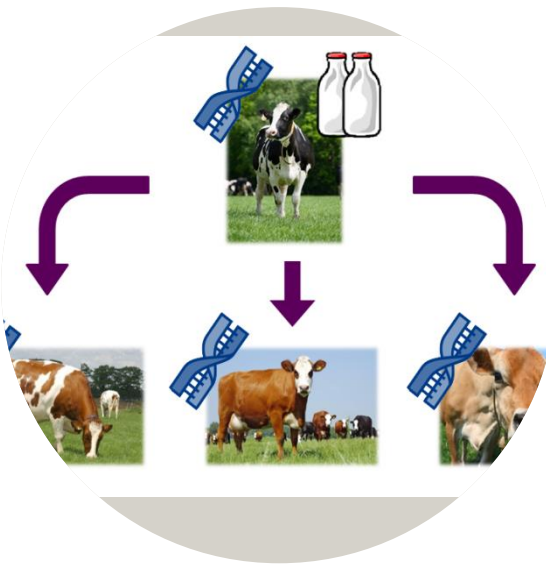


Genomic prediction using information from multiple populations

Yvonne Wientjes

INRA

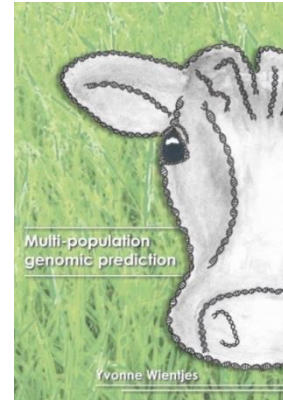
28 October 2019



Who am I?

PhD - Wageningen University & Research

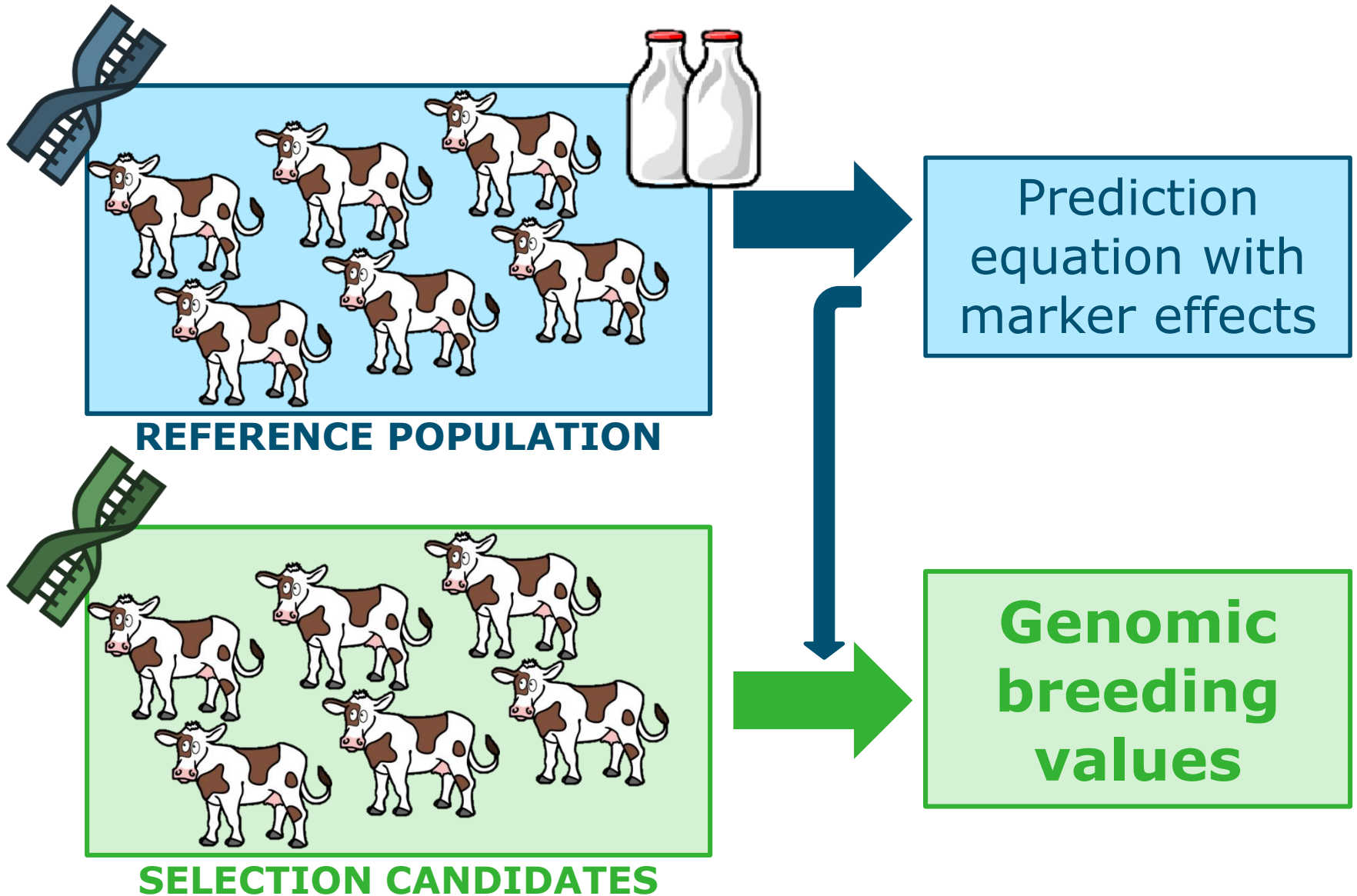
- Animal Breeding and Genomics
- 'Multi-population genomic prediction'



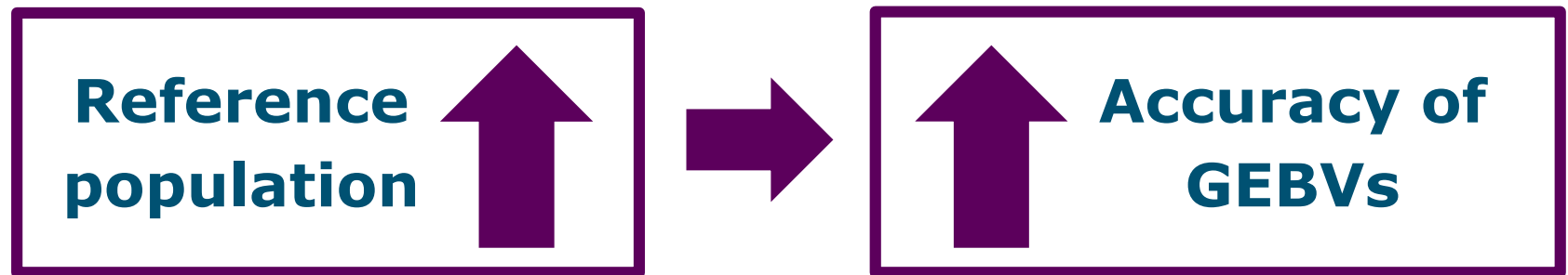
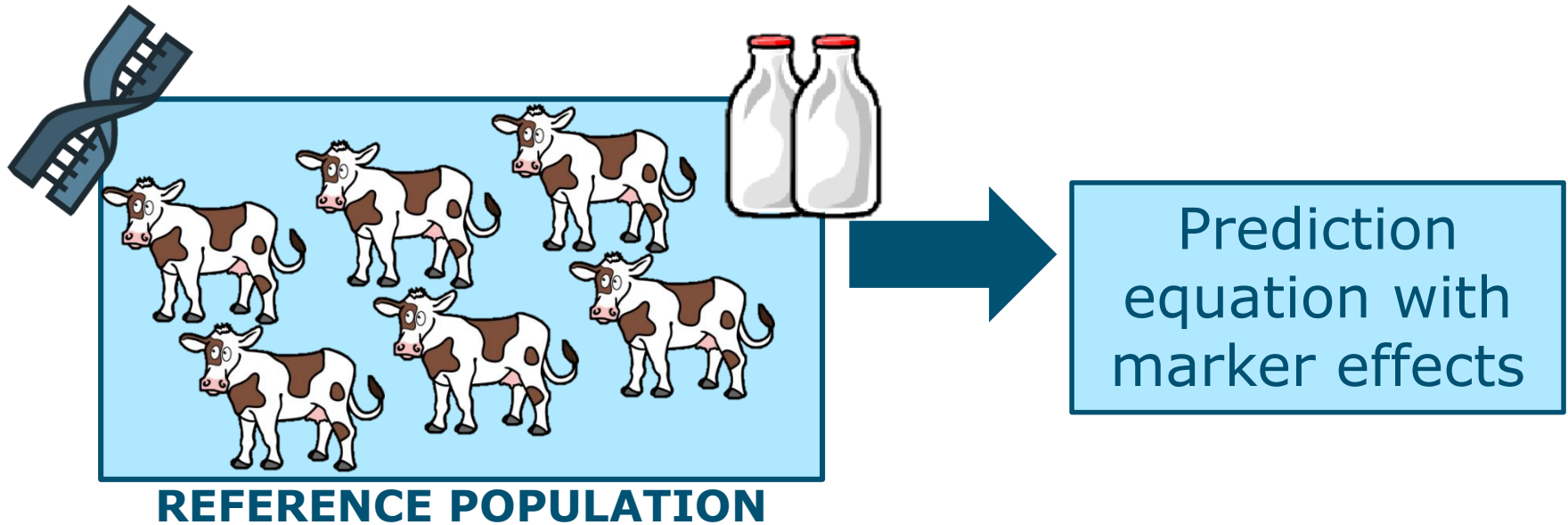
Postdoc - Wageningen University & Research

- Genomic prediction for crossbred performance
- Long-term effects of Genomic Selection

