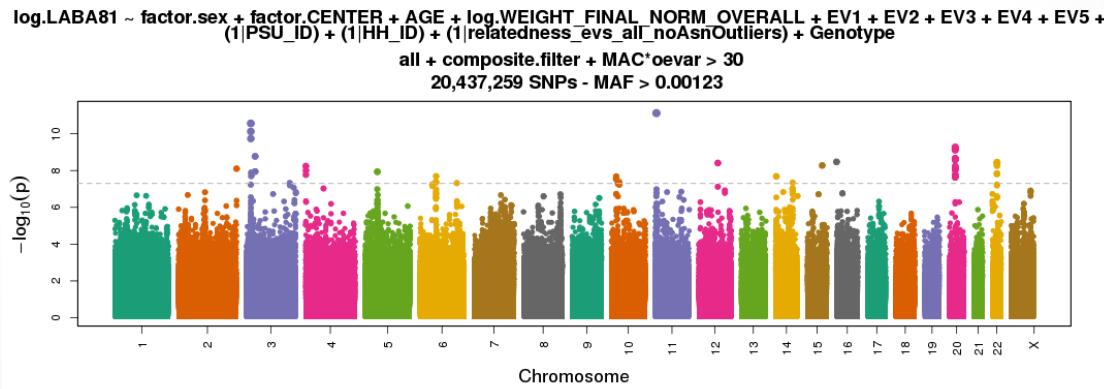
# HCHS/SOL: Admixture and Association Mapping

## Phenotype: Albumin-to-creatinine ratio

Admixture  
Mapping



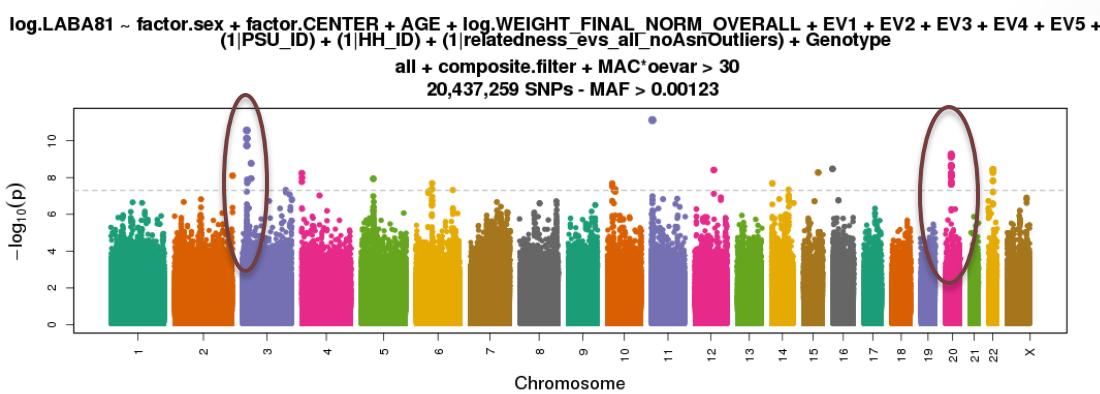
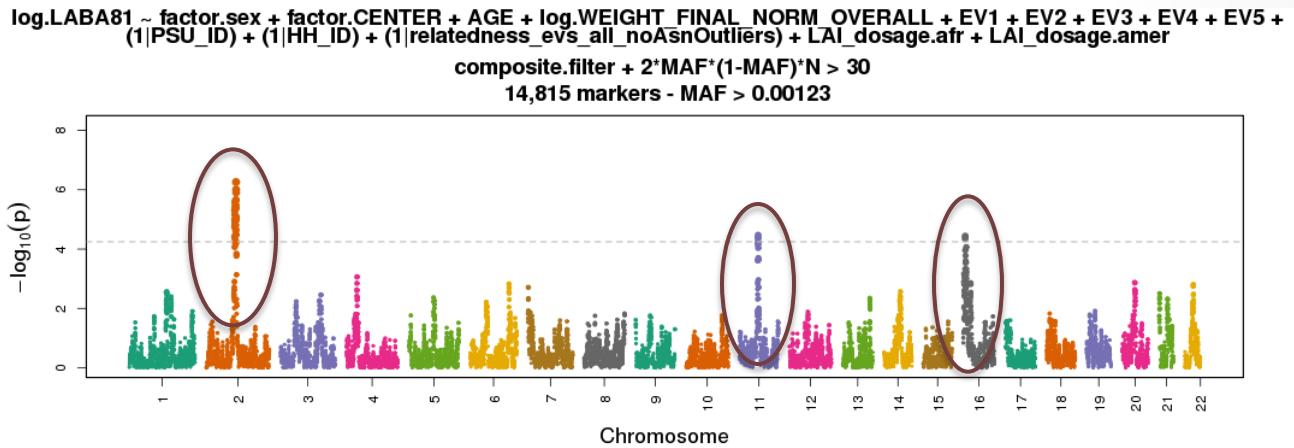
Association  
Mapping



