

# Hidden structure in polygenic scores and the challenge of disentangling ancestry interactions in admixed populations

Joint work with R. Mandla, Z. Shi, B. Paşaniuc and I. Mathieson

Alan J. Aw  
Postdoctoral Fellow  
Department of Genetics  
University of Pennsylvania

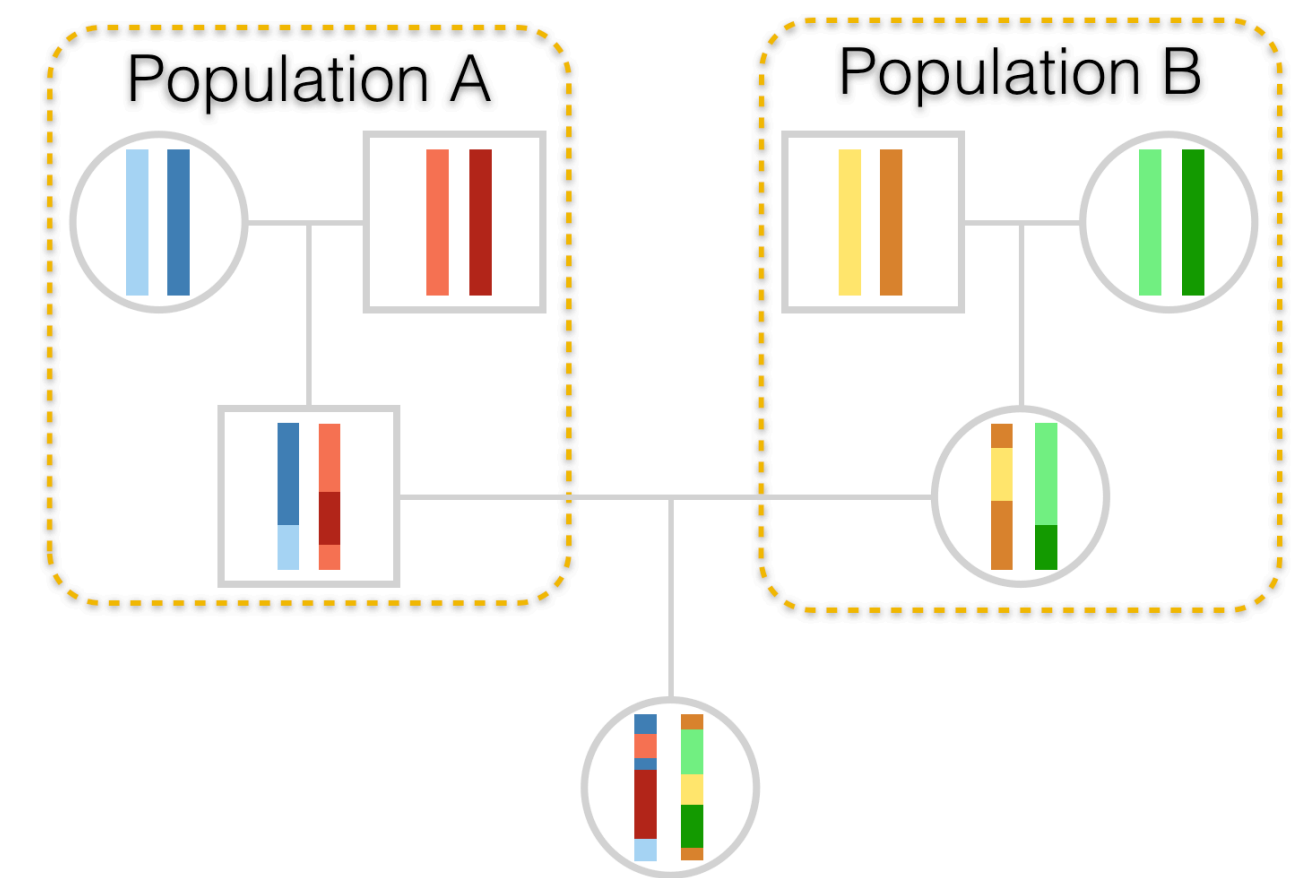


Image Source: Rodney Dyer

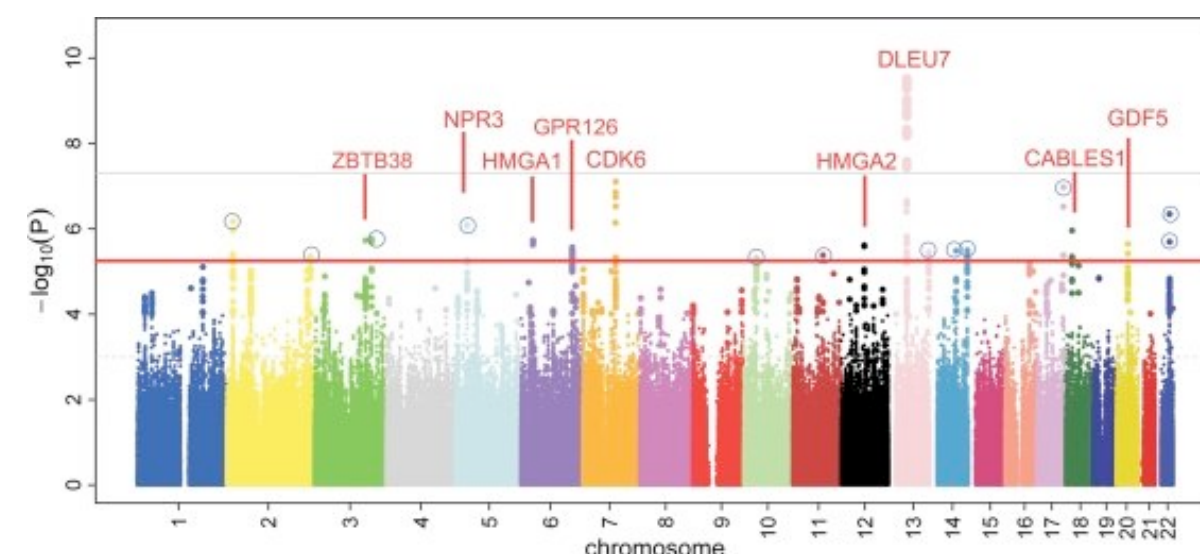
# Complex traits and poor portability

- Complex traits (e.g., height) are influenced by networks of genes that act in concert to regulate expression
- Polygenic scores trained in one population port poorly into other populations

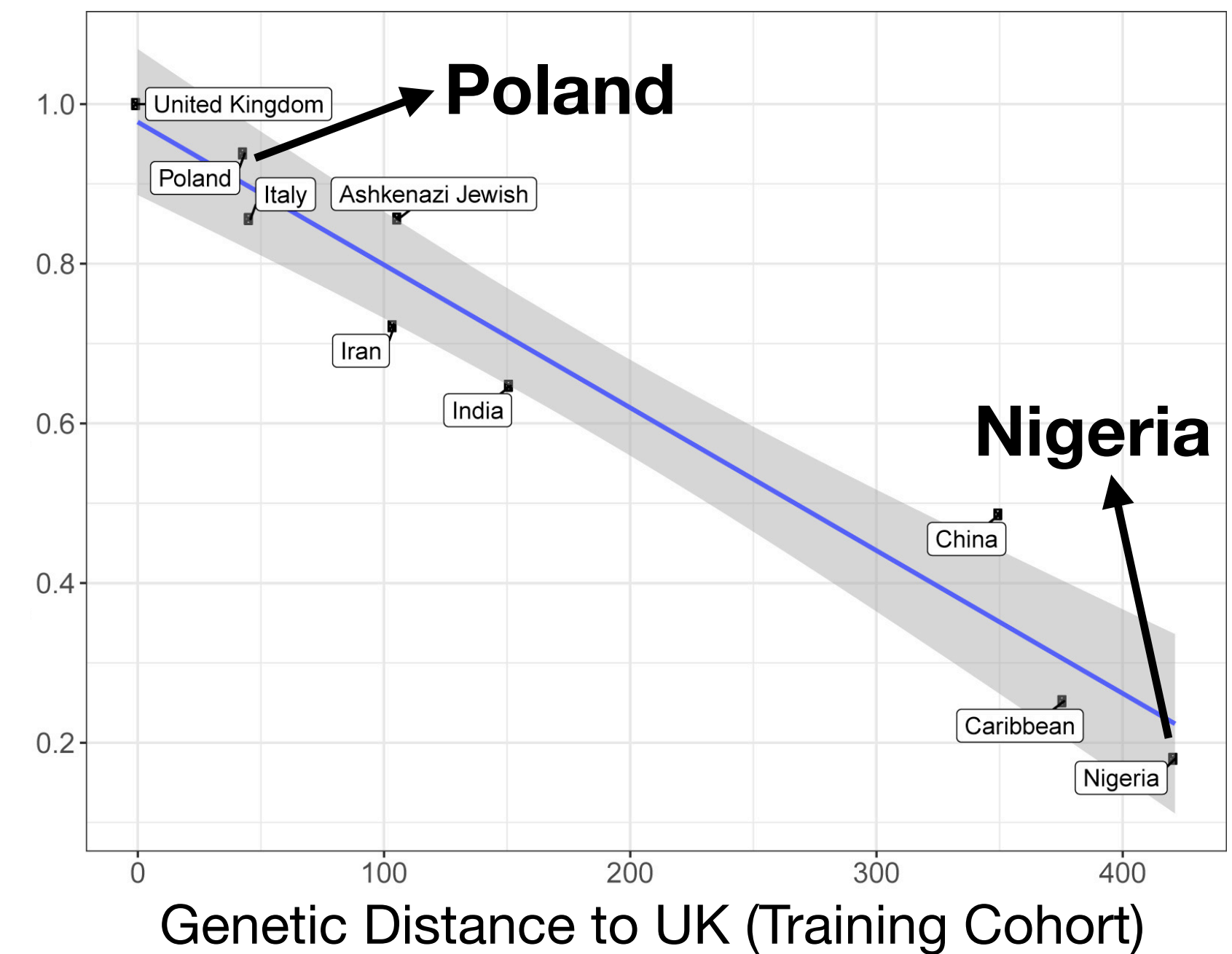
$$\text{PGS}_{\text{Height}} = \beta_1 x_1 + \dots + \beta_{5000} x_{5000}$$

Allelic Dosages

Variant Effect Sizes  
(obtained from European GWAS,  
trained on European samples)