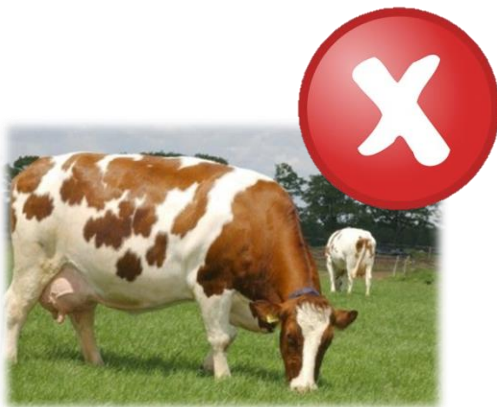
Genomic selection



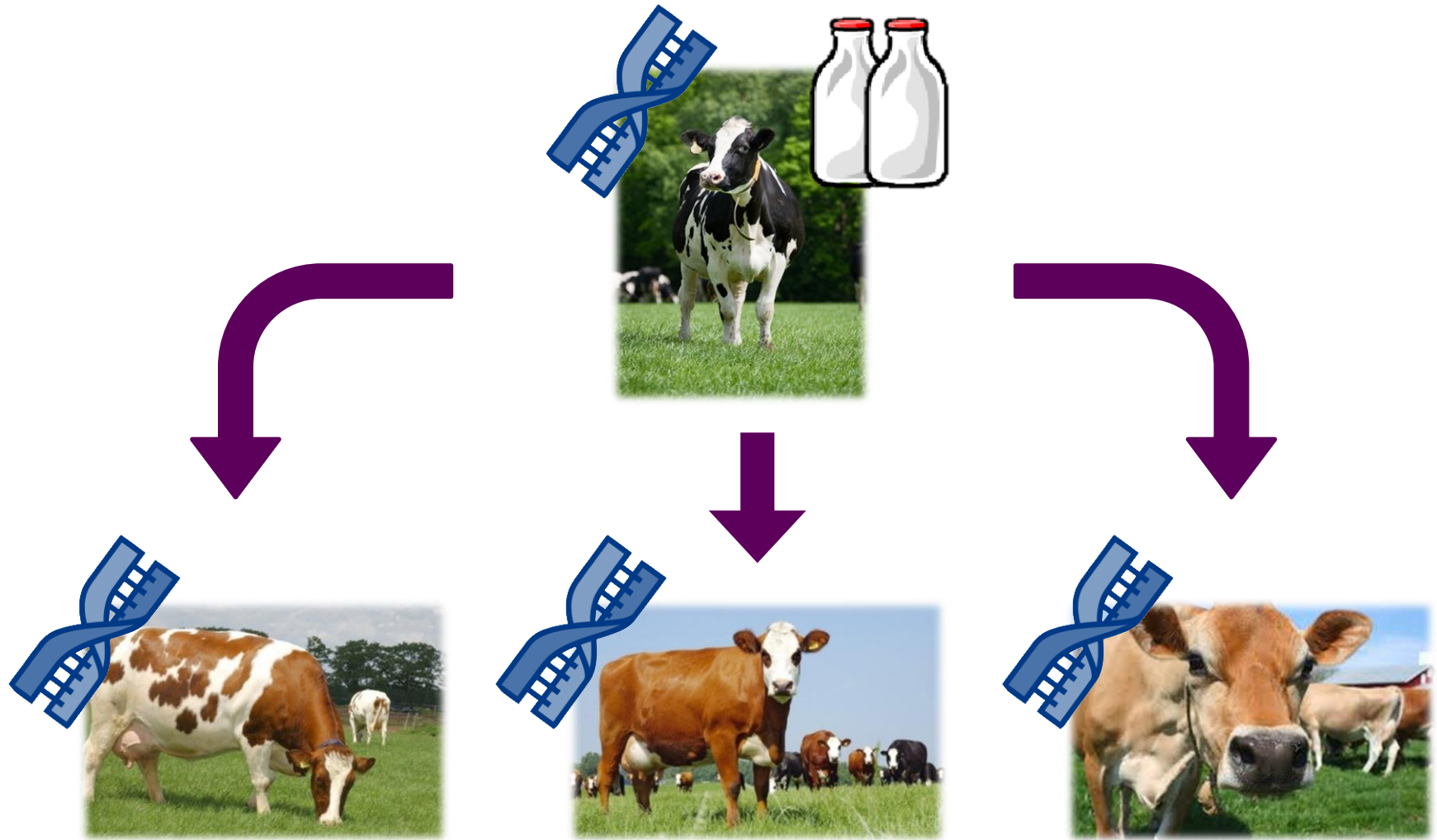
Reference population



Size of reference population



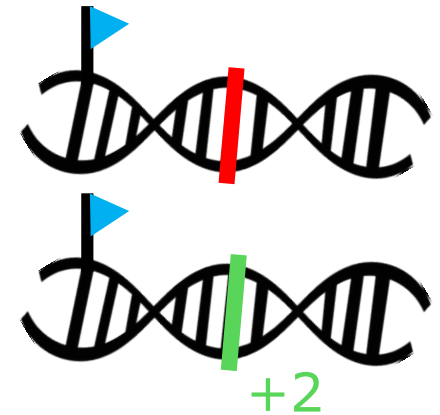
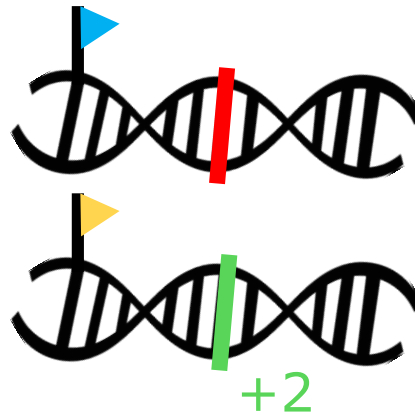
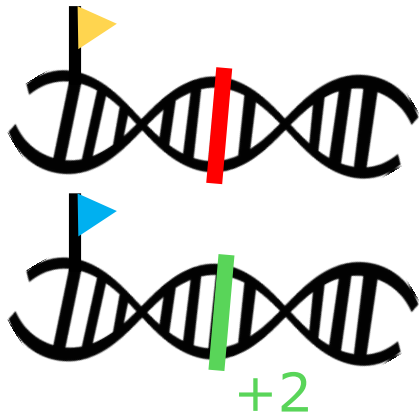
Use information across populations?



Differences between populations

- Linkage disequilibrium between markers and causal variants

Differences between populations in LD



Differences between populations

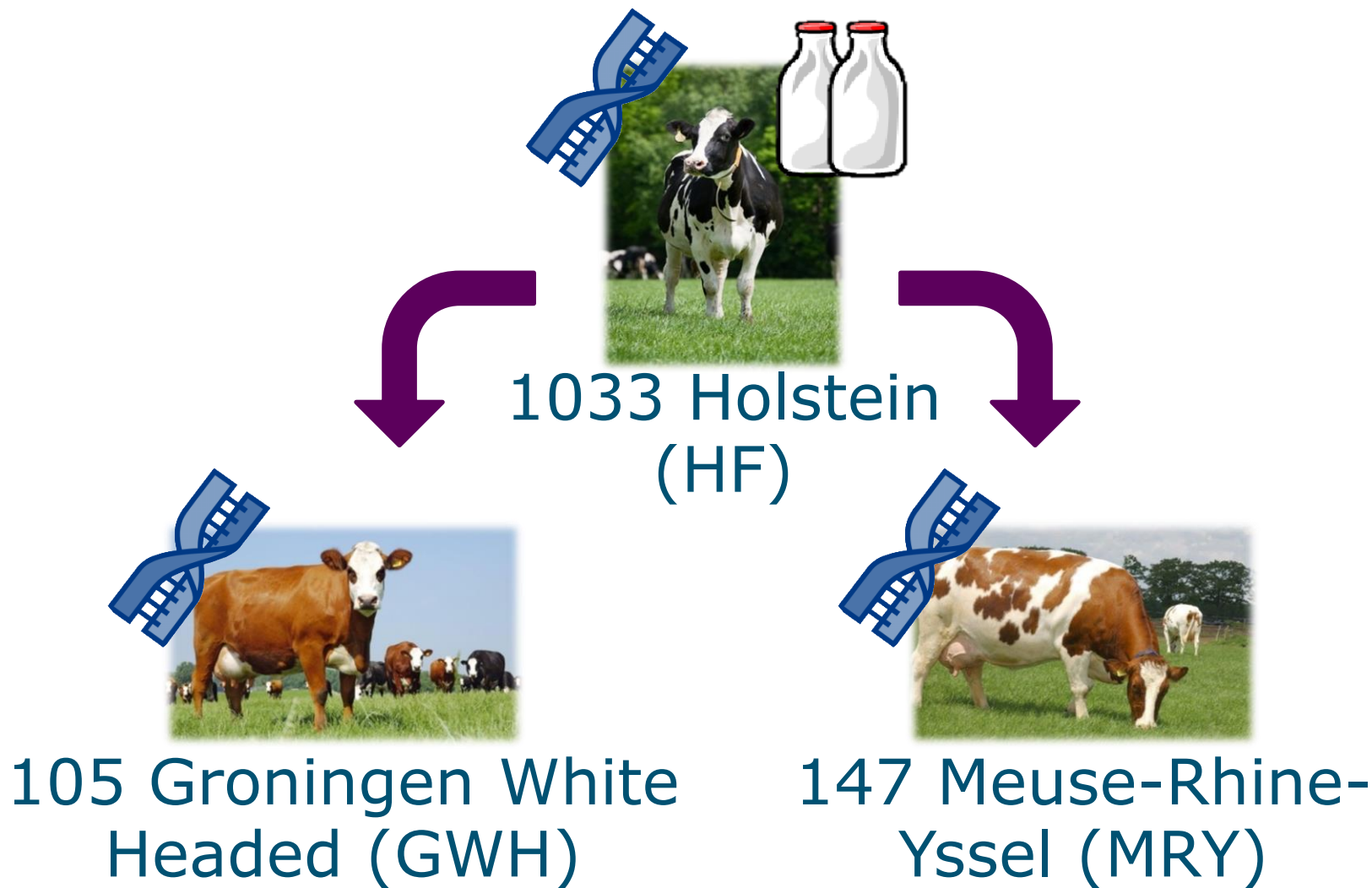
- Linkage disequilibrium between markers and causal variants
 - Allele frequencies of causal variants
 - Effects of causal variants
 - Environment different
 - Non-additive effects
- *Genetic correlation between populations*