# HCHS/SOL: Admixture and Association Mapping

## Phenotype: Albumin-to-creatinine ratio

Admixture  
Mapping



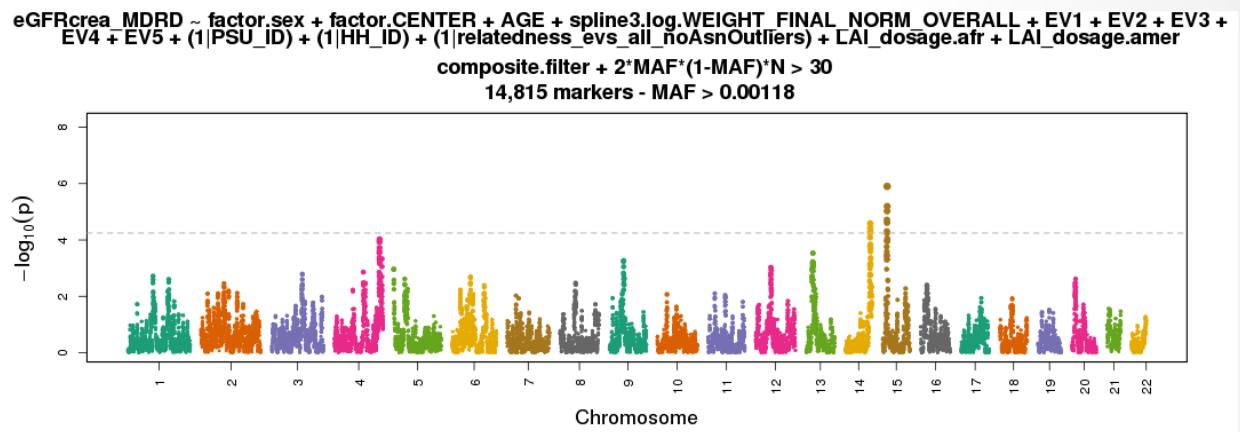
Association  
Mapping

Different  
signals for  
the two  
methods

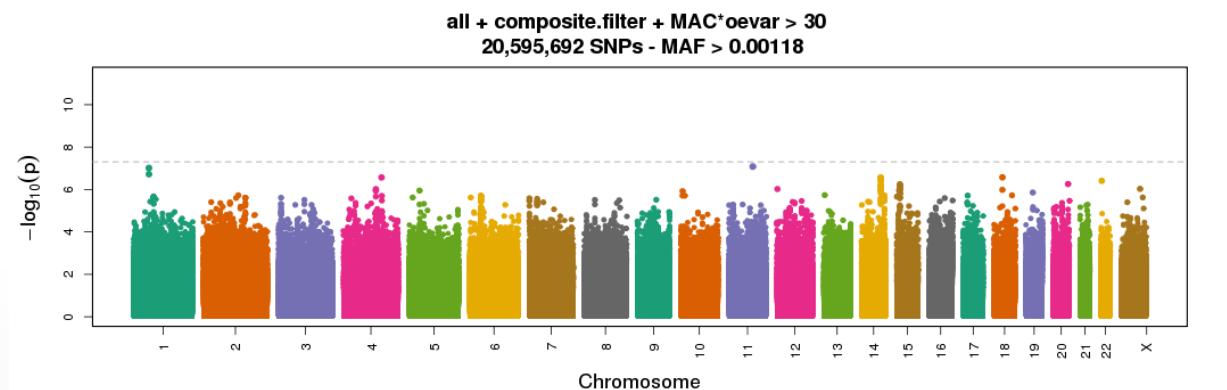
# HCHS/SOL: Admixture and Association Mapping