Relative predictive performance  
with UK

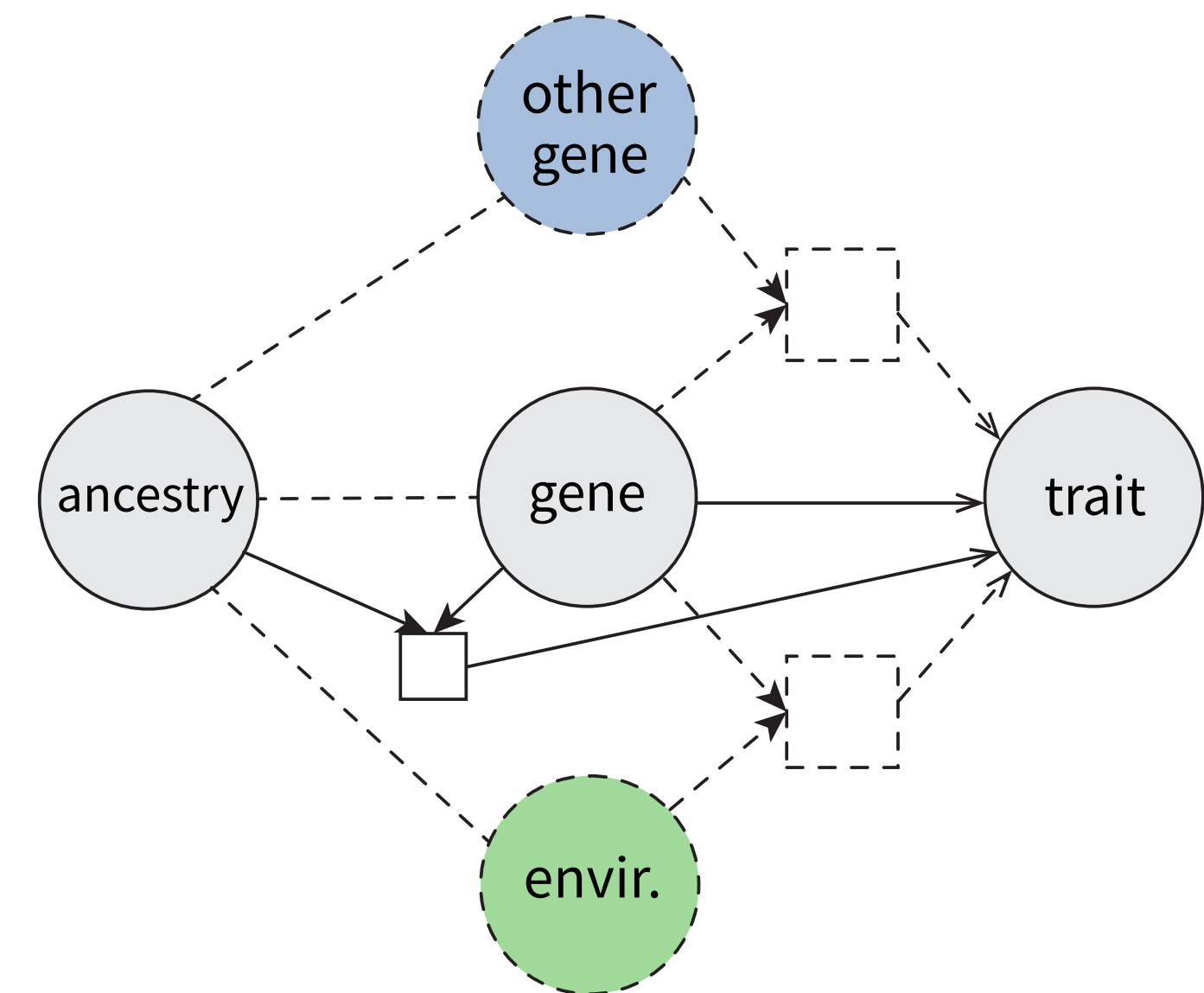
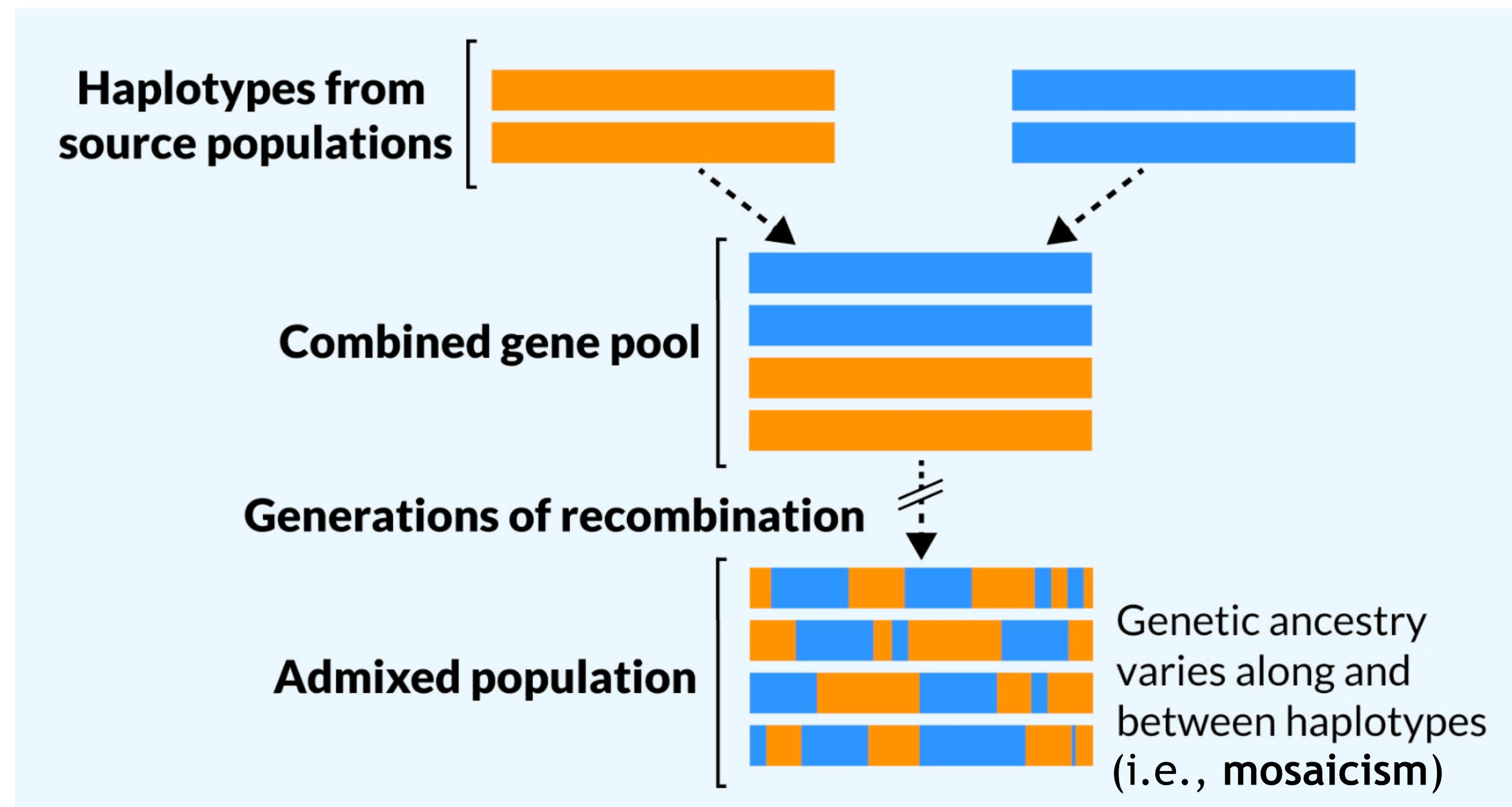


# What explains the poor portability?

1. Recent work (Hou et al., 2023 *Nat. Genet.*; Hu et al., 2025 *Nat. Genet.*) suggests high similarity in causal effects across ancestries
2. Differences in linkage disequilibrium (LD) patterns and allele frequencies between ancestries
3. Interactions (Gene-by-gene [GxG] and Gene-by-environment [GxE])
  - How can causal effects be highly similar in spite of interactions?

# The role of admixed populations

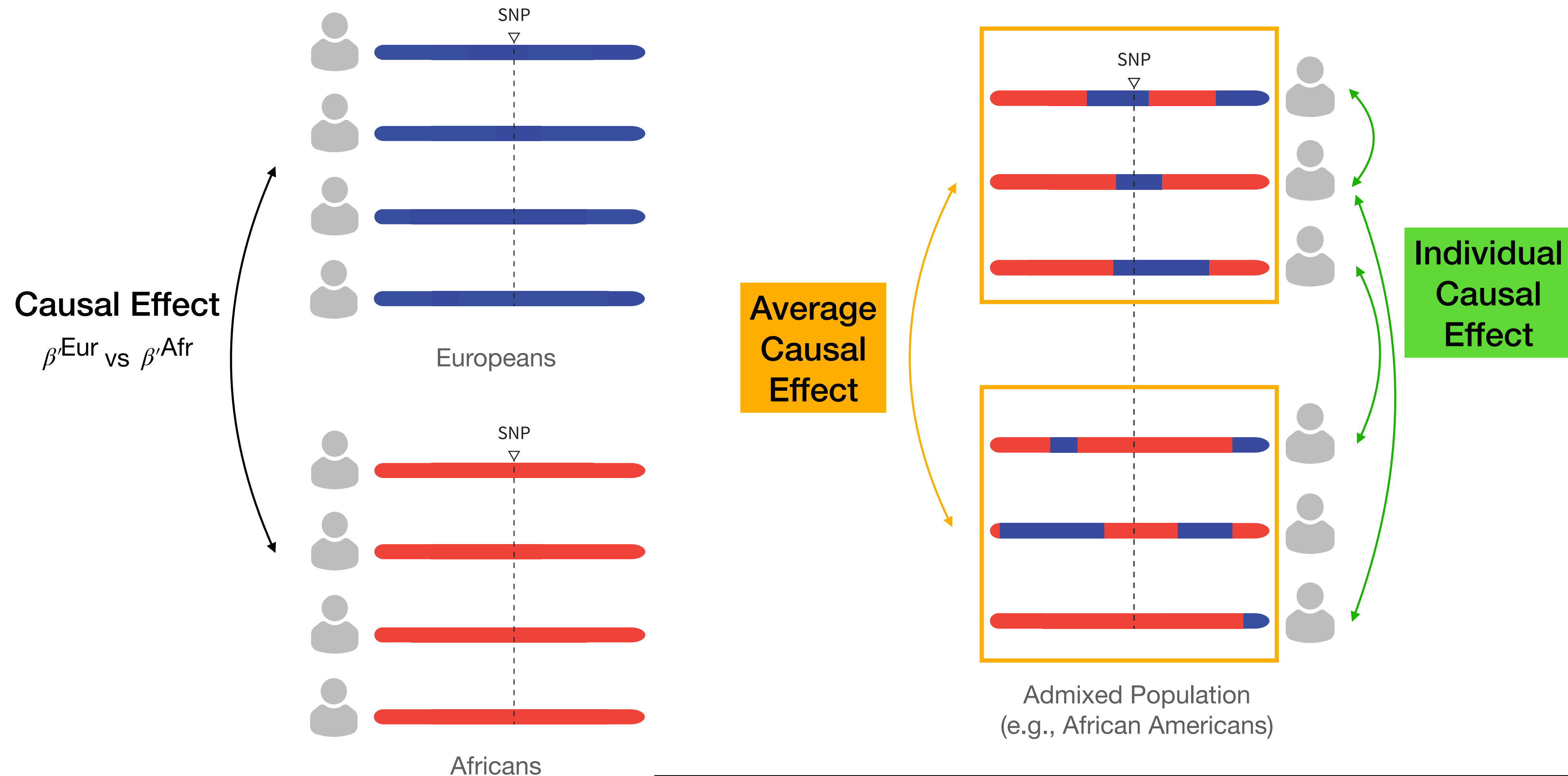
- Ancestry mosaicism in admixed individuals can capture differences in allele frequencies and environmental exposures



Build Statistical Models of Gene-by-Ancestry (GxA) Interactions



# Causal effects are similar *between what?*

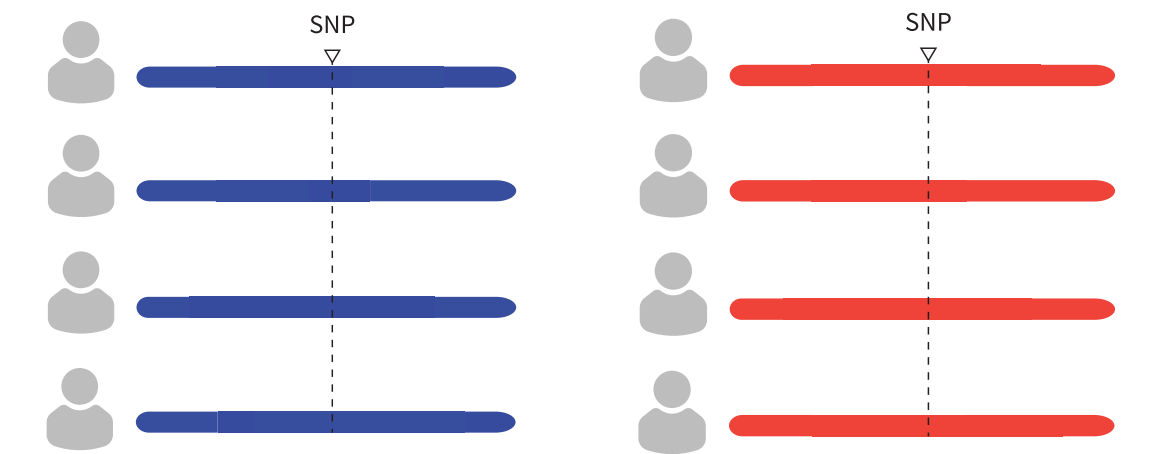


Hou et al. (2023) and Hu et al. (2025): Average causal effects are highly similar across local ancestries