Differences between populations

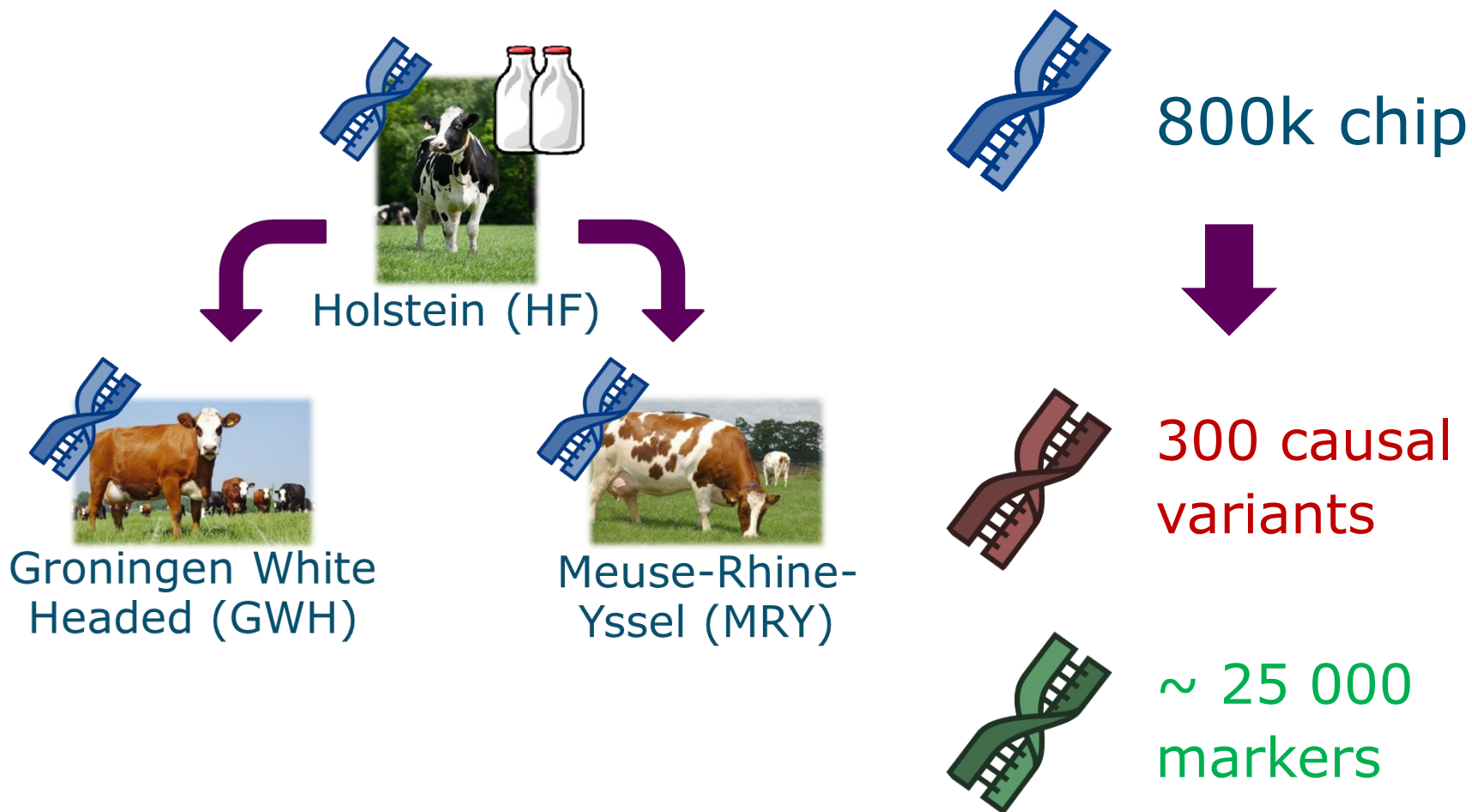
- Linkage disequilibrium between markers and causal variants
- Allele frequencies of causal variants
- Effects of causal variants
 - Environment different
 - Non-additive effects

→ *Genetic correlation between populations*
- Close family relationships are absent

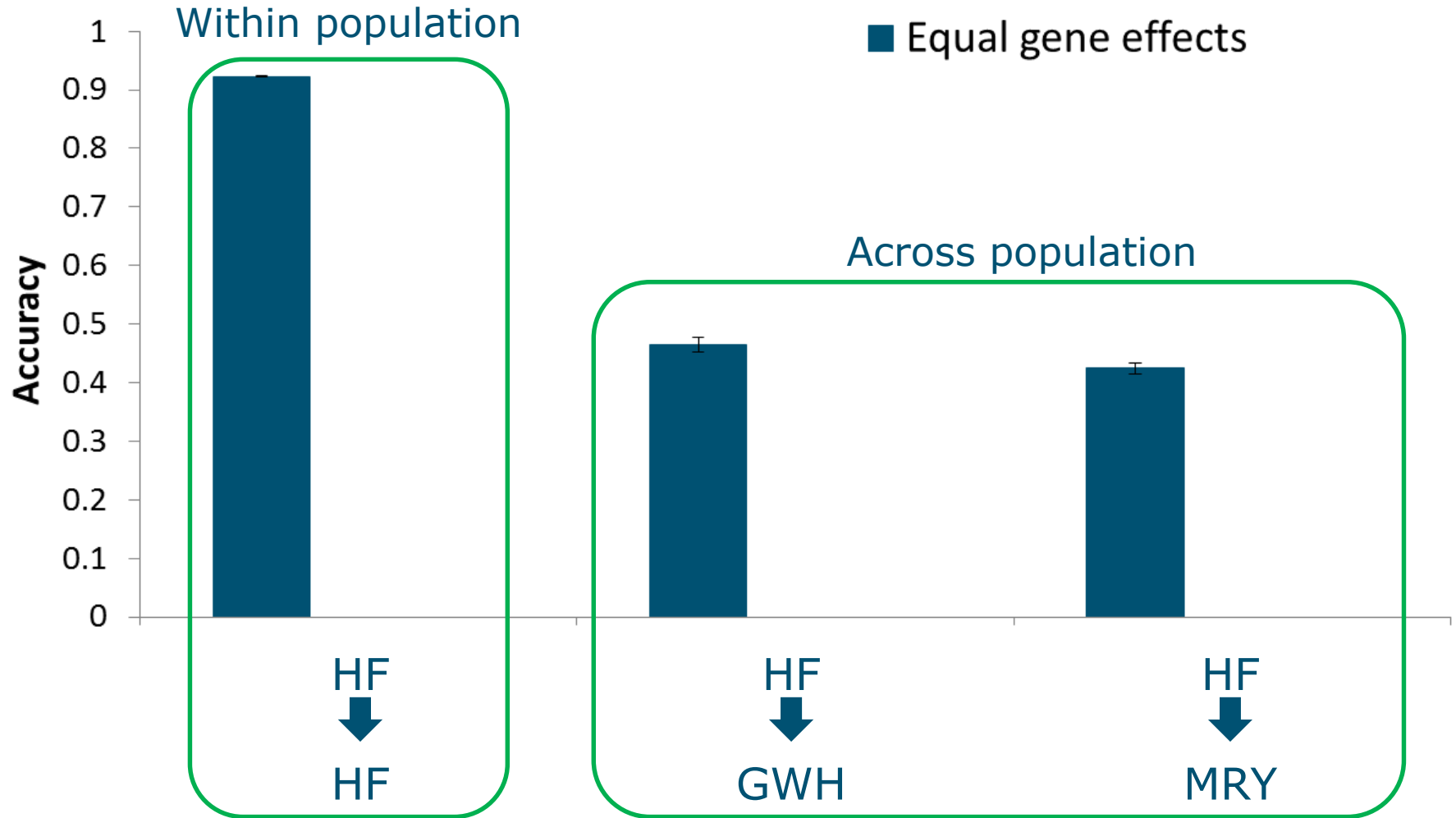
Across-population genomic prediction



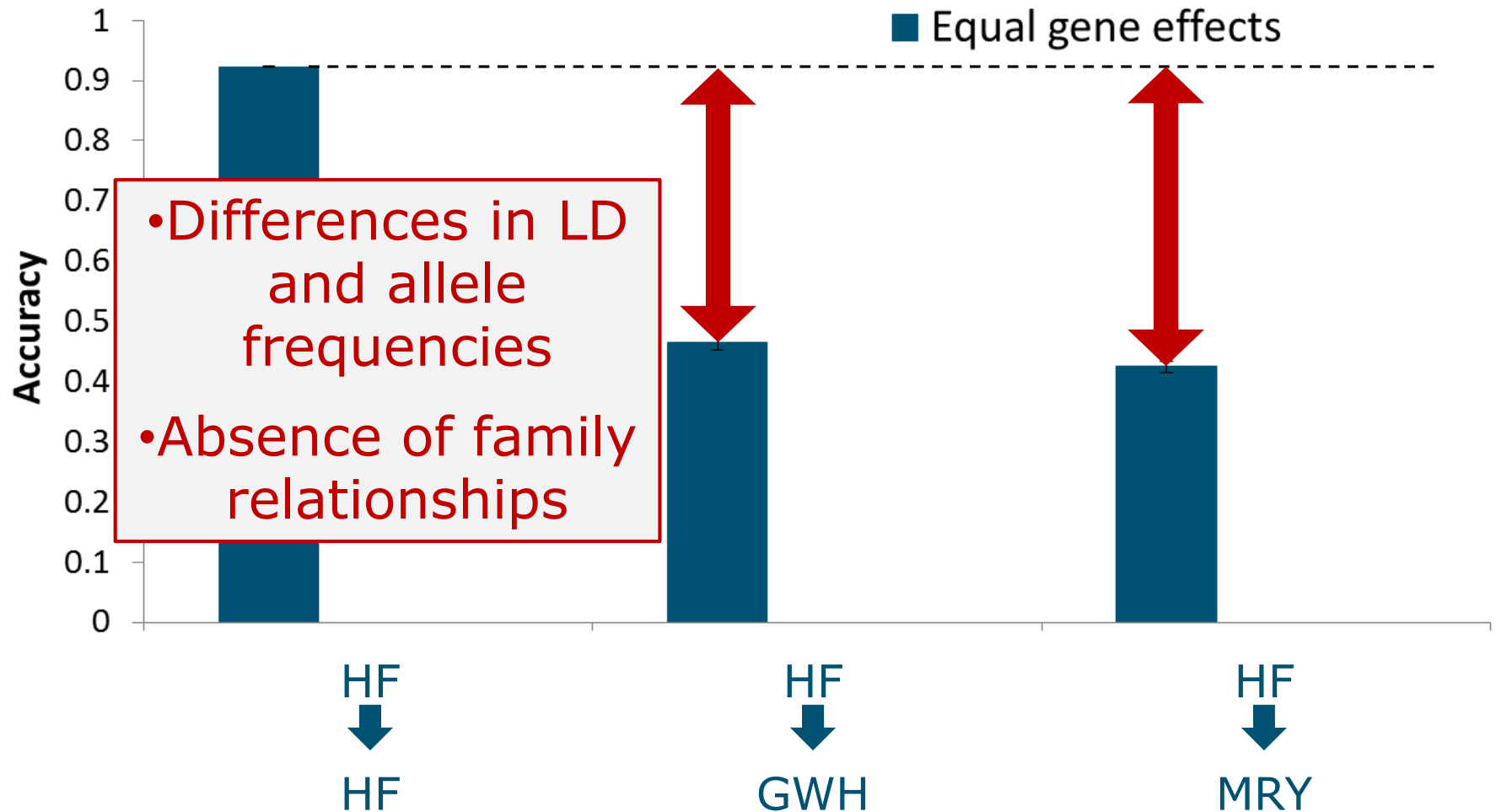
Across-population genomic prediction



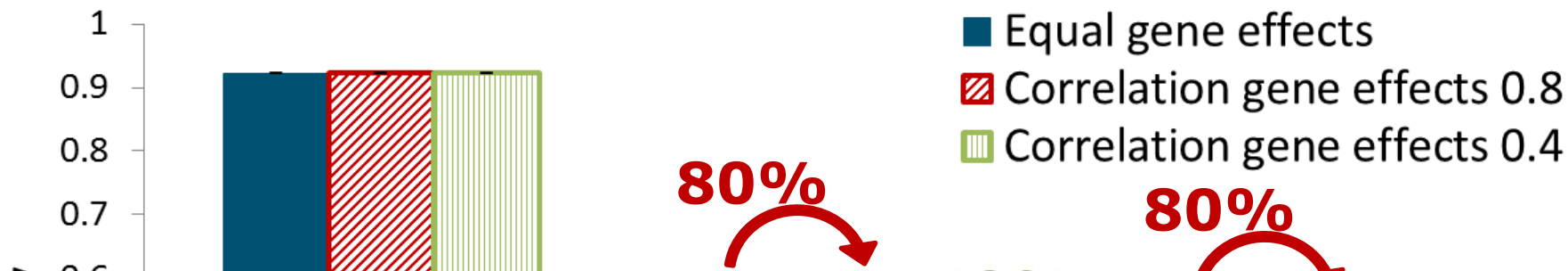
Across-population genomic prediction



Across-population genomic prediction



Across-population genomic prediction



- ❖ Accuracy of across-population genomic prediction is much lower than within-population genomic prediction
- ❖ Differences in effects of causal variants reduce accuracy