## Phenotype: Estimated Glomerular Filtration rate

Admixture  
Mapping



Association  
Mapping

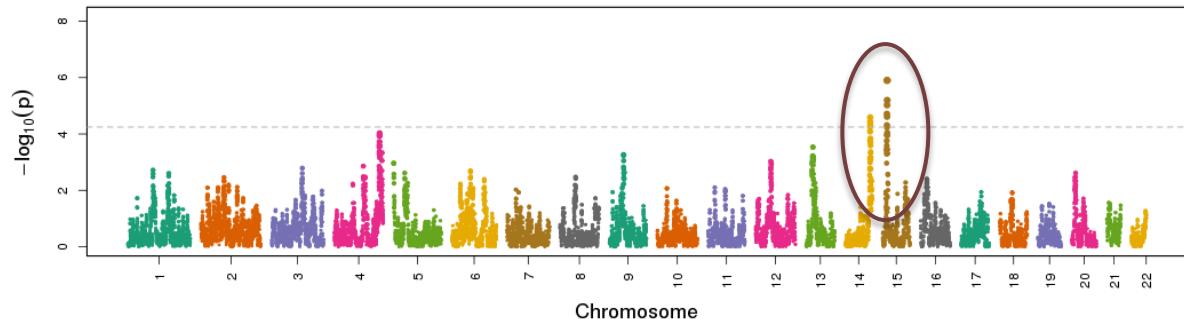


# HCHS/SOL: Admixture and Association Mapping

## Phenotype: Estimated Glomerular Filtration rate

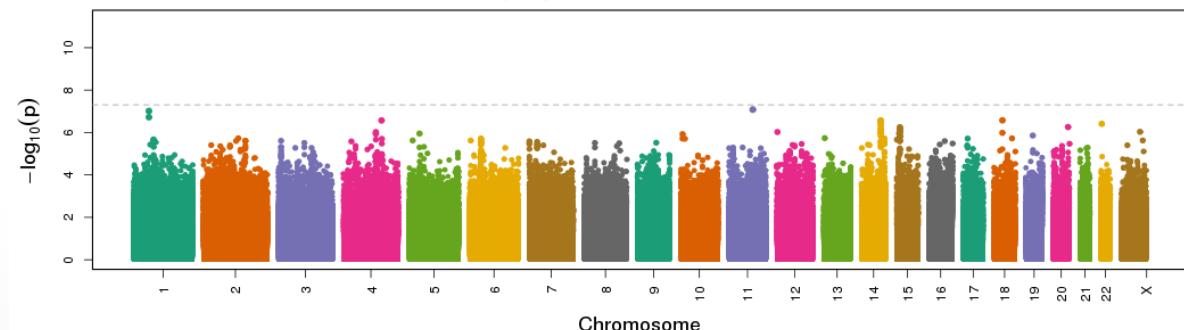
Admixture  
Mapping

eGFRcrea\_MDRD ~ factor.sex + factor.CENTER + AGE + spline3.log.WEIGHT\_FINAL\_NORM\_OVERALL + EV1 + EV2 + EV3 +  
EV4 + EV5 + (1|PSU\_ID) + (1|HH\_ID) + (1|relatedness\_evs\_all\_noAsnOutliers) + LAI dosage.afr + LAI dosage.amer  
composite.filter + 2\*MAF\*(1-MAF)\*N > 30  
14,815 markers - MAF > 0.00118