# Base model of causal effects

- Ancestral non-admixed population **causal effect sizes** follow a bivariate normal distribution:

$$\begin{bmatrix} \beta'^{\text{Eur}} \\ \beta'^{\text{Afr}} \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma'^{\text{Eur}^2} & \tau' \\ \tau' & \sigma'^{\text{Afr}^2} \end{bmatrix} \right)$$



$$\text{Variance: } \sigma'^{\text{Eur}^2} = \frac{r^2}{2 \sum_{j=1}^p f_j^{\text{Eur}} (1 - f_j^{\text{Eur}})} \quad \sigma'^{\text{Afr}^2} = \frac{r^2}{2 \sum_{j=1}^p f_j^{\text{Afr}} (1 - f_j^{\text{Afr}})}$$

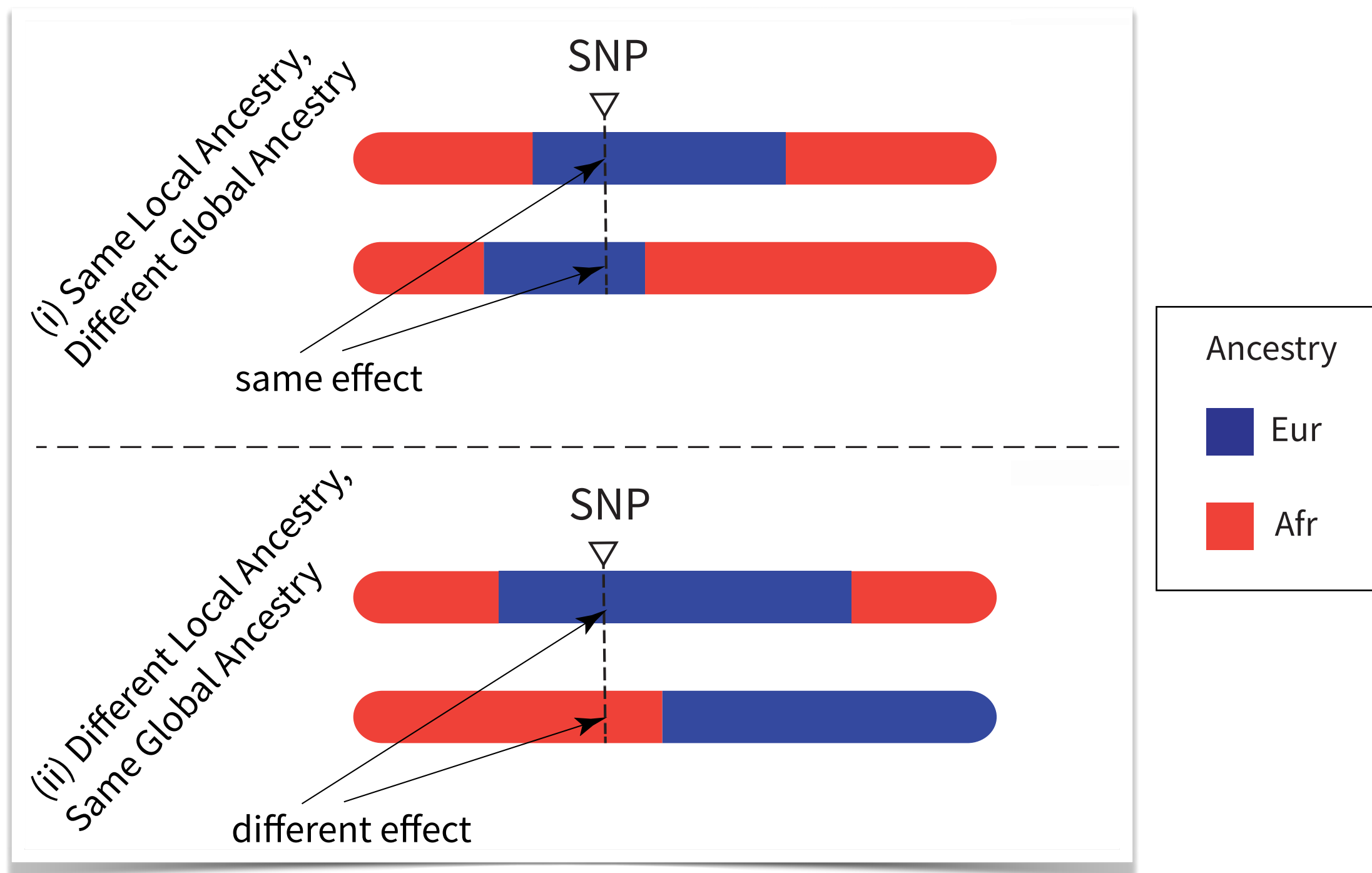
$$\text{Covariance: } \tau' = \frac{\rho r^2}{2 \sqrt{\sum_{j=1}^p f_j^{\text{Eur}} (1 - f_j^{\text{Eur}})} \sqrt{\sum_{j=1}^p f_j^{\text{Afr}} (1 - f_j^{\text{Afr}})}}$$

$$\text{Causal effect correlation} = \tau' / \sqrt{\sigma'^{\text{Eur}^2} \sigma'^{\text{Afr}^2}} = \rho$$

# Two models of gene-by-ancestry interaction

## Local Model

Captures interactions in *cis*

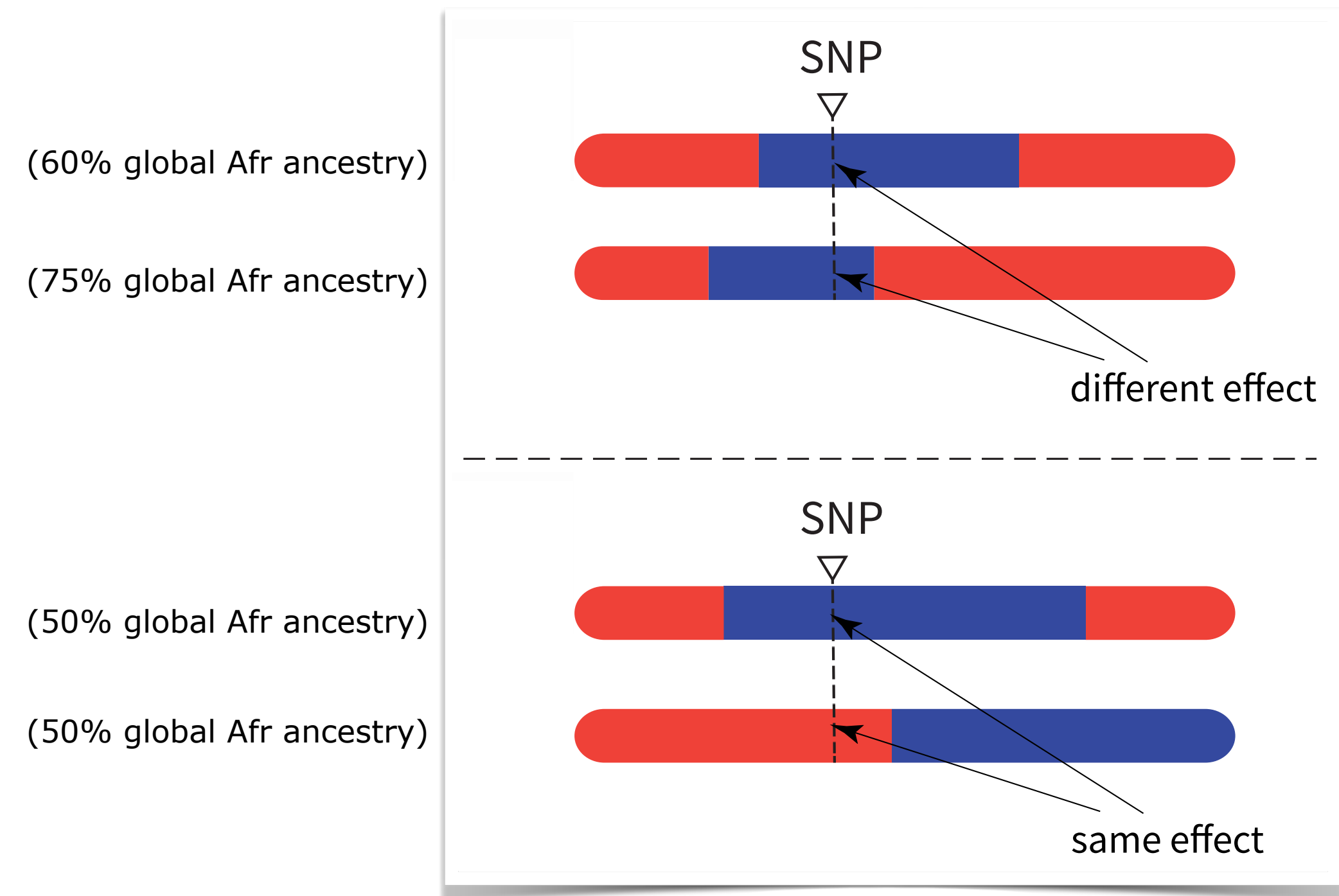


$$\beta'_{\text{SNP}} = \beta'^{\text{Afr}} a + \beta'^{\text{Eur}} (1 - a)$$

$$a = \begin{cases} 1 & \text{if ancestry is Afr} \\ 0 & \text{if ancestry is Eur} \end{cases}$$

## Global Model

Captures interactions in *trans* and GxE

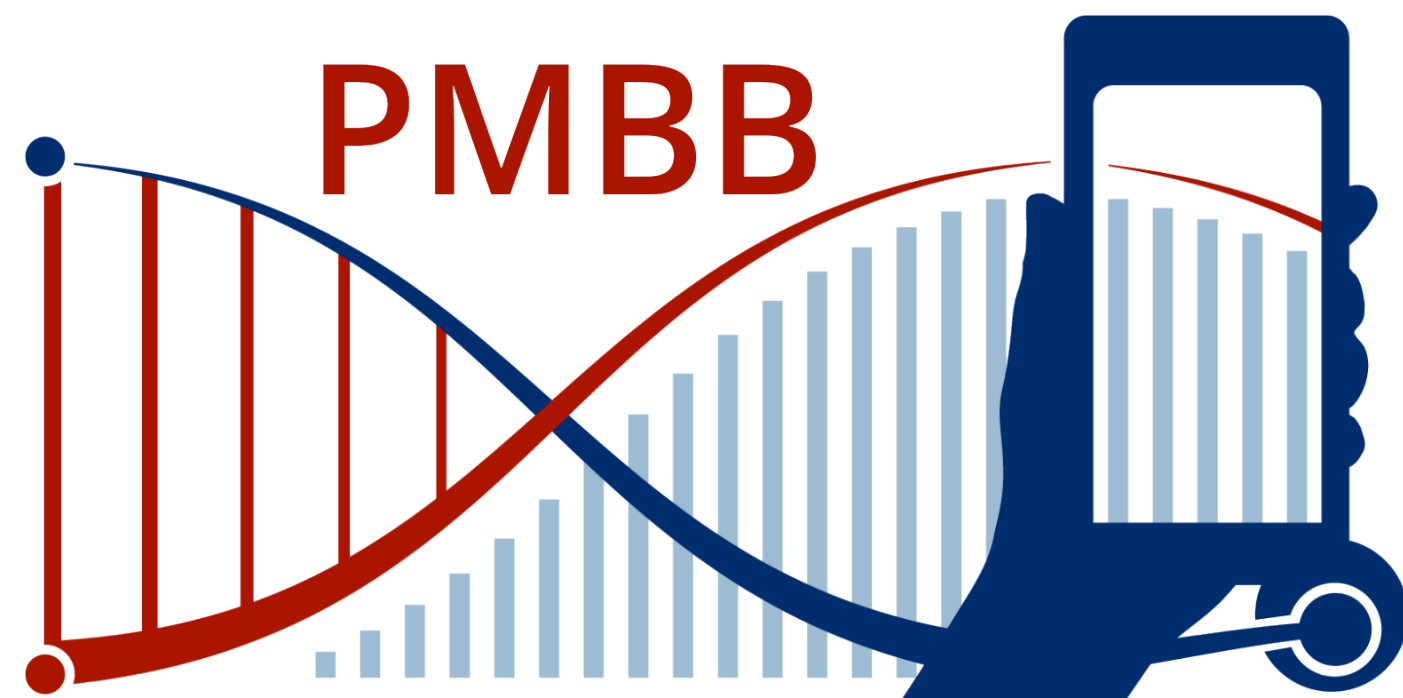


$$\beta'_{\text{SNP}} = \beta'^{\text{Afr}} \bar{a} + \beta'^{\text{Eur}} (1 - \bar{a})$$

$\bar{a}$  = genome-wide/global Afr ancestry

# Questions

1. What do the local and global models imply about individual and average causal effects?
2. Can polygenic scores differentiate the global and local models?