Effect of properties causal variants



Holstein



Jersey

Causal variants

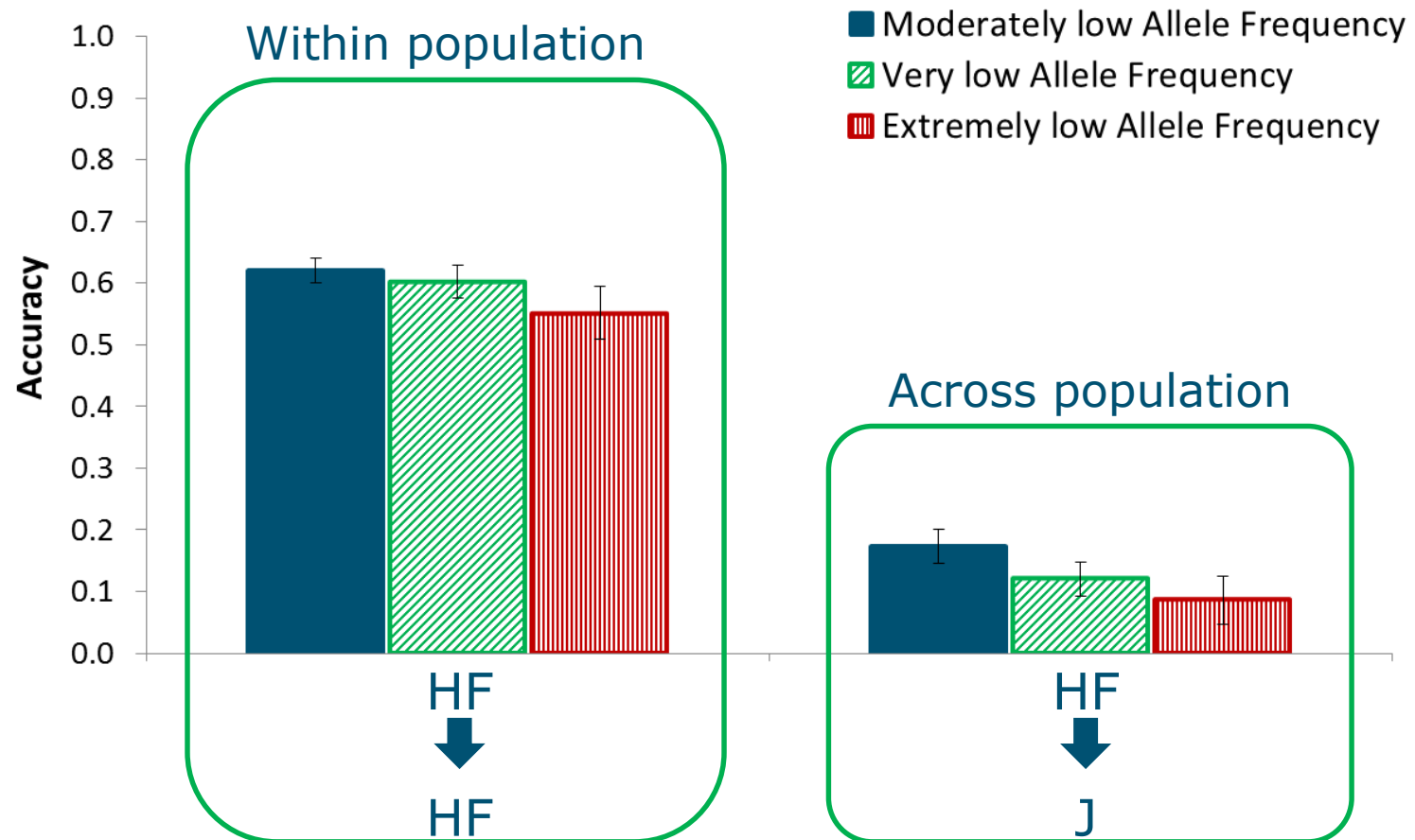
Moderately low
allele frequency
(~1 in 8)

Very low allele
frequency
(~1 in 13)

Extremely low
allele frequency
(~1 in 63)

Accuracies of predicting breeding values

Gene effects randomly sampled

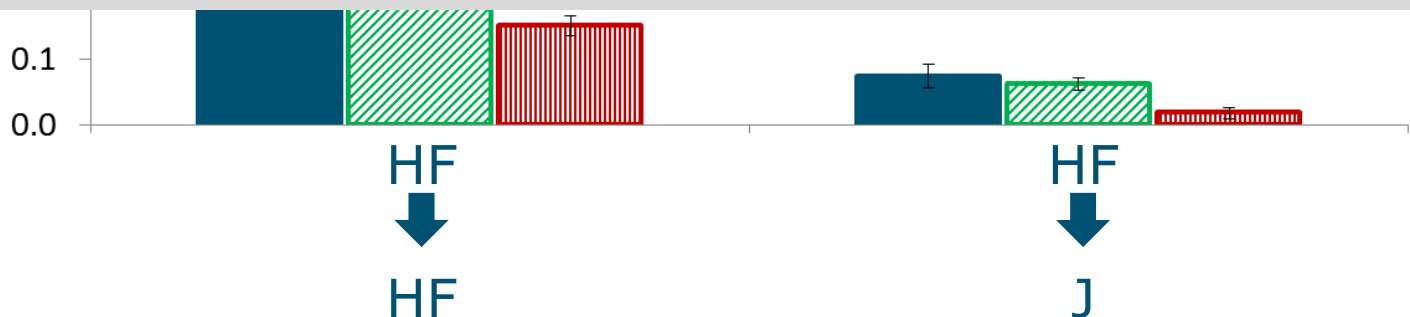


Accuracies of predicting breeding values

Larger effects for genes with lower frequency



Accuracy of genomic prediction depends on genetic architecture of trait



Multi-population genomic prediction

Accuracy of across-population genomic prediction is low...

, but what if we combine populations in one reference population?

