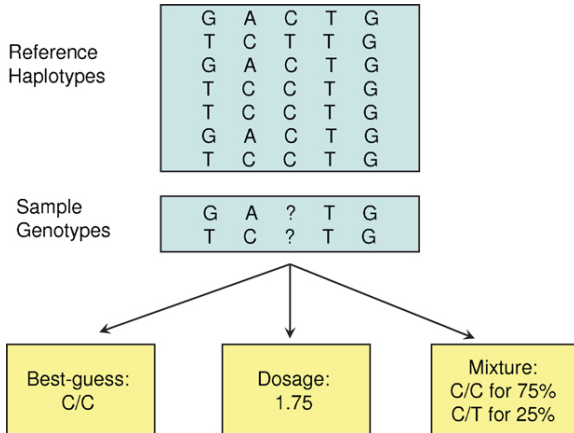# The Impact of Imputation Quality on Family-Based Analysis

Mahdi Mir (UCLA & SSGAC)
Alexander Young (UCLA & SSGAC)

Oct 2024

# Imputation from a Reference Panel



Hard Call = Best Guess     Dosages $= \mathbb{E}[\hat{G}]$

$0 \leq$ Info Score $\leq 1$     $\uparrow$ Imputation Quality $\iff \uparrow$ Info Score

Source: Zheng et al. (2011)

# Reference Based Imputation VS Mandelian Imputation

- Mandellian imputation as done in SNIPAR (Young et.al, Nature Genetics 2022) is very different from reference based imputation.

- Reference based imputation is not taking into account the relationships between the individuals.
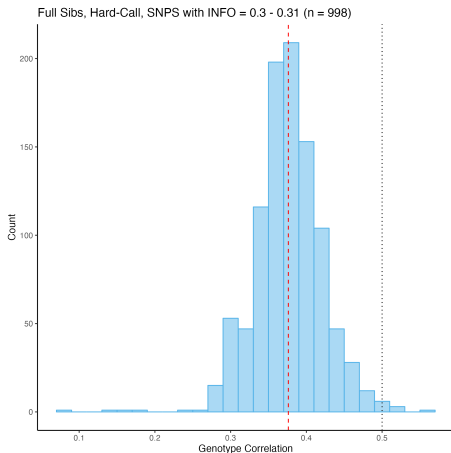
# Motivation

- Family-based research designs rely on special properties of the data.

- Low-quality imputed genotypes may not work for family-based analysis.

- Understaning the impact of imputation quality on downstream analysis by comparing to the WGS data
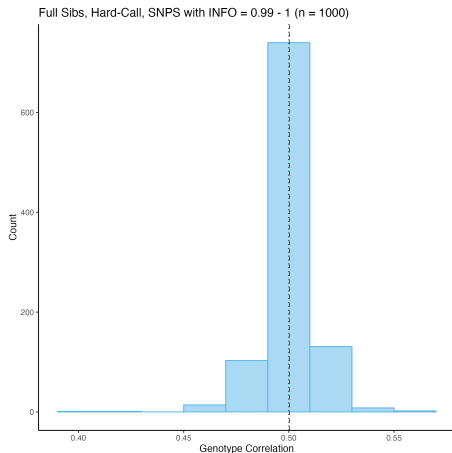
# Correlation Analysis in UK Biobank

- UK Biobank Imputed Data

  - White British subsample

  - $19K$ sibling pairs

  - $4K$ parent-offspring pairs

  - SNPs with $MAF > 1\%$

- Howe et.al (2022) Sib-GWAS used low-quality (Info Score $> 0.3$) imputed SNPs.