## **Penn Medicine Biobank**

10,000 genotyped African Americans

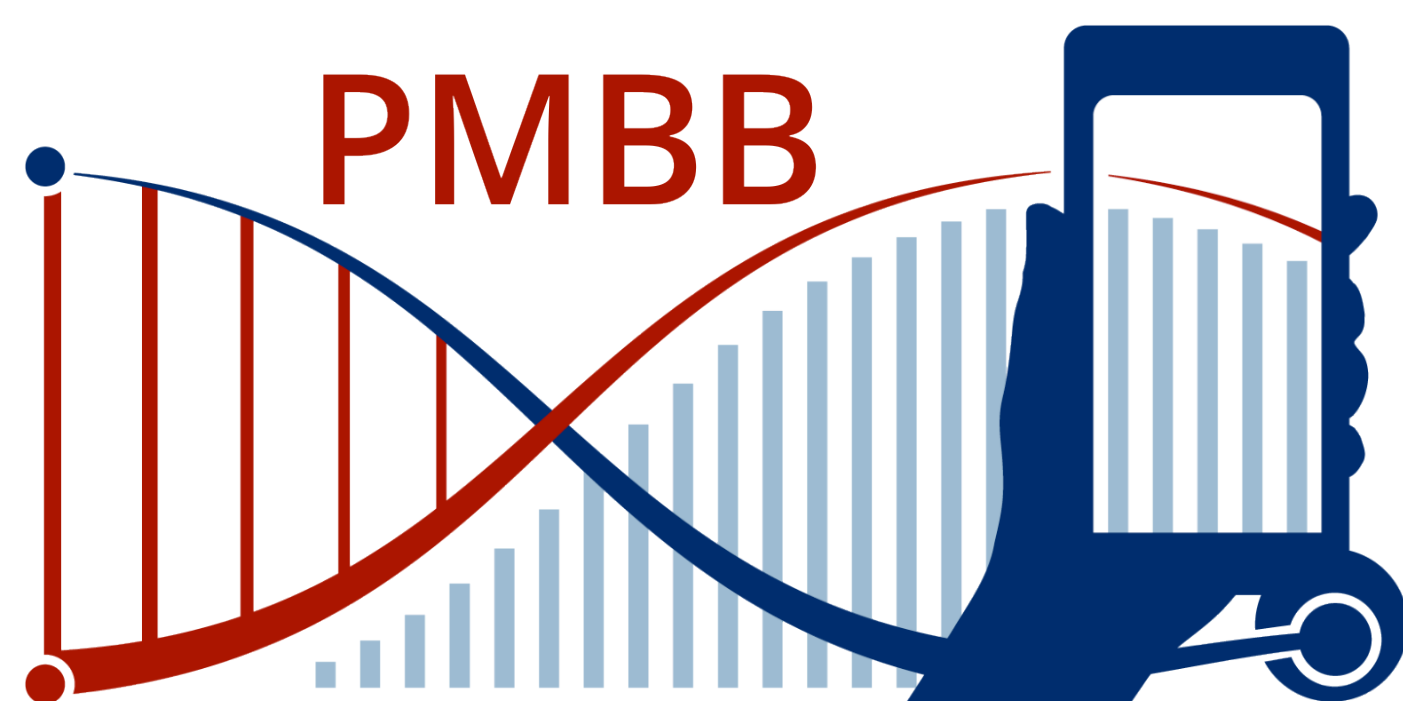
30,000 genotyped European Americans

6 quantitative traits



# Questions

1. What do the local and global models imply about individual and average causal effects?
2. Can polygenic scores differentiate the global and local models?



## **Penn Medicine Biobank**

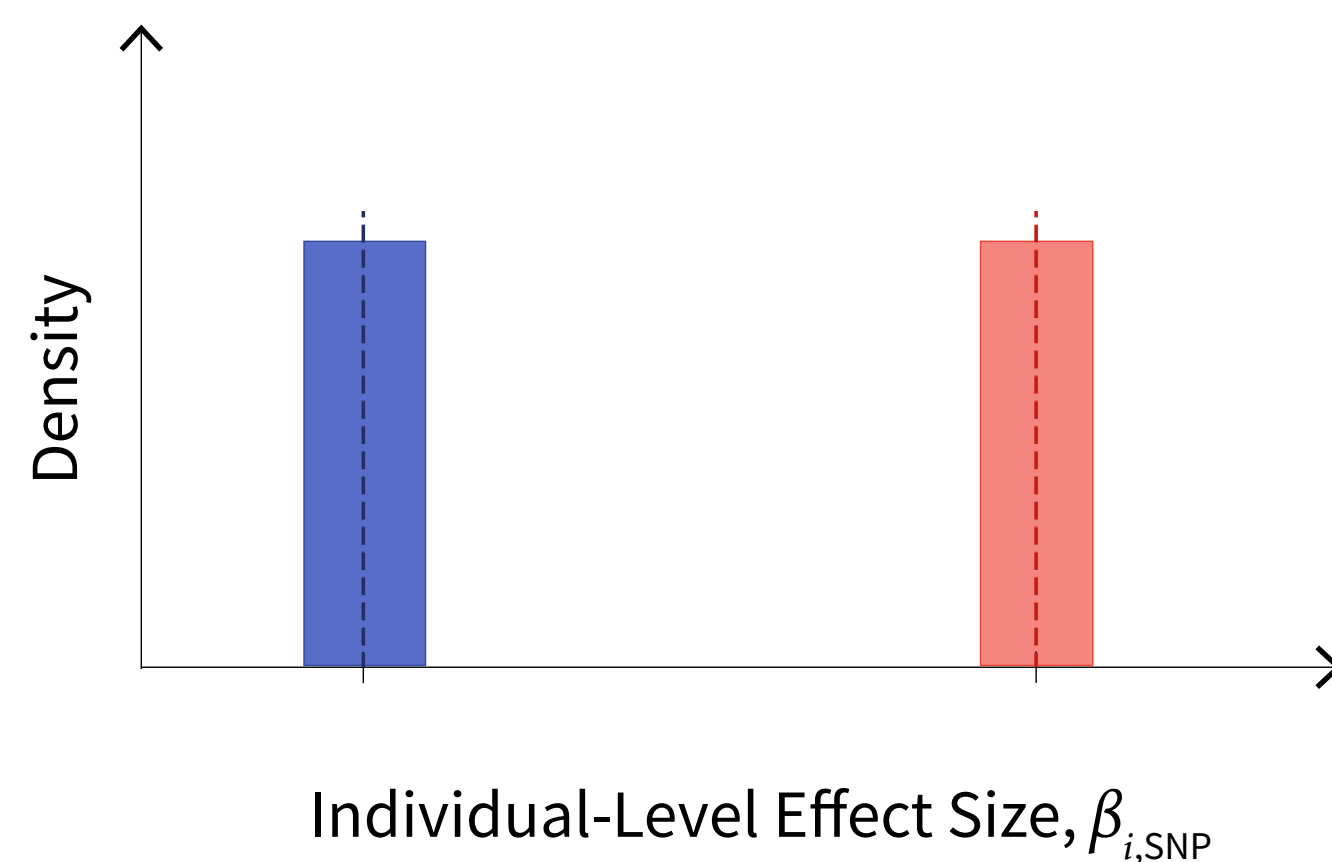
10,000 genotyped African Americans  
30,000 genotyped European Americans  
6 quantitative traits

# Individual Causal Effect

Global model implies high variability in individual effect

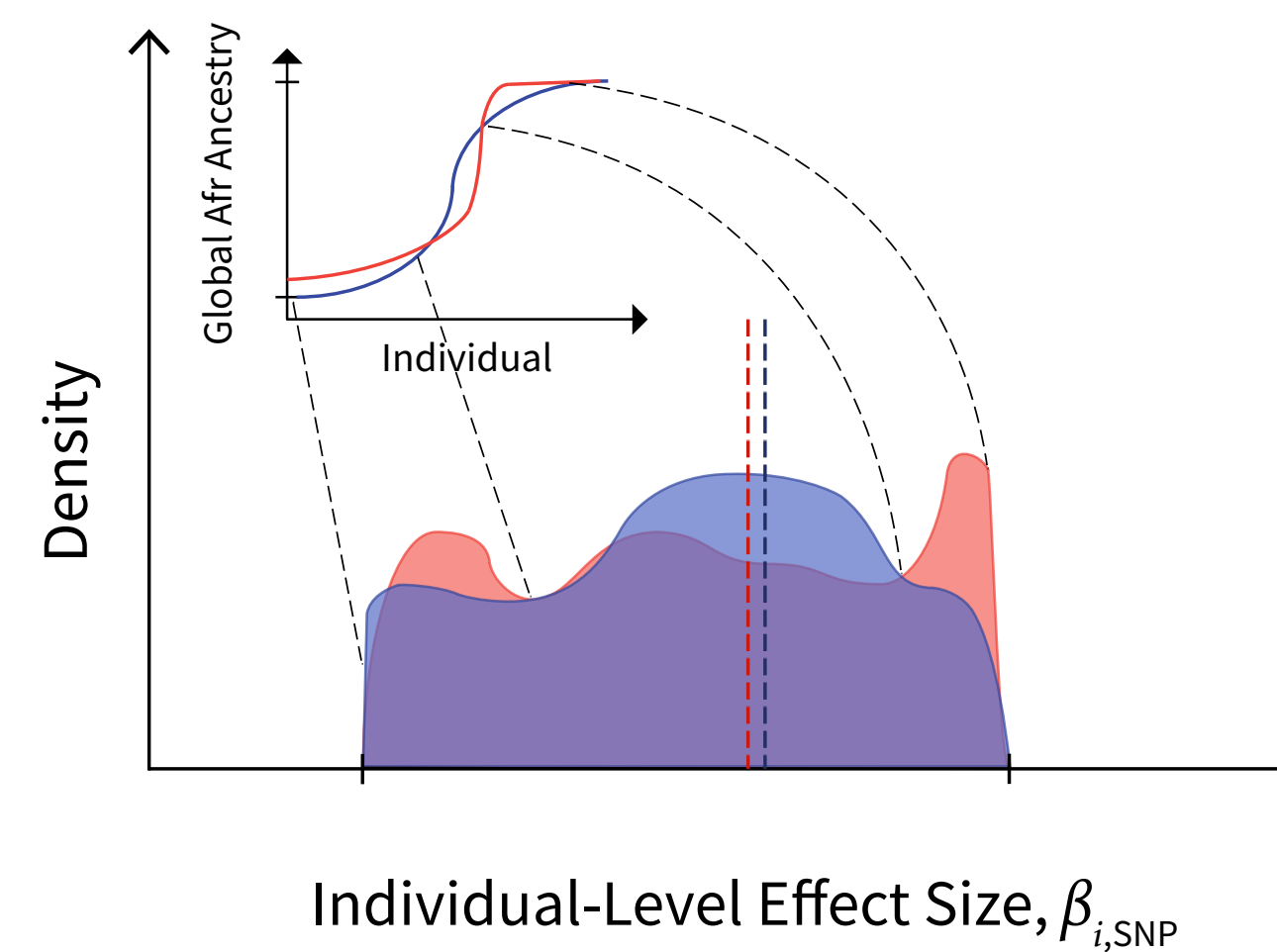
$$\beta'_{\text{SNP}} = \beta'^{\text{Afr}} a + \beta'^{\text{Eur}} (1 - a)$$

Local Model



$$\beta'_{\text{SNP}} = \beta'^{\text{Afr}} \bar{a} + \beta'^{\text{Eur}} (1 - \bar{a})$$