Multi-population genomic prediction

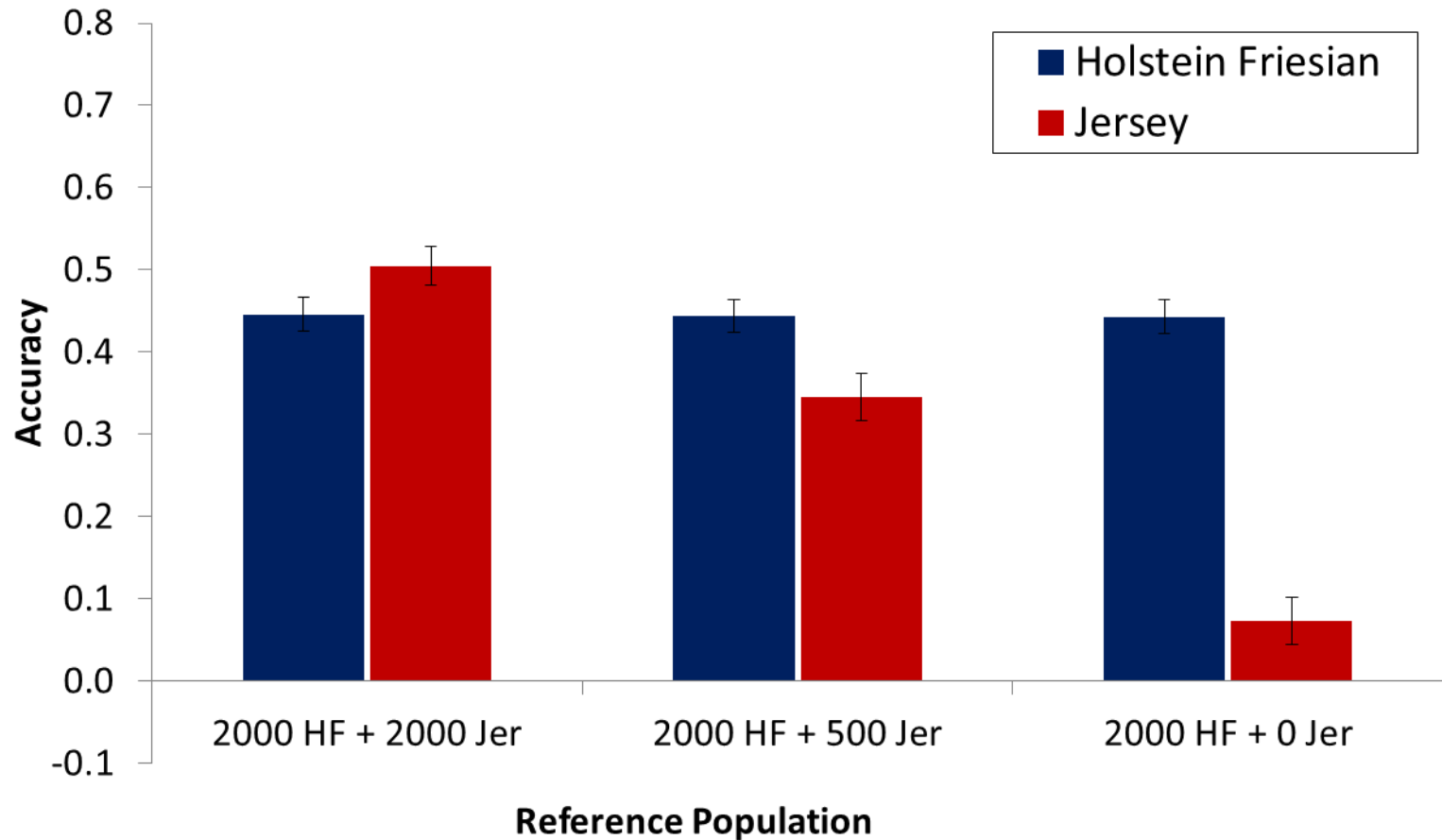


Prediction equation with marker effects

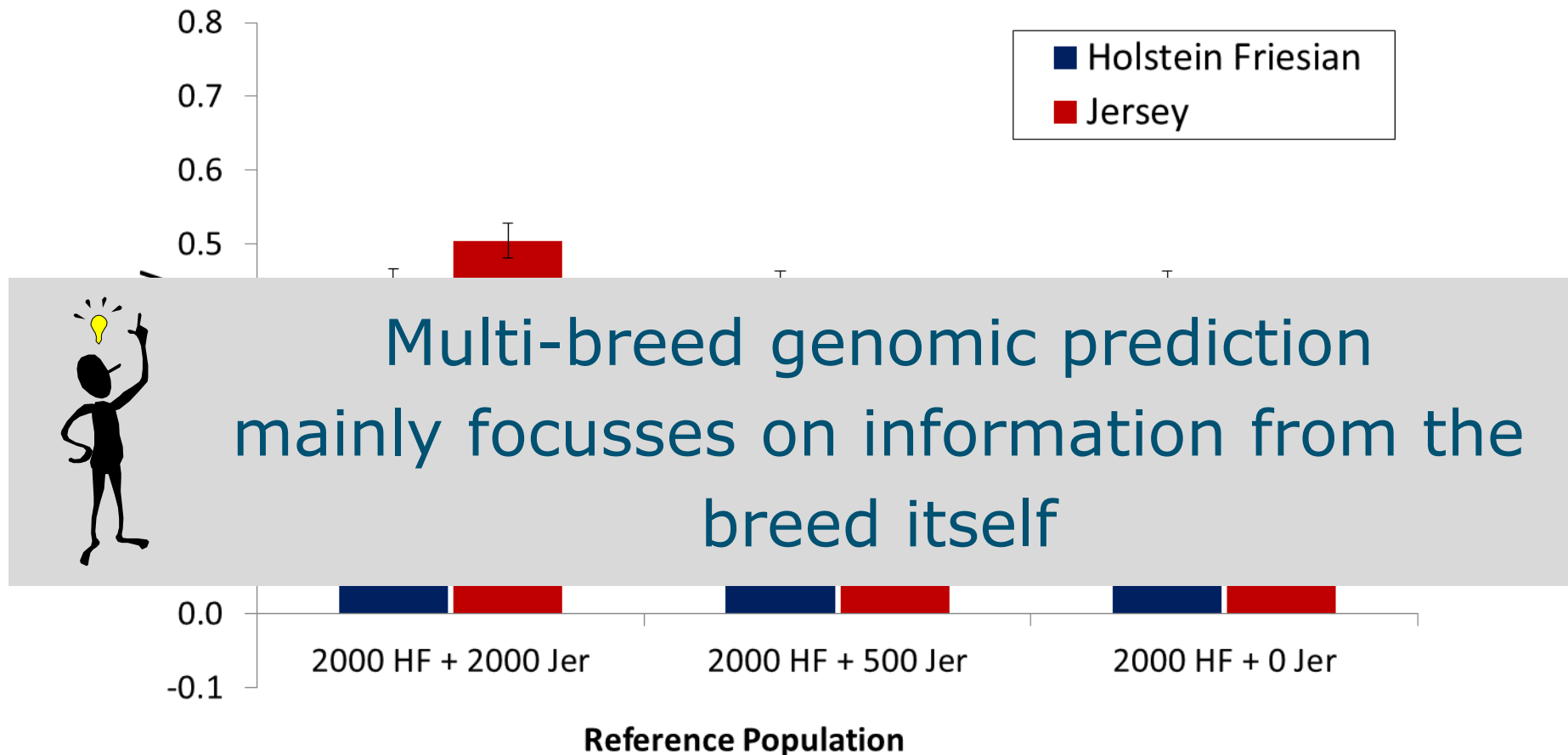


Genomic breeding values

Multi-population genomic prediction



Multi-population genomic prediction

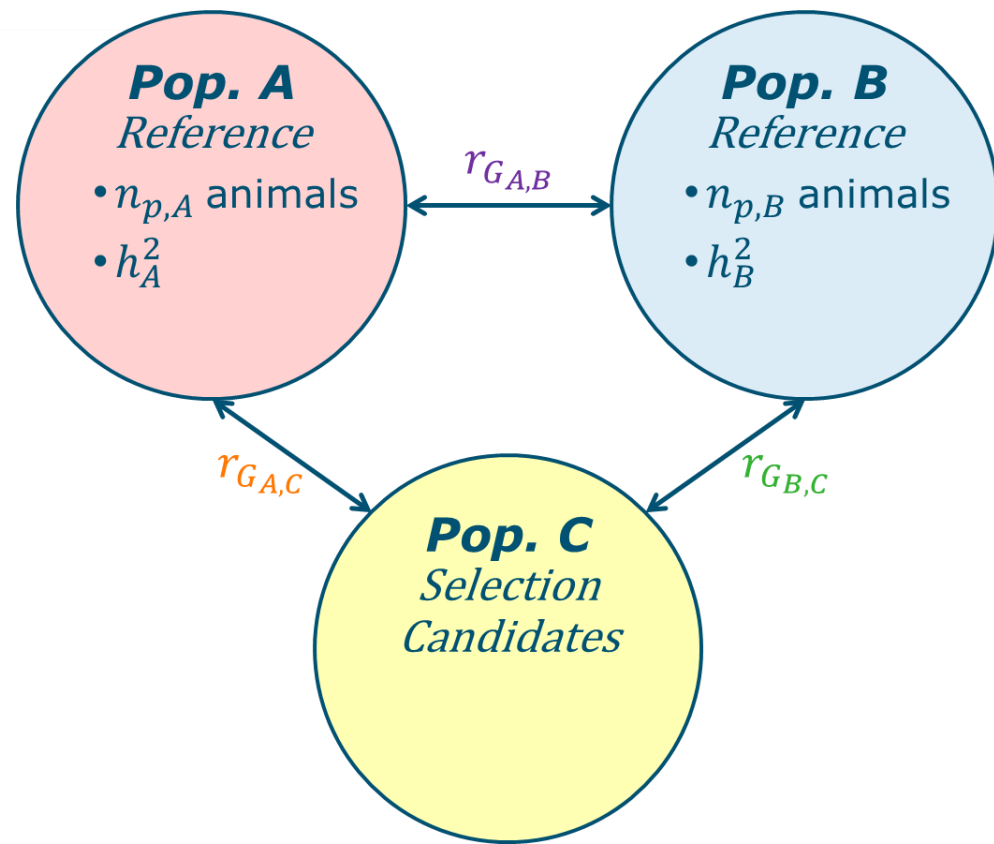


Can we predict accuracy?

Important when designing breeding programs

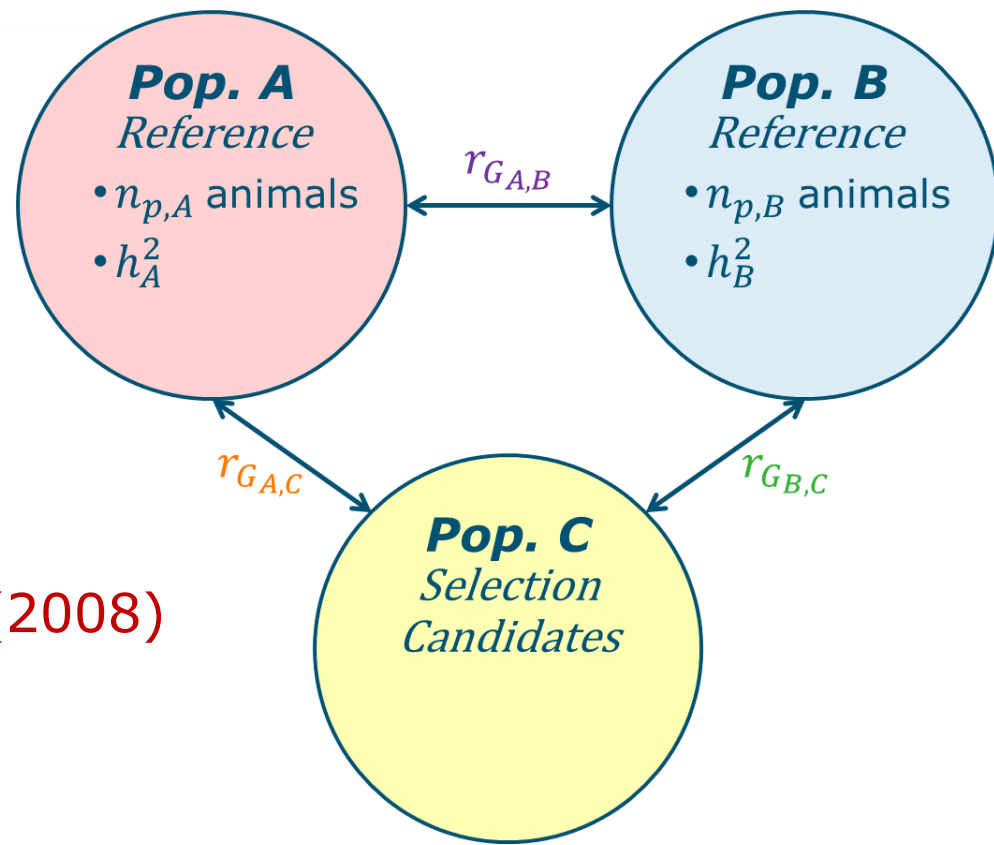
- Who to genotype?

Derived equation



$$r_{GEBV_{A+B,C}} = \sqrt{\begin{bmatrix} r_{G_{A,C}} \sqrt{\frac{h_A^2}{M_{e_{A,C}}}} & r_{G_{B,C}} \sqrt{\frac{h_B^2}{M_{e_{B,C}}}} \end{bmatrix} \begin{bmatrix} \frac{h_A^2}{M_{e_{A,C}}} + \frac{1}{n_{p,A}} & r_{G_{A,B}} \frac{\sqrt{h_A^2 h_B^2}}{\sqrt{M_{e_{A,C}} M_{e_{B,C}}}} \\ r_{G_{A,B}} \frac{\sqrt{h_A^2 h_B^2}}{\sqrt{M_{e_{A,C}} M_{e_{B,C}}}} & \frac{h_B^2}{M_{e_{B,C}}} + \frac{1}{n_{p,B}} \end{bmatrix}^{-1} \begin{bmatrix} r_{G_{A,C}} \sqrt{\frac{h_A^2}{M_{e_{A,C}}}} \\ r_{G_{B,C}} \sqrt{\frac{h_B^2}{M_{e_{B,C}}}} \end{bmatrix}}$$

Derived equation



1 population in reference
 → Equal to Daetwyler *et al.* (2008)

$$r_{GEBV_{A+B,C}} = \sqrt{\begin{bmatrix} r_{G_{A,C}} \sqrt{\frac{h_A^2}{M_{e_{A,C}}}} & r_{G_{B,C}} \sqrt{\frac{h_B^2}{M_{e_{B,C}}}} \end{bmatrix} \begin{bmatrix} \frac{h_A^2}{M_{e_{A,C}}} + \frac{1}{n_{p,A}} & \frac{r_{G_{A,B}} \sqrt{h_A^2 h_B^2}}{\sqrt{M_{e_{A,C}} M_{e_{B,C}}}} \\ r_{G_{A,B}} \frac{\sqrt{h_A^2 h_B^2}}{\sqrt{M_{e_{A,C}} M_{e_{B,C}}}} & \frac{h_B^2}{M_{e_{B,C}}} + \frac{1}{n_{p,B}} \end{bmatrix}^{-1} \begin{bmatrix} r_{G_{A,C}} \sqrt{\frac{h_A^2}{M_{e_{A,C}}}} \\ r_{G_{B,C}} \sqrt{\frac{h_B^2}{M_{e_{B,C}}}} \end{bmatrix}}$$

M_e across populations

'Effective number of estimated effects'