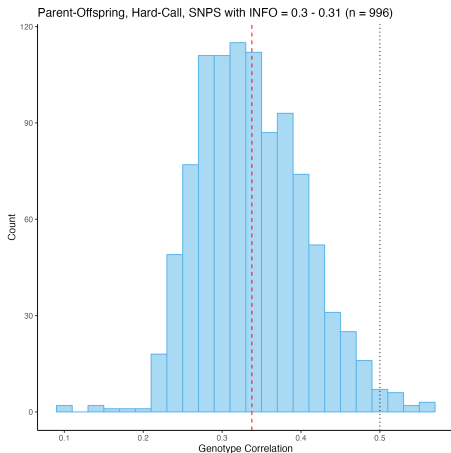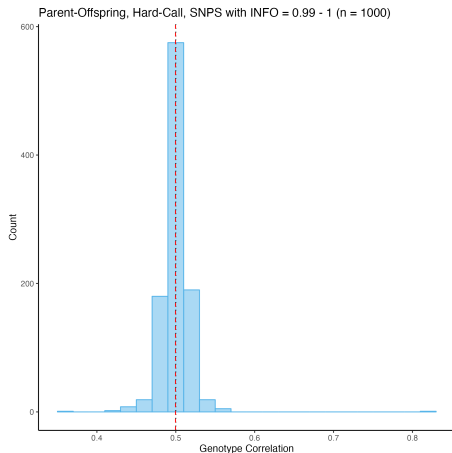# Correlations Distribution - Full Siblings

## Low Quality Imputed



Full Sibs, Hard-Call, SNPS with INFO = 0.3 - 0.31 (n = 998)

## High Quality Imputed



Full Sibs, Hard-Call, SNPS with INFO = 0.99 - 1 (n = 1000)

# Correlations Distribution - Parent-Offspring

## Low Quality Imputed



Parent-Offspring, Hard-Call, SNPS with INFO = 0.3 - 0.31 (n = 996)

## High Quality Imputed



Parent-Offspring, Hard-Call, SNPS with INFO = 0.99 - 1 (n = 1000)

# Correlation Analysis Conditional on IBD states

- Quantitative genetics theory tells us the correlation between siblings' genotypes depends on their IBD state.

- IBD state records how many allels they share by descent from their parents.

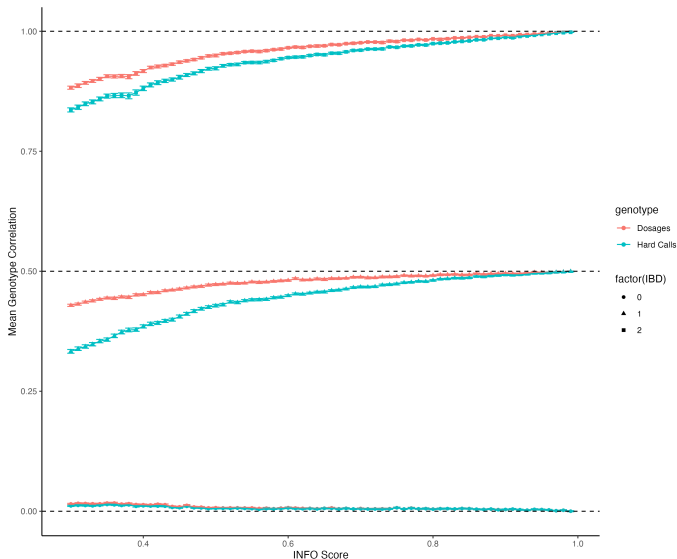- Suppose $i$ and $j$ are siblings. Then in theory [under random-mating] we have:

$$Corr(G_i, G_j | IBD = 0) = 0$$

$$Corr(G_i, G_j | IBD = 1) = 0.5$$

$$Corr(G_i, G_j | IBD = 2) = 1$$

## As a Function of Info Score

# Next Steps

- Using recently released WGS data from UKB. We are interesed to see what is the the downstream effect of using low-quality imputed genotypes in Family-Based analysis.

- We can do that by comparing the results of Family-Based analysis using imputed genotypes and the WGS data.

# Conclusion

- Genotyps imputed from a reference panel do not preserve Mandellian laws except for the very highest quality imputed variants.

- This is worse for best guess (Hard Calls) genotypes than for dosages.

- We are interested in developing reference-based imputation methods that take into account the relationships between the individuals.

# Thank You!