Global Model

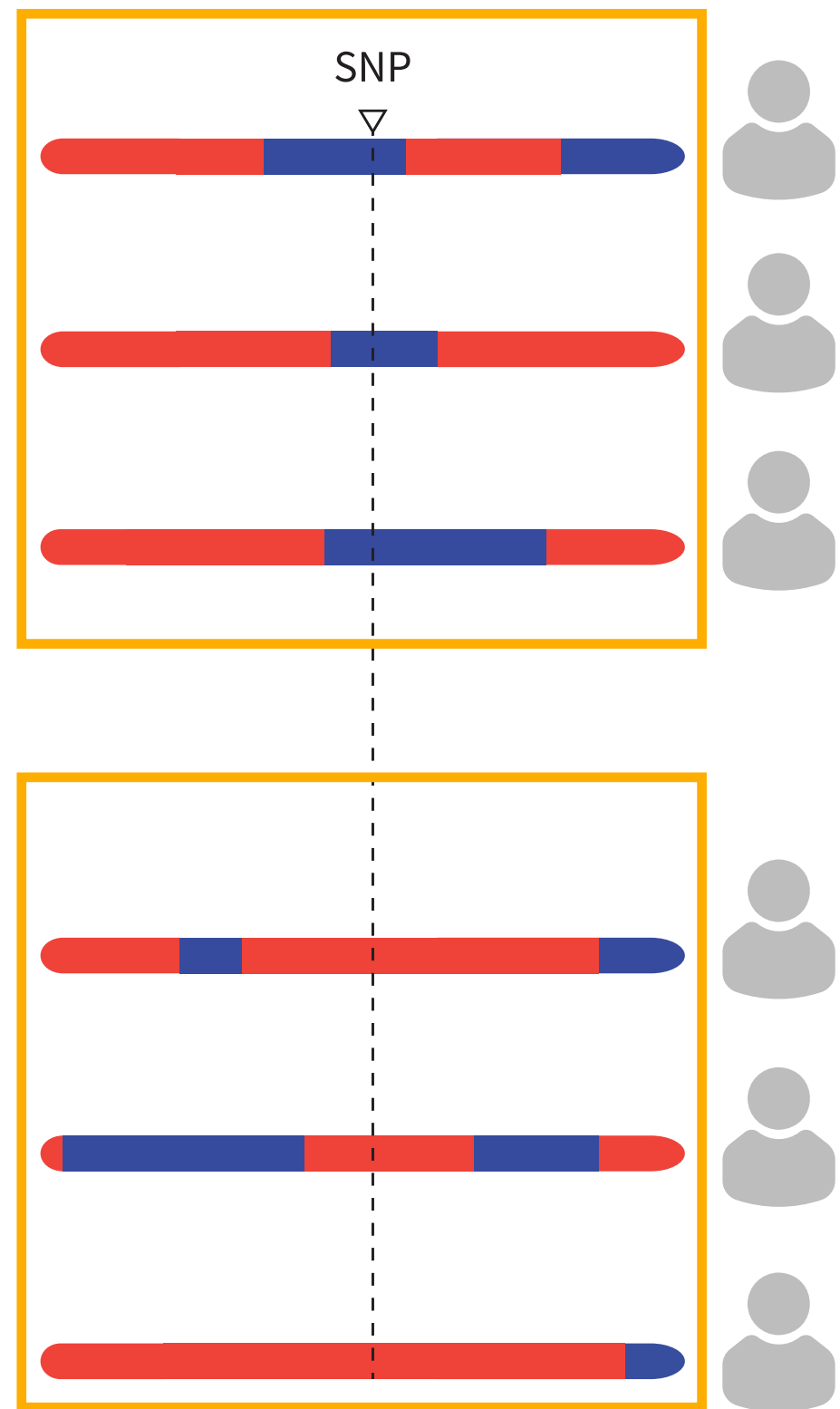


Local Ancestry  
at SNP

— Eur

— Afr

# Average Causal Effect – Local Model

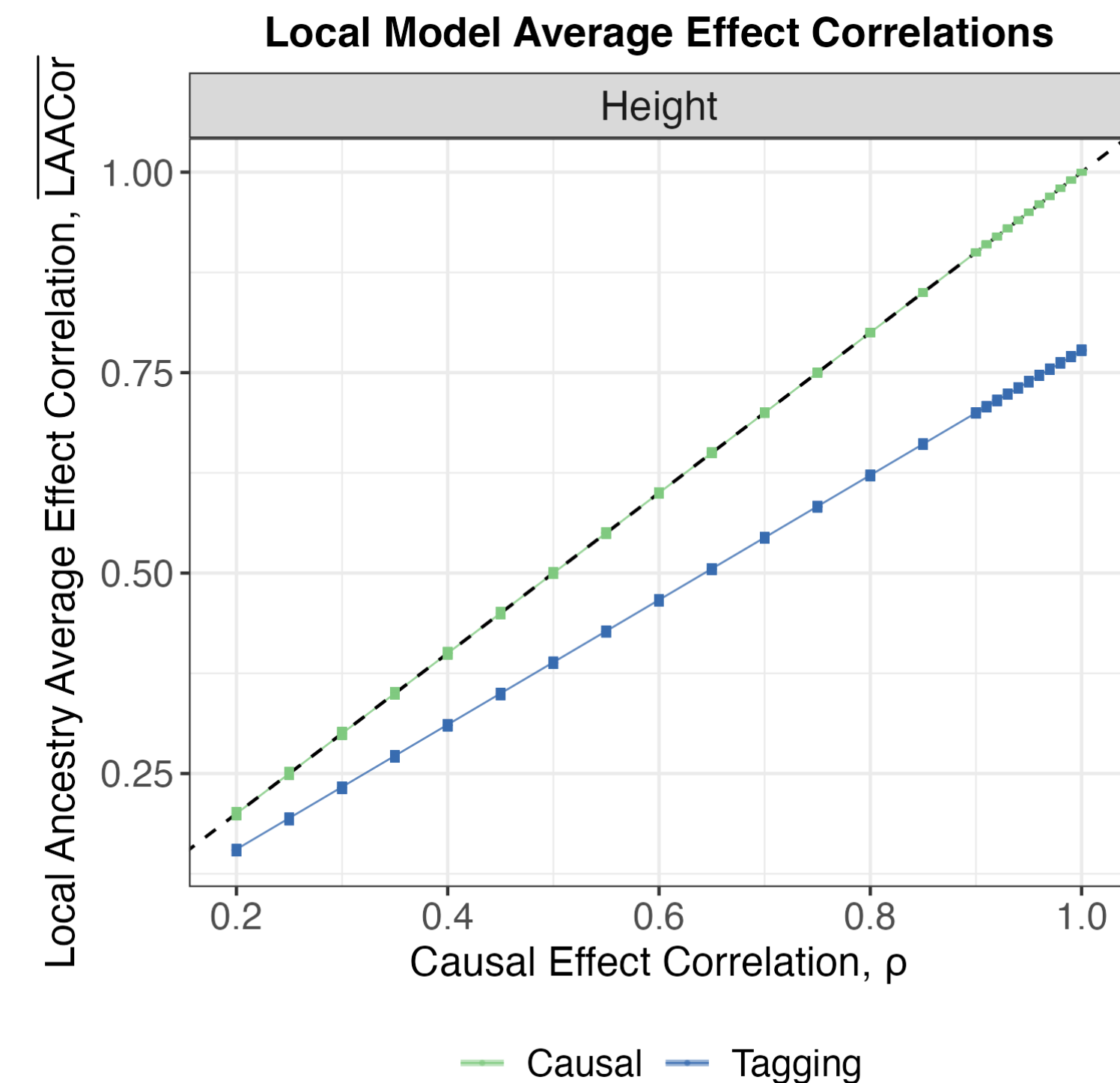


Under local model, distribution of average causal effect is just the distribution of causal effects in base model:

$$\begin{bmatrix} \bar{\beta}' |_{LA=Eur} \\ \bar{\beta}' |_{LA=Afr} \end{bmatrix} \stackrel{d}{=} \begin{bmatrix} \beta'^{Eur} \\ \beta'^{Afr} \end{bmatrix}$$

LAACor = *Local Ancestry Average Causal Effect Correlation*

$$\overline{LAACor'}_{Loc} \stackrel{def}{=} \frac{\text{cov}(\bar{\beta}' |_{LA=Eur}, \bar{\beta}' |_{LA=Afr})}{\sqrt{\text{var}(\bar{\beta}' |_{LA=Eur}) \text{var}(\bar{\beta}' |_{LA=Afr})}} = \rho$$



# Average Causal Effect – Global Model

## Global model produces high average causal effect similarity

**Proposition 4.2** (Joint Distribution of Local Ancestry Average Causal Effects). *Under the local model, the joint distribution of average causal effects is the same as the original joint distribution of causal effects in the base model, Eq. (2). Let  $\beta'_{\cdot j}|_{LA=Afr}$  and  $\beta'_{\cdot j}|_{LA=Eur}$  denote the African and European local ancestry average causal effects under the global model. The joint distribution of these quantities is*

$$\begin{bmatrix} \beta'_{\cdot j}|_{LA=Afr} \\ \beta'_{\cdot j}|_{LA=Eur} \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} u'_j & w'_j \\ w'_j & v'_j \end{bmatrix} \right),$$

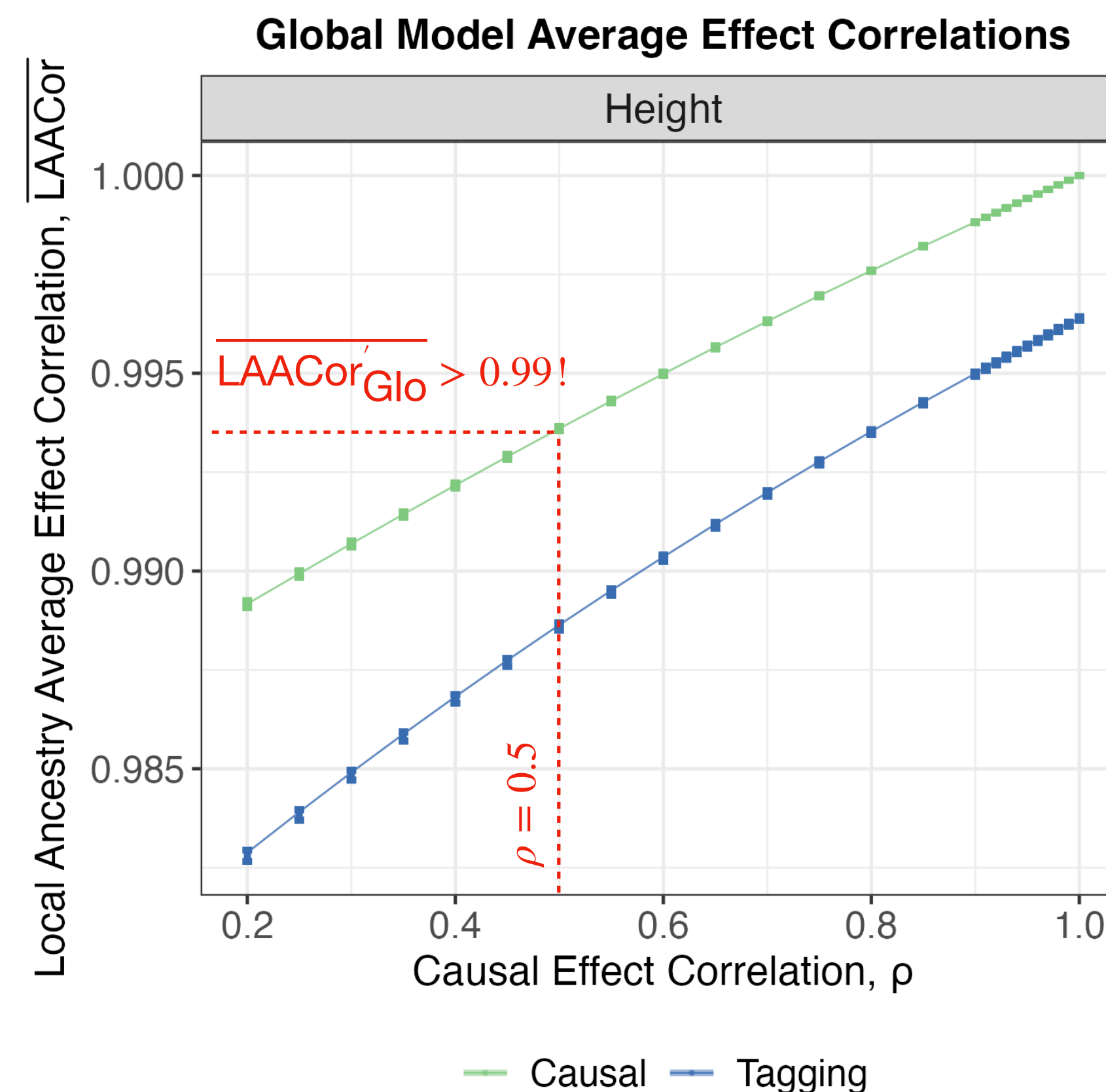
where

$$\begin{aligned} u'_j &= \sigma_{Eur}^{\prime 2} \omega_{1,j}^{\prime 2} + 2\tau' \omega'_{1,j} \omega'_{2,j} + \sigma_{Afr}^{\prime 2} \omega_{2,j}^{\prime 2} \\ v'_j &= \sigma_{Eur}^{\prime 2} \omega_{3,j}^{\prime 2} + 2\tau' \omega'_{3,j} \omega'_{4,j} + \sigma_{Afr}^{\prime 2} \omega_{4,j}^{\prime 2} \\ w'_j &= \sigma_{Eur}^{\prime 2} \omega'_{1,j} \omega'_{3,j} + \tau' (\omega'_{2,j} \omega'_{3,j} + \omega'_{1,j} \omega'_{4,j}) + \sigma_{Afr}^{\prime 2} \omega'_{2,j} \omega'_{4,j} \end{aligned}$$