Pop. B  $M_e = 4$

Pop. C  $M_e = 5$

Pop. B - C  $M_e = 18$

M_e across populations

'Effective number of estimated effects'

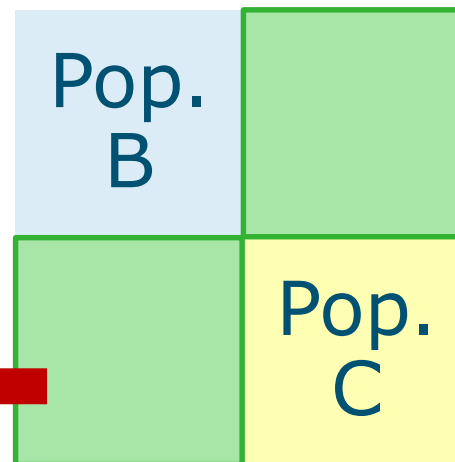
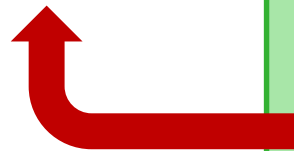
Pop. B  $M_e = 4$

Pop. C  $M_e = 5$

Pop. B - C  $M_e = 18$

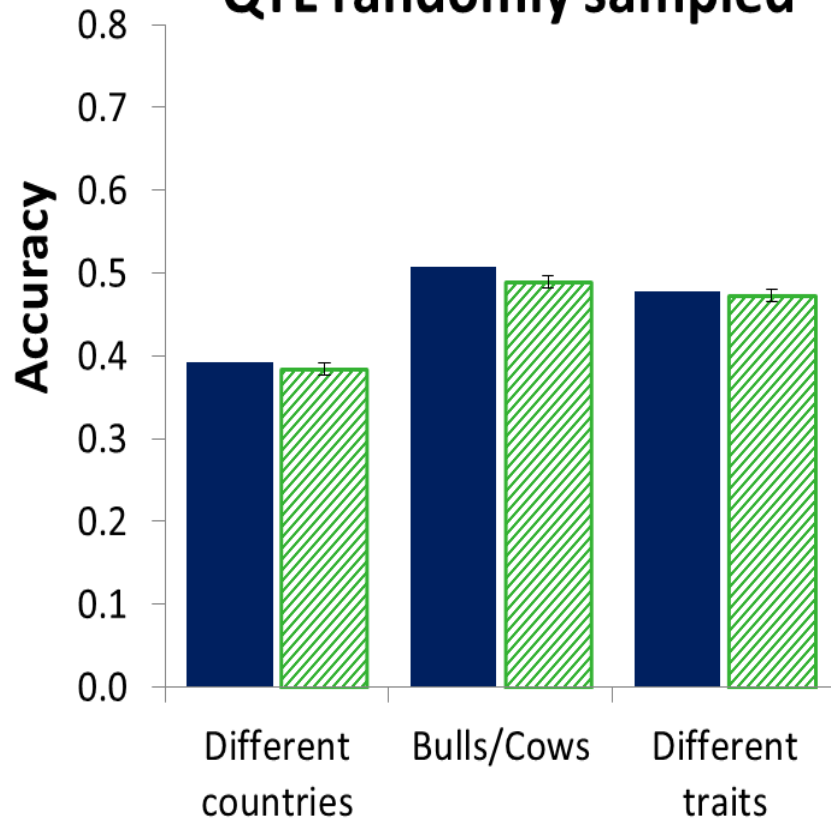
Relationship matrix

$$M_{e_{B,C}} = \frac{1}{Var(\mathbf{G})}$$



Results of validation

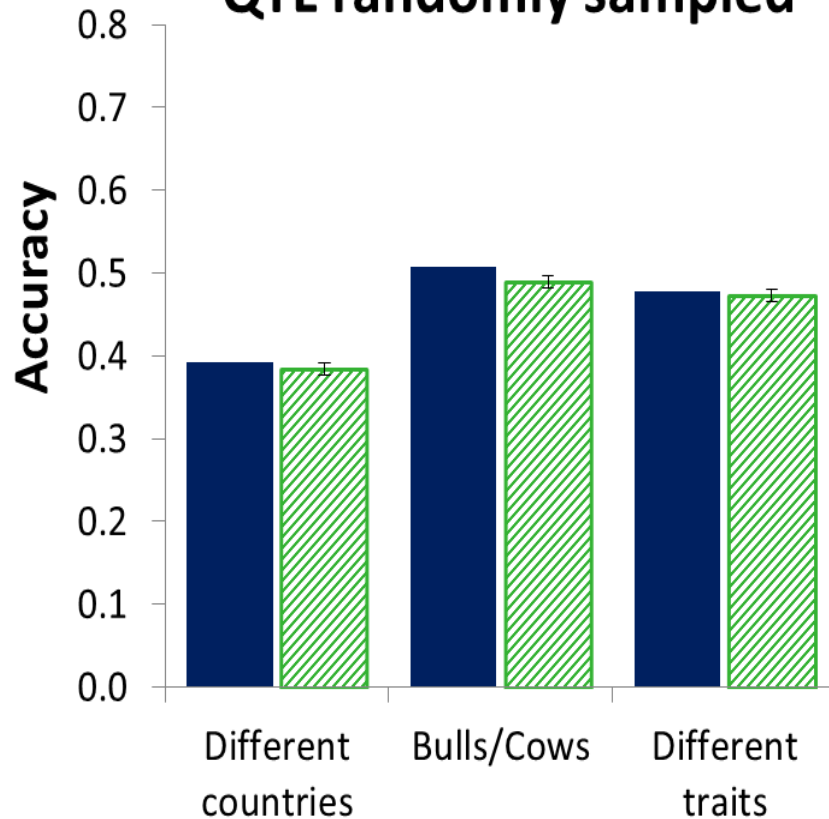
QTL randomly sampled



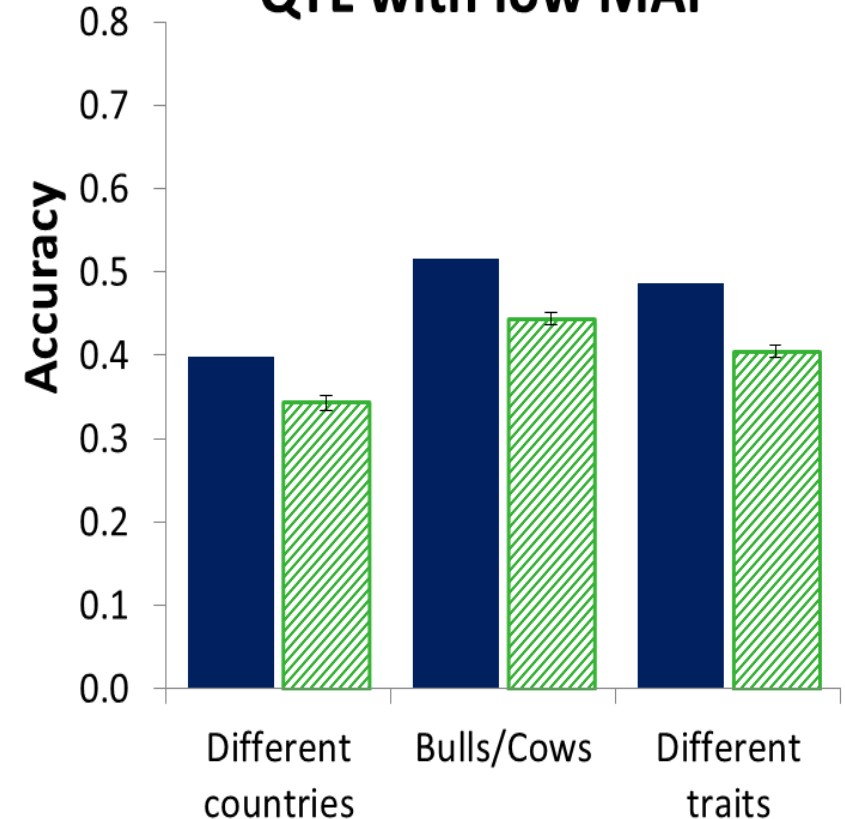
■ Predicted accuracy ▨ Empirical accuracy

Results of validation

QTL randomly sampled



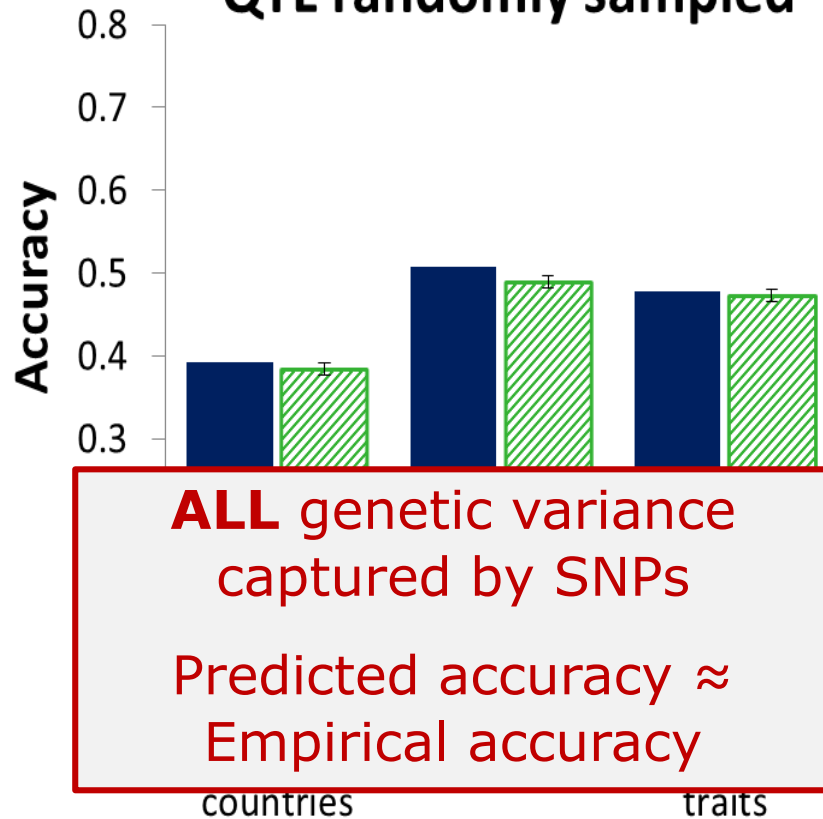
QTL with low MAF



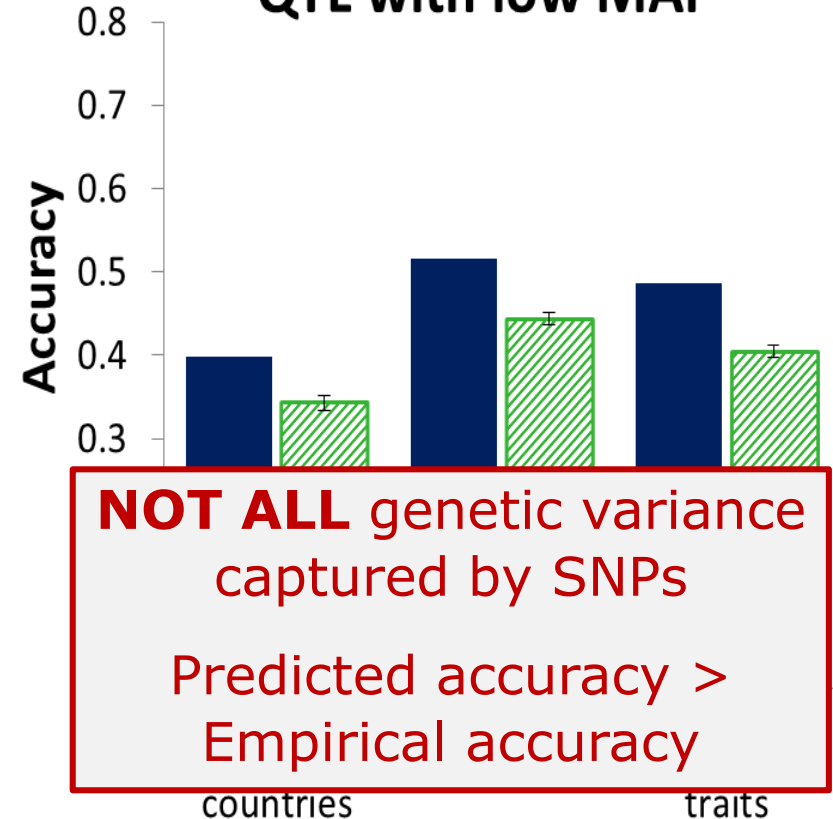
■ Predicted accuracy ▨ Empirical accuracy