Only admixture mapping  
provides significant results

all + composite.filter + MAC\*oevar > 30  
20,595,692 SNPs - MAF > 0.00118

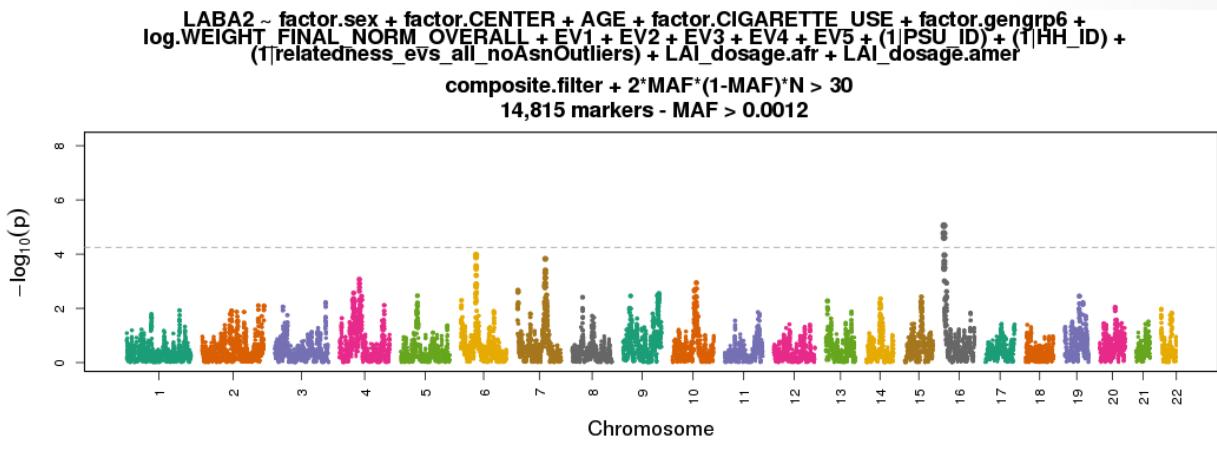


Association  
Mapping

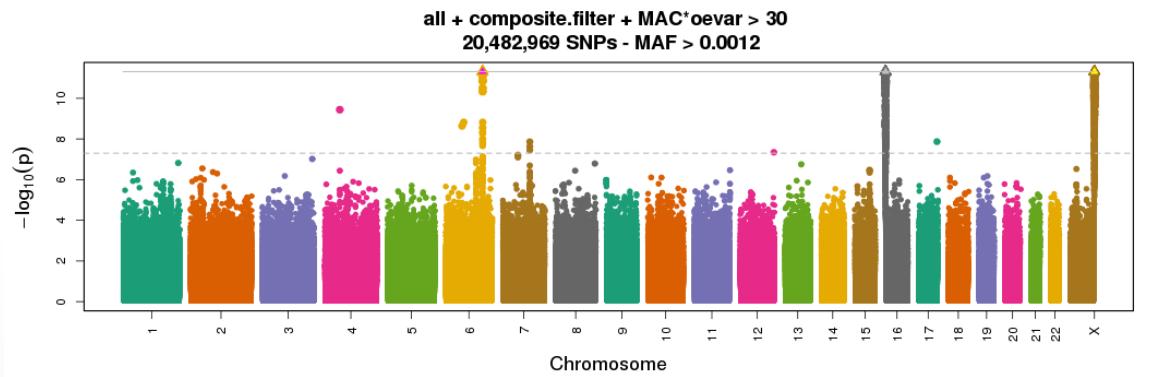
# HCHS/SOL: Admixture and Association Mapping

## Phenotype: Red blood cell count

Admixture  
Mapping



Association  
Mapping



# HCHS/SOL: Admixture and Association Mapping

## Phenotype: Red blood cell count

Admixture  
Mapping

association mapping  
gives more significant  
results

Association  
Mapping