are terms in the covariance matrix, with quantities  $\sigma_{Afr}^{\prime 2}$ ,  $\sigma_{Eur}^{\prime 2}$  and  $\tau'$  defined in Eqs. (3)-(5), and quantities  $\omega'_{1,j}$ ,  $\omega'_{2,j}$ ,  $\omega'_{3,j}$ ,  $\omega'_{4,j}$  defined in Supplementary Material Subsection S8 (Box C) depending only on the haplotype and local ancestry matrices.

**LAACor = Local Ancestry Average Causal Effect Correlation**

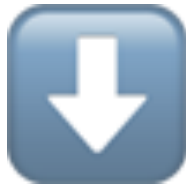

$$\Rightarrow \overline{\text{LAACor}'_{\text{Glo}}} \approx \frac{\sum_{j=1}^p w'_j}{\sqrt{\sum_{j=1}^p u'_j} \sqrt{\sum_{j=1}^p v'_j}}$$





# Summary of Q1

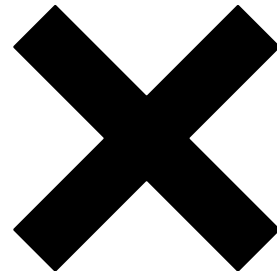
- High (average) causal effect similarity does not rule out variability in individual causal effect

Model	Individual Causal Effect Variability	Average Causal Effect Similarity
Local Model		Same as causal effect correlation $\rho$
Global Model		Can be very high, despite small $\rho$

## Example: Height

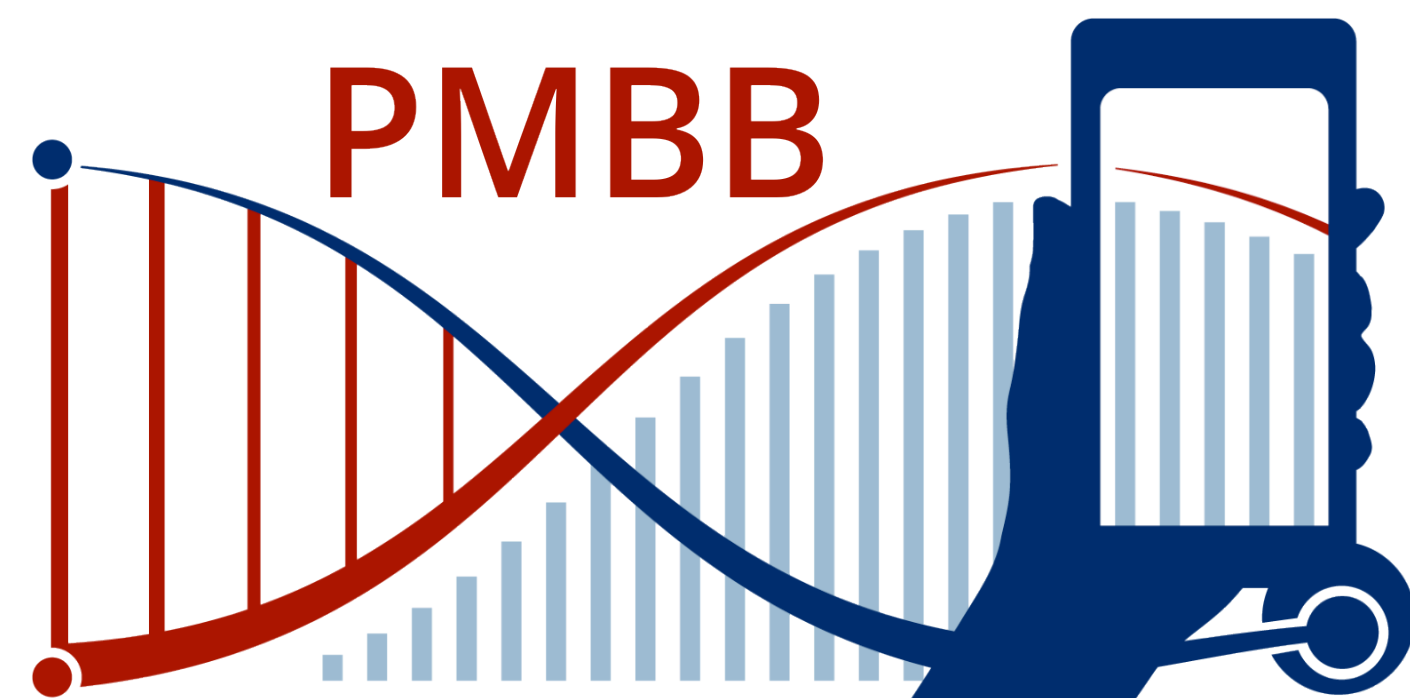
Hou et al. (2023): LAACor is about 0.94

Hu et al. (2025): LAACor lies in  $[0.9, 1]$ , 95% CI contains 1

Local Model		
Global Model		
	Low $\rho$	High $\rho$

# Questions

1. What do the local and global models imply about individual and average causal effects?
- 2. Can polygenic scores differentiate the global and local models?**



## **Penn Medicine Biobank**

10,000 genotyped African Americans  
30,000 genotyped European Americans  
6 quantitative traits

# Polygenic scores (computed on tagging variants)

$$\beta^{\text{Eur}} = \beta'^{\text{Eur}} \cdot \text{LD}^{\text{Eur}} \cdot \sqrt{\frac{f'^{\text{Eur}} (1 - f'^{\text{Eur}})}{f^{\text{Eur}} (1 - f^{\text{Eur}})}}$$

Diagram illustrating the formula for the polygenic score in Europeans ( $\beta^{\text{Eur}}$ ).

The formula is:  $\beta^{\text{Eur}} = \beta'^{\text{Eur}} \cdot \text{LD}^{\text{Eur}} \cdot \sqrt{\frac{f'^{\text{Eur}} (1 - f'^{\text{Eur}})}{f^{\text{Eur}} (1 - f^{\text{Eur}})}}$

Annotations:

- $\beta'^{\text{Eur}}$ : Causal Effect
- $\text{LD}^{\text{Eur}}$ : Causal and Tagging Variant LD
- $f^{\text{Eur}}$ : Causal Variant Allele Frequency
- $f'^{\text{Eur}}$ : Tagging Variant Allele Frequency

$$\beta^{\text{Afr}} = \beta'^{\text{Afr}} \cdot \text{LD}^{\text{Afr}} \cdot \sqrt{\frac{f'^{\text{Afr}} (1 - f'^{\text{Afr}})}{f^{\text{Afr}} (1 - f^{\text{Afr}})}}$$

See: Zaidi (2020) [blog post]; Vukcevic et al. (2011) *AJHG*

# Polygenic scores (computed on tagging variants)

- **Standard, or Total, polygenic score:** assign European effect sizes to all alleles
- **Partial polygenic score:** restrict to genomic chunks of European ancestry only

Individual Genotype and Local Ancestry	SNP 1	SNP 2	SNP 3	SNP 4
	1	0	1	0
	1	1	0	1

Ancestry

Eur

Afr

$$\text{TotPGS} = \beta_1 \times \left( \frac{1 - f_1^{\text{Afr}}}{1 - f_1^{\text{Eur}}} \right) + \beta_2 \times \left( \frac{-f_2^{\text{Eur}}}{1 - f_2^{\text{Eur}}} \right) + \beta_3 \times \left( \frac{1 - f_3^{\text{Afr}}}{-f_3^{\text{Afr}}} \right) + \beta_4 \times \left( \frac{-f_4^{\text{Eur}}}{1 - f_4^{\text{Afr}}} \right)$$

$$\text{ParPGS} = \beta_1 \times \left( \frac{1 - f_1^{\text{Eur}}}{1 - f_1^{\text{Eur}}} \right) + \beta_2 \times \left( \frac{-f_2^{\text{Eur}}}{1 - f_2^{\text{Eur}}} \right) + \beta_3 \times \left( \frac{1 - f_3^{\text{Afr}}}{-f_3^{\text{Afr}}} \right) + \beta_4 \times \left( \frac{-f_4^{\text{Eur}}}{1 - f_4^{\text{Afr}}} \right)$$

(More on Partial PGS: Sun et al., 2024 *Nat. Comm.*; Marnetto et al., 2020 *Nat. Comm.*; Bitarello and Mathieson, 2020 G3)



# Causal Variants Known: Partial PGS differentiates the two models (but Total PGS does not)

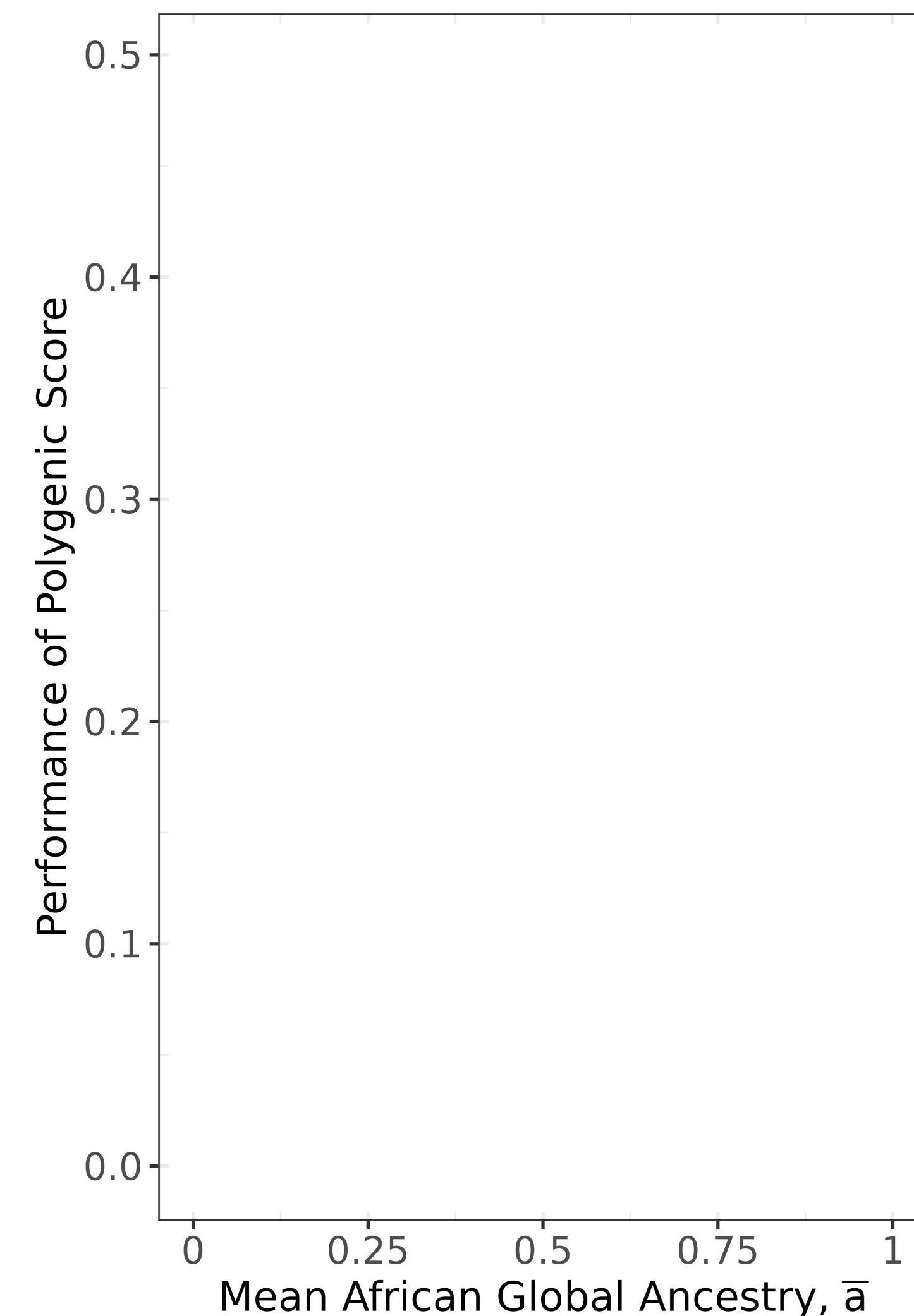
- Partial PGS performance declines **cubically** in global ancestry under the global model, but declines **linearly** under the local model