Results of validation

QTL randomly sampled



QTL with low MAF



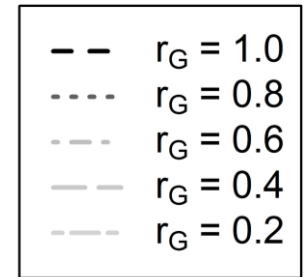
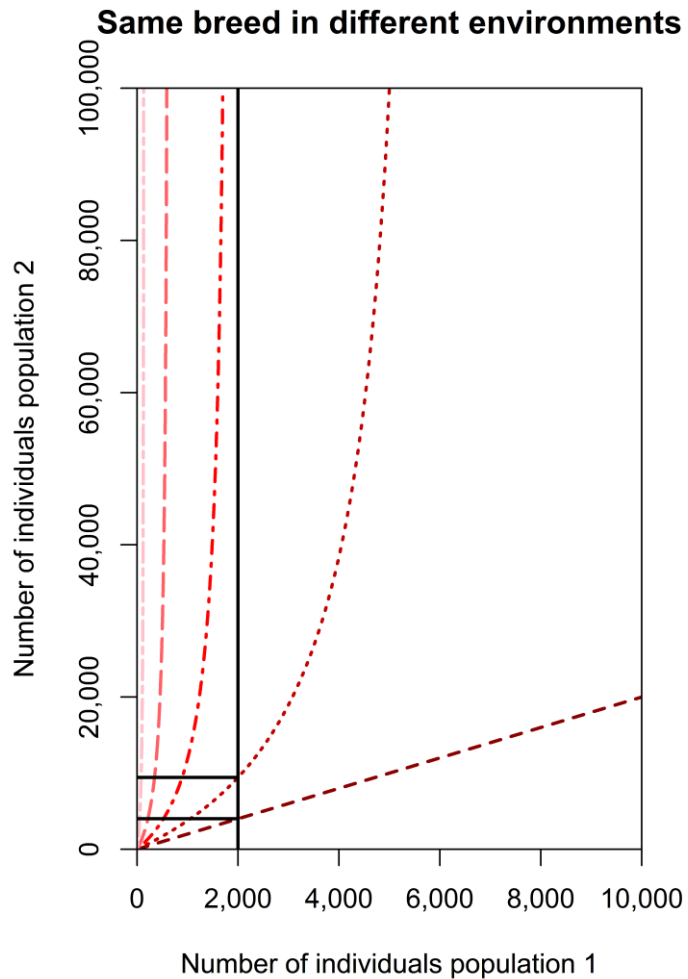
■ Predicted accuracy ■ Empirical accuracy

Combining populations

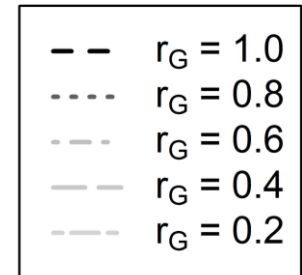
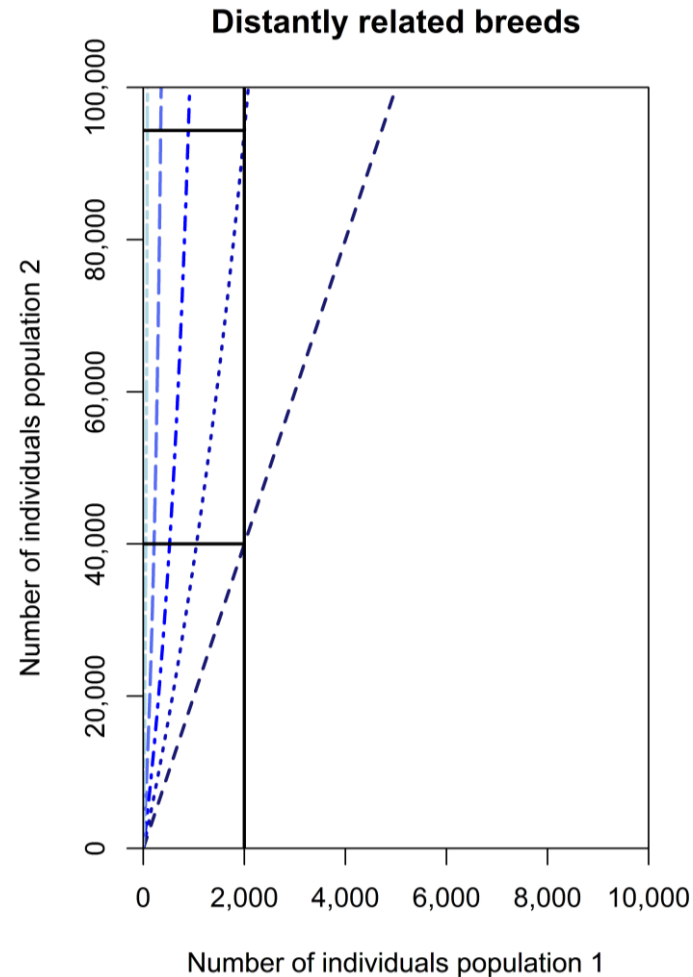
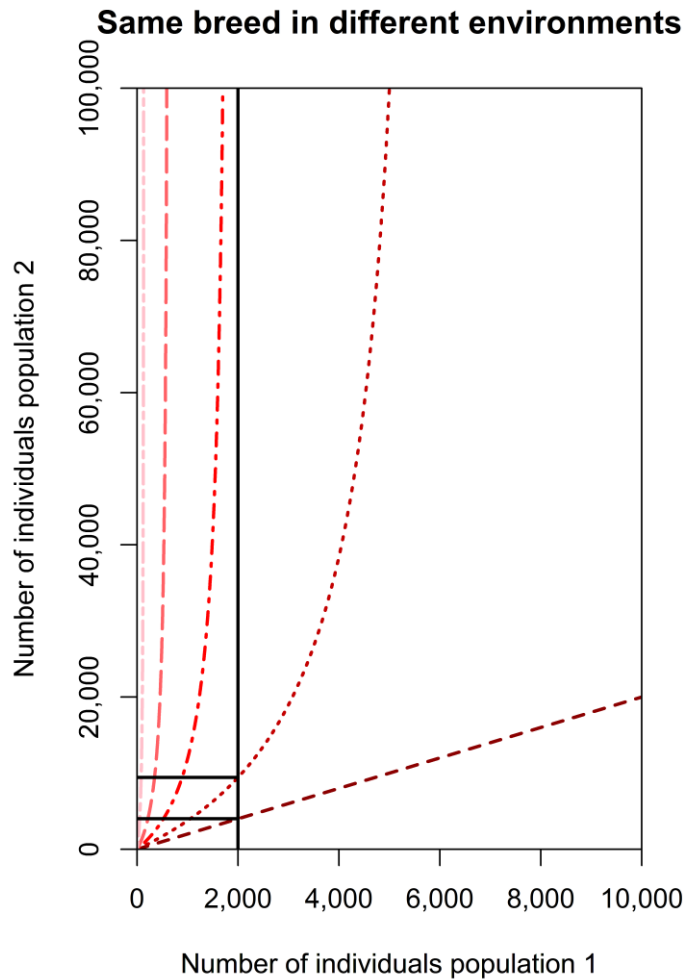
Populations of:

- Different countries
- Different breeds
- Bulls/Cows
- Measured for different traits
- Different generations
- ...

Combining populations?



Combining populations?



Combining populations?

Beneficial, when:

- Closely related populations
- Population itself is small
- A large number of individuals is added

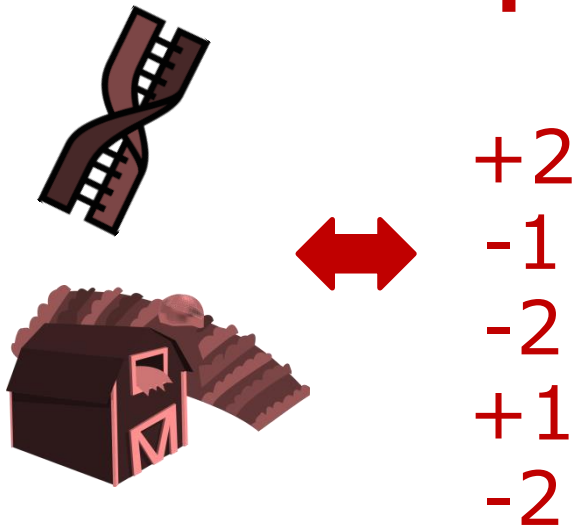


**Depends on genetic correlation
between populations!!**

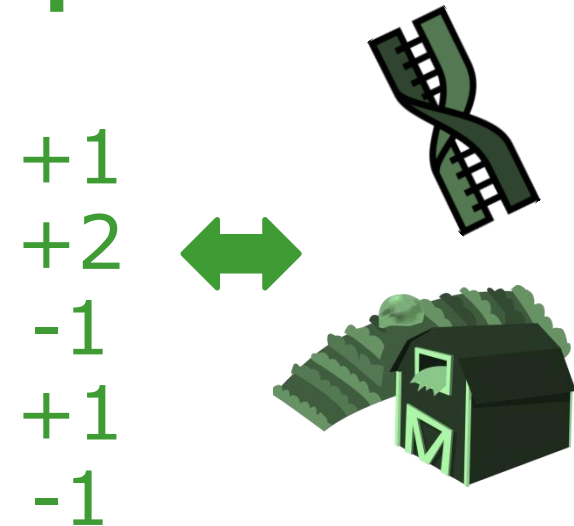
Genetic correlation between populations (r_g)

'Correlation between breeding values of two individuals with the same genotype in the two populations'

Pop. 1



Pop. 2



Estimation of r_g

- Multi-trait model
 - Each population different trait
- REML

Requires relationships between populations

Relationships between populations

Theoretically defined relationships:

$$\mathbf{G}_{\text{New}} = \begin{bmatrix} \frac{\mathbf{z}_1 \mathbf{z}'_1}{\sum 2p_{1i}(1 - p_{1i})} & \frac{\mathbf{z}_1 \mathbf{z}'_2}{\sqrt{\sum 2p_{1i}(1 - p_{1i})} \sqrt{\sum 2p_{2i}(1 - p_{2i})}} \\ \frac{\mathbf{z}_2 \mathbf{z}'_1}{\sqrt{\sum 2p_{1i}(1 - p_{1i})} \sqrt{\sum 2p_{2i}(1 - p_{2i})}} & \frac{\mathbf{z}_2 \mathbf{z}'_2}{\sum 2p_{2i}(1 - p_{2i})} \end{bmatrix}$$

Relationships between populations

Theoretically defined relationships:

$$\mathbf{G}_{\text{New}} = \begin{bmatrix} \frac{\mathbf{Z}_1 \mathbf{Z}'_1}{\sum 2p_{1i}(1 - p_{1i})} & \frac{\mathbf{Z}_1 \mathbf{Z}'_2}{\sqrt{\sum 2p_{1i}(1 - p_{1i})} \sqrt{\sum 2p_{2i}(1 - p_{2i})}} \\ \frac{\mathbf{Z}_2 \mathbf{Z}'_1}{\sqrt{\sum 2p_{1i}(1 - p_{1i})} \sqrt{\sum 2p_{2i}(1 - p_{2i})}} & \frac{\mathbf{Z}_2 \mathbf{Z}'_2}{\sum 2p_{2i}(1 - p_{2i})} \end{bmatrix}$$

VanRaden – method 1

Relationships between populations

Theoretically defined relationships:

$$G_{\text{New}} = \begin{bmatrix} \frac{Z_1 Z'_1}{\sum 2p_{1i}(1-p_{1i})} & \frac{Z_1 Z'_2}{\sqrt{\sum 2p_{1i}(1-p_{1i})} \sqrt{\sum 2p_{2i}(1-p_{2i})}} \\ \frac{Z_2 Z'_1}{\sqrt{\sum 2p_{1i}(1-p_{1i})} \sqrt{\sum 2p_{2i}(1-p_{2i})}} & \frac{Z_2 Z'_2}{\sum 2p_{2i}(1-p_{2i})} \end{bmatrix}$$

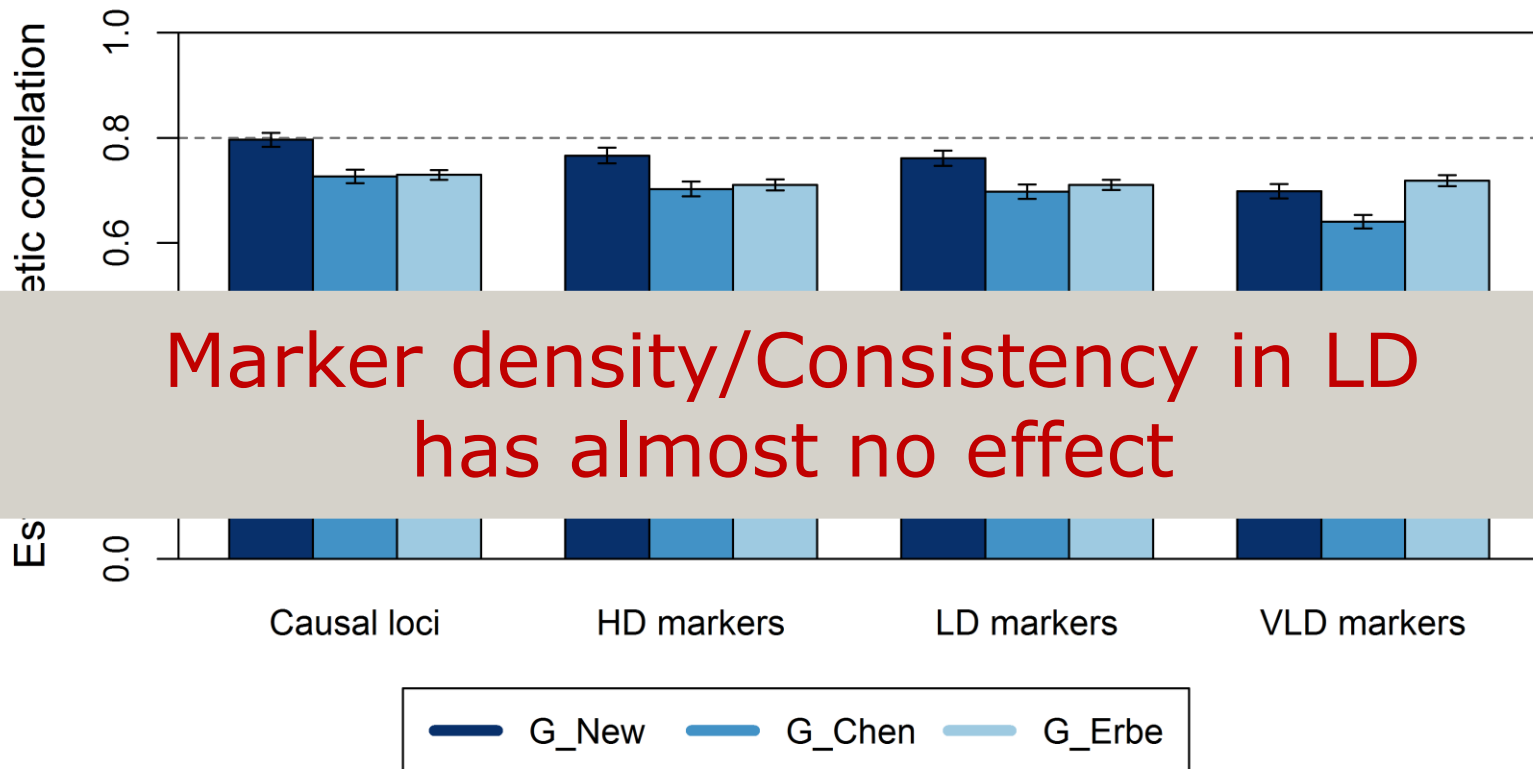
G_Chen: Scaling factor across populations:

$$\sum 2 \sqrt{p_{1i}(1-p_{1i})p_{2i}(1-p_{2i})}$$

G_Erbe: Base when populations separated

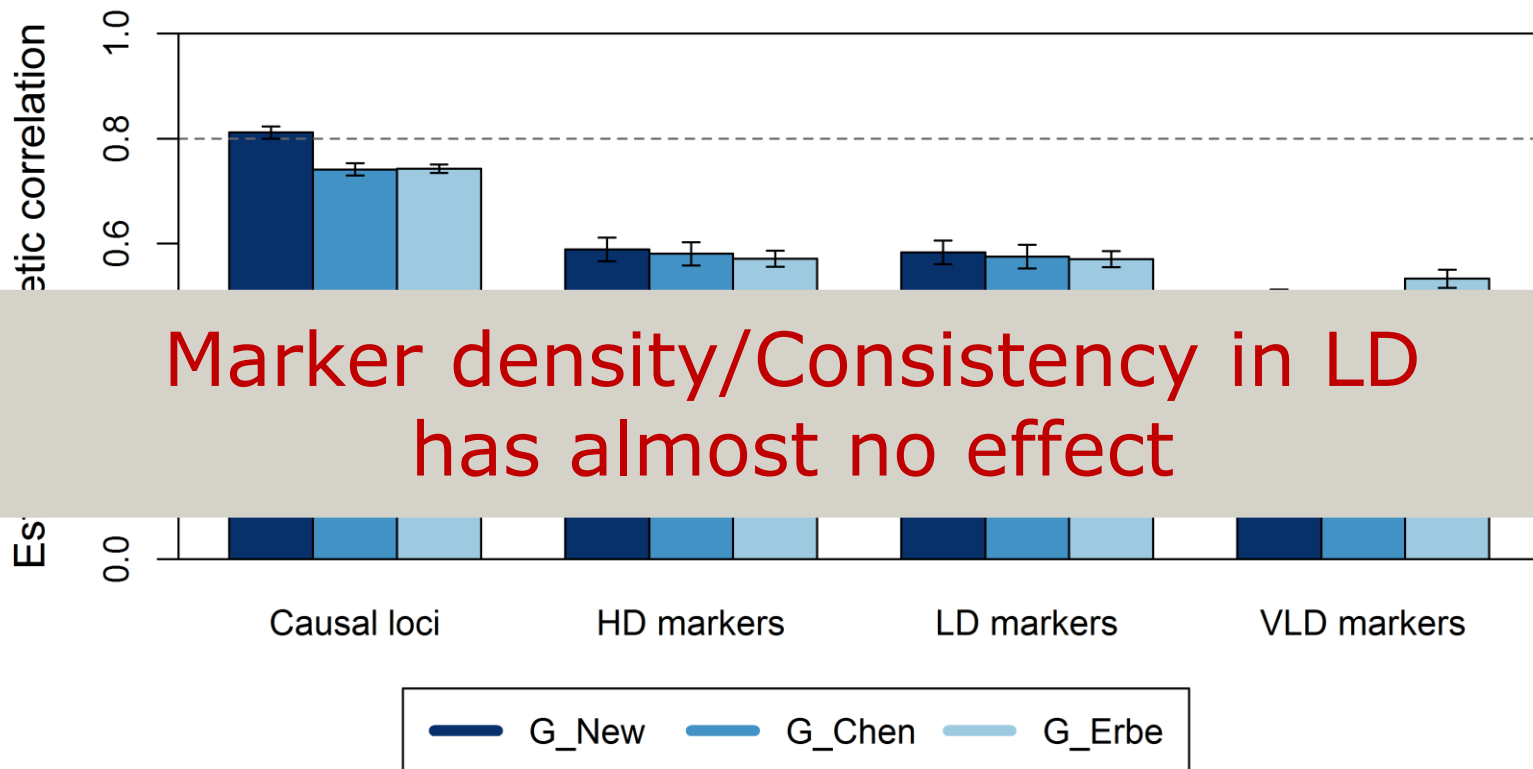
Estimated genetic correlation

Allele frequency differences between populations of causal loci REPRESENTED by markers



Estimated genetic correlation

Allele frequency differences between populations of causal loci NOT REPRESENTED by markers



Conclusion

Combining populations beneficial when:

- Closely related populations
- Population itself is small
- Many individuals are added

Genetic correlation between populations:

- Important parameter
- Unbiasedly estimated with **G**_New
 - markers represent properties causal variants



Acknowledgements



Mario Calus
Roel Veerkamp
Piter Bijma
Henk Bovenhuis
Pascal Duenk
Jeremie Vandenplas



Ben Hayes
Mike Goddard



Chris Schrooten

Financial support