**Global Model:**  $\mathbb{E}[\text{cor}^2(\text{ParPGS}, y)] \approx r^2(1 - \bar{a})(1 - \bar{a} + \rho\bar{a})^2$

Correlation in causal effects between ancestries

(Mean African) Global ancestry

Partial PGS Performance vs Global African Ancestry



# Causal Variants Known: Partial PGS differentiates the two models (but Total PGS does not)

- Partial PGS performance declines **cubically** in global ancestry under the global model, but declines **linearly** under the local model

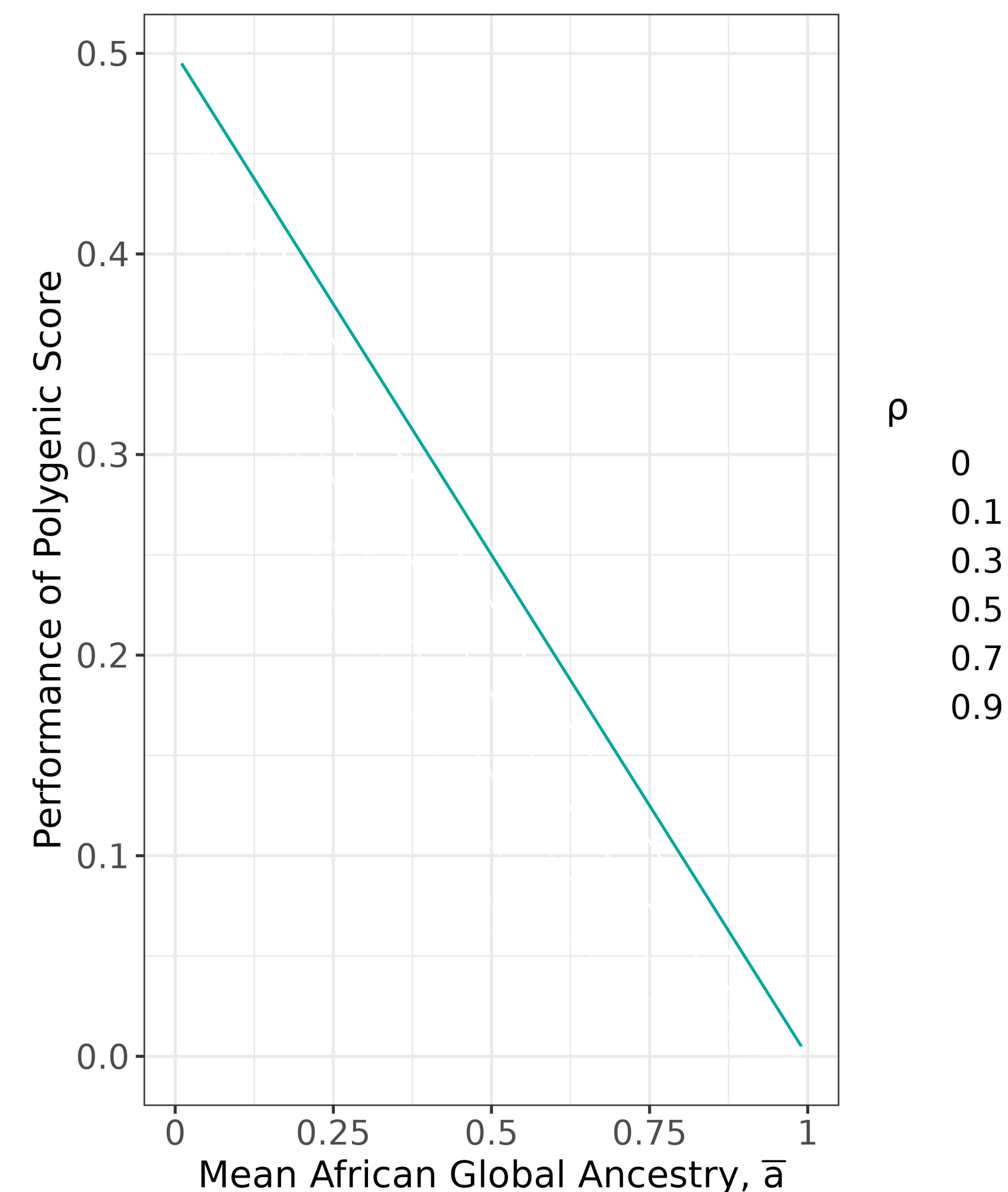
Correlation in tagging effects between ancestries

Global Model:  $\mathbb{E}[\text{cor}^2(\text{ParPGS}, y)] \approx r^2(1 - \bar{a})(1 - \bar{a} + \rho\bar{a})^2$

Local Model:  $\mathbb{E}[\text{cor}^2(\text{ParPGS}, y)] \approx r^2(1 - \bar{a})$

(Mean African) Global ancestry

Partial PGS Performance vs Global African Ancestry



# Causal Variants Known: Partial PGS differentiates the two models (but Total PGS does not)

- Partial PGS performance declines **cubically** in global ancestry under the global model, but declines **linearly** under the local model

Correlation in causal effects between ancestries

Global Model:  $\mathbb{E}[\text{cor}^2(\text{ParPGS}, y)] \approx r^2(1 - \bar{a})(1 - \bar{a} + \rho\bar{a})^2$

Local Model:  $\mathbb{E}[\text{cor}^2(\text{ParPGS}, y)] \approx r^2(1 - \bar{a})$

(Mean African) Global ancestry

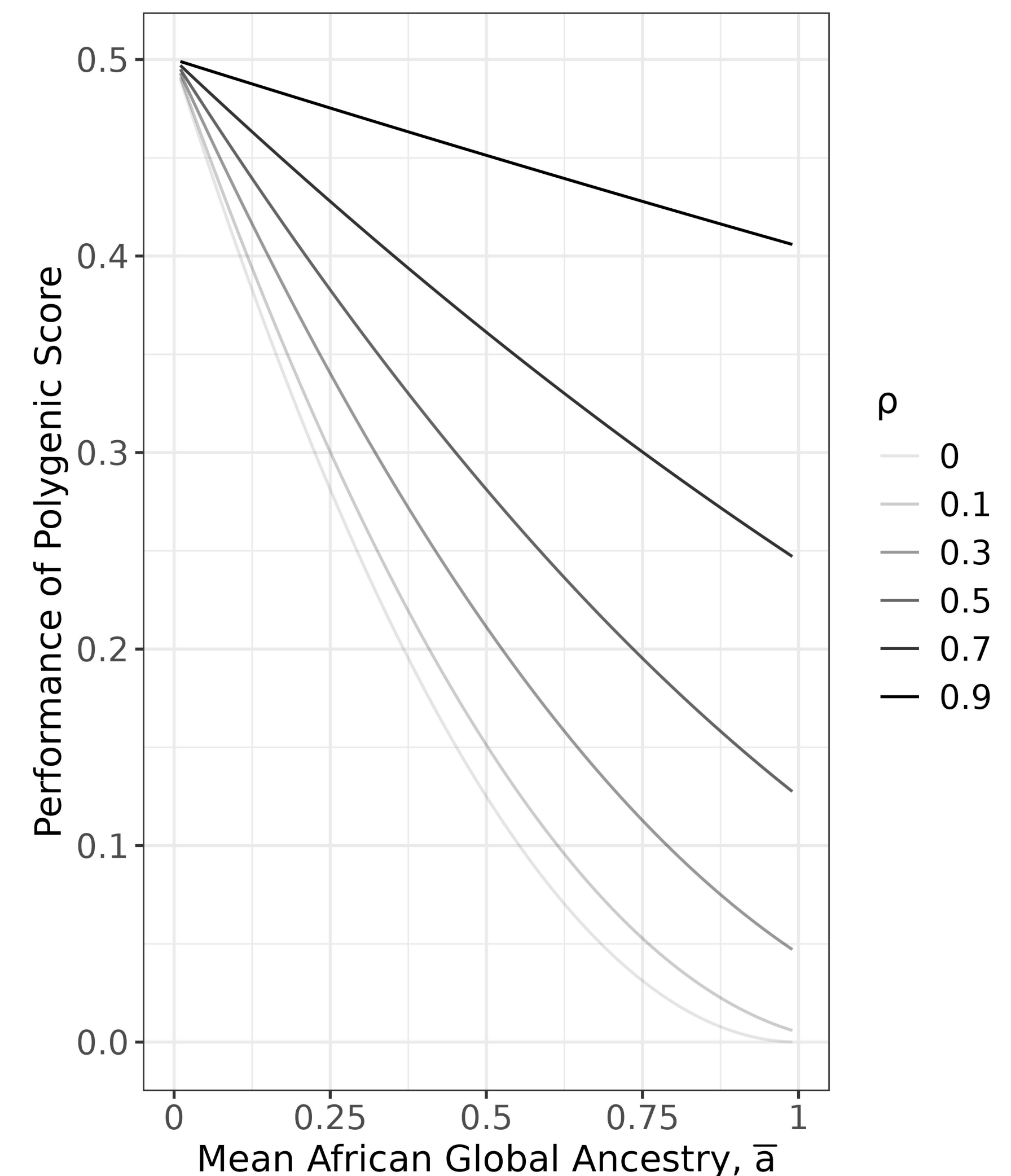
- Total PGS performance declines **quadratically** in global ancestry under either model

Correlation in causal effects between ancestries

$$\mathbb{E}[\text{cor}^2(\text{TotPGS}, y)] \approx r^2(1 - \bar{a} + \rho\bar{a})^2$$

(Mean African) Global ancestry

Total PGS Performance vs Global African Ancestry



# Causal Variants Known: Partial PGS differentiates the two models (but Total PGS does not)

- Partial PGS performance declines **cubically** in global ancestry under the global model, but declines **linearly** under the local model

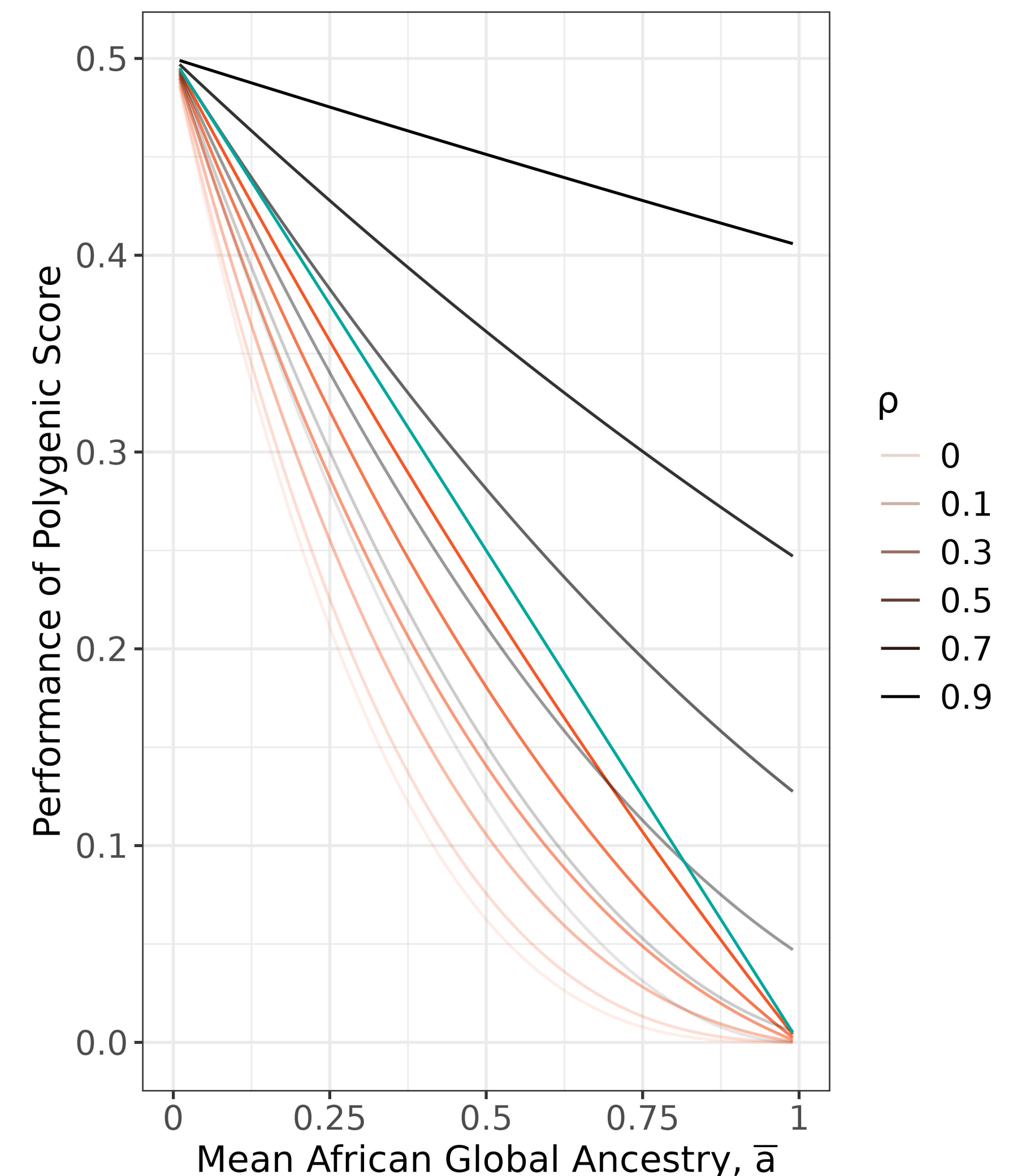
**Global Model:**  $\mathbb{E}[\text{cor}^2(\text{ParPGS}, y)] \approx r^2(1 - \bar{a})(1 - \bar{a} + \rho\bar{a})^2$

**Local Model:**  $\mathbb{E}[\text{cor}^2(\text{ParPGS}, y)] \approx r^2(1 - \bar{a})$

- Total PGS performance declines **quadratically** in global ancestry under either model

$$\mathbb{E}[\text{cor}^2(\text{TotPGS}, y)] \approx r^2(1 - \bar{a} + \rho\bar{a})^2$$

PGS Performance vs Global African Ancestry

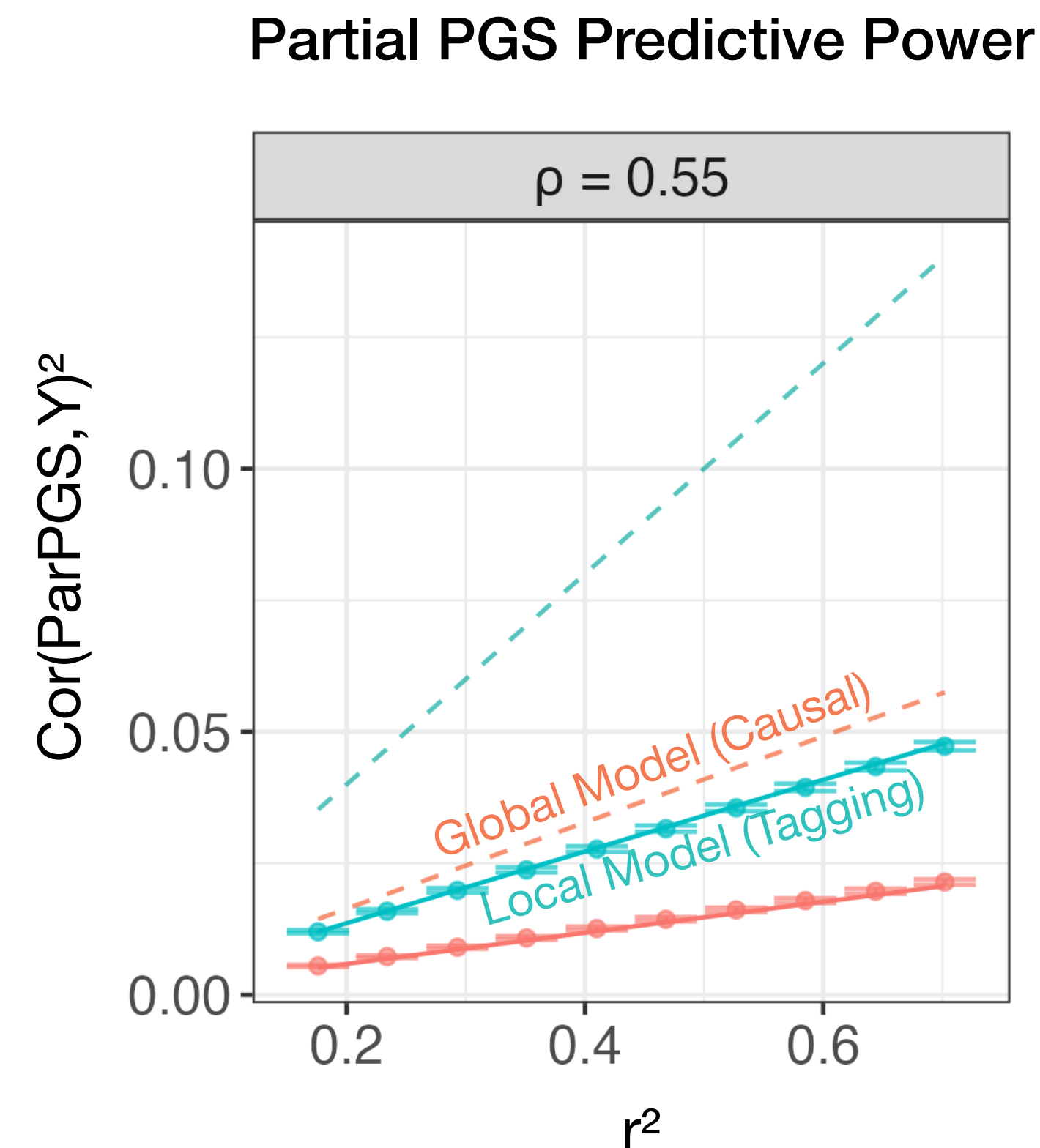




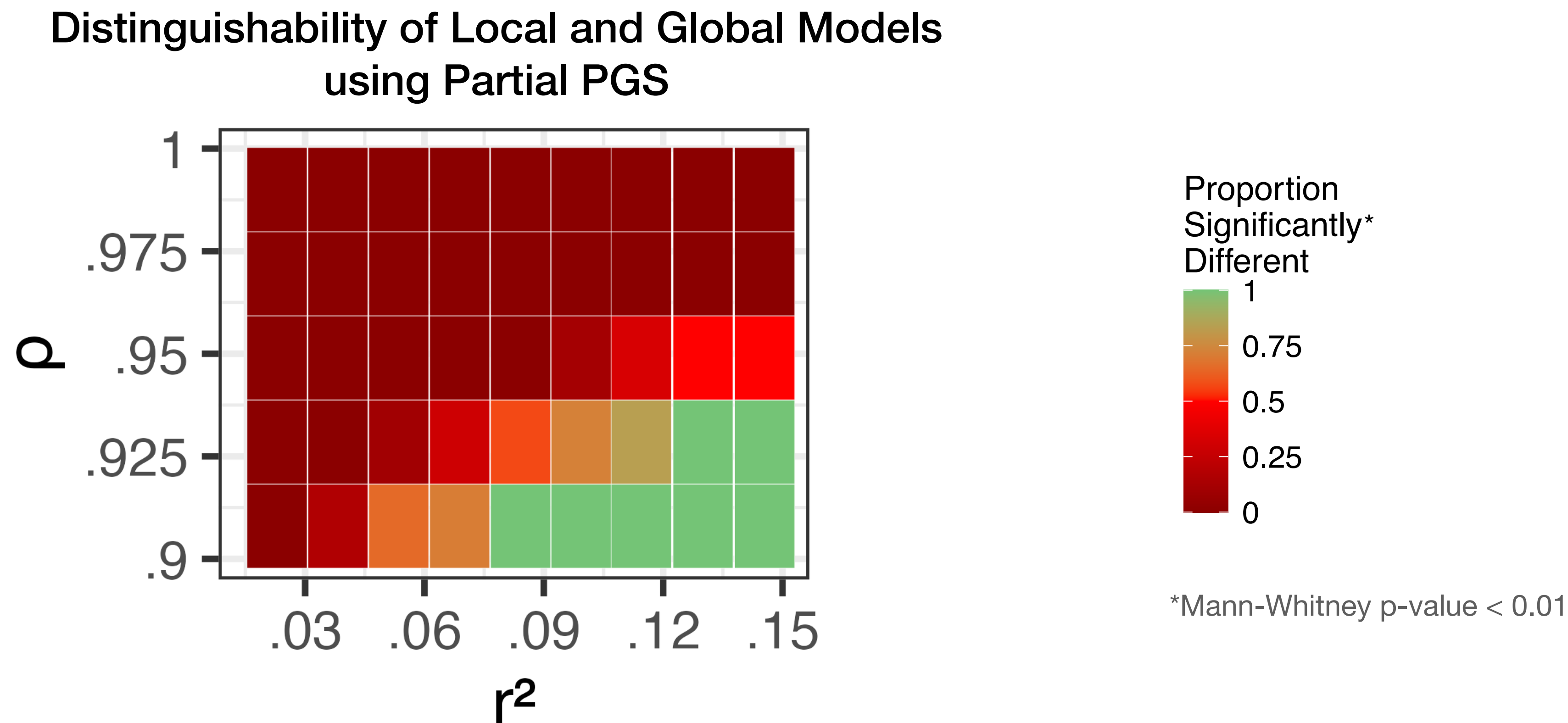
# Causal Variants Unknown: Heterogeneity in LD and allele frequencies hinders differentiation of models

Both  $\mathbb{E}[\text{cor}(TotPGS, y)^2]$  and  $\mathbb{E}[\text{cor}(ParPGS, y)^2]$  depend on causal-tagging LD and causal allele frequencies

LD and causal AF heterogeneities may produce differences in the two models that resemble analytical differences



# High causal effect correlation also hinders distinguishability of local and global models in general



# Summary of Q2

## Causal Variants Known (Ideal)

- Can differentiate local and global models using ParPGS

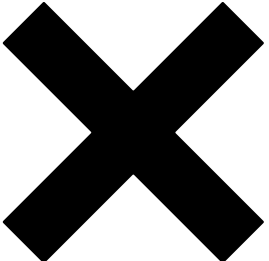
## Causal Variants Unknown (Realistic)

- Unknown differences in LD patterns and allele frequencies hinder differentiation

## High $\rho$ ?

- Difficult to differentiate local and global models

Example: Height

Local Model		
Global Model	$\text{Cor}(\text{ParPGS}, Y)^2$ cubic in global ancestry	
	Low $\rho$	High $\rho$

Assuming all polygenic score variants are causal:

50% contribution of global model

# Conclusion

- Models of GxA interaction are consistent with:
  - poor cross-ancestry portability
  - high (average) causal effect similarity across ancestries
- Fine-mapping causal variants helps differentiate the two models in future work