

Comparison between SNP array and imputed data to estimate population structure in horse breeds

<u>Chessari G.</u>, Reich P., Criscione A., Falker-Gieske C., Mastrangelo S., Tumino S., Bordonaro S., Marletta D. and Tetens J.







Origin of *Equus caballus*

The domestication began ~5000 years ago in the *Eurasian steppe*, followed by multiple events that led to the definition of different genetic roots. Horse species has played a pivotal role in human evolution, involving both *production* aspects and *recreational* and *sporting* aspects.

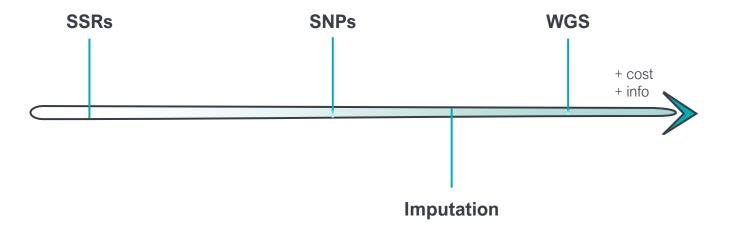
- Assessing breeding values
- Establishing management plans
- Safeguarding native breeds
- Improving of existing breeds



Origin of *Equus caballus*

The domestication began ~5000 years ago in the *Eurasian steppe*, followed by multiple events that led to the definition of different genetic roots. Horse species has played a pivotal role in human evolution, involving both *production* aspects and *recreational* and *sporting* aspects.

- Assessing breeding values
- Establishing management plans
- Safeguarding native breeds
- Improving of existing breeds

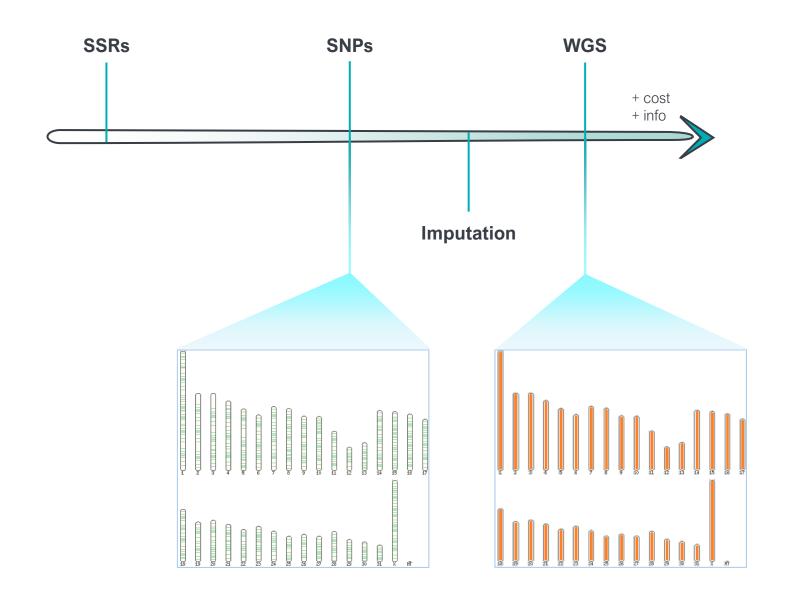




Origin of *Equus caballus*

The domestication began ~5000 years ago in the *Eurasian steppe*, followed by multiple events that led to the definition of different genetic roots. Horse species has played a pivotal role in human evolution, involving both *production* aspects and *recreational* and *sporting* aspects.

- Assessing breeding values
- Establishing management plans
- Safeguarding native breeds
- Improving of existing breeds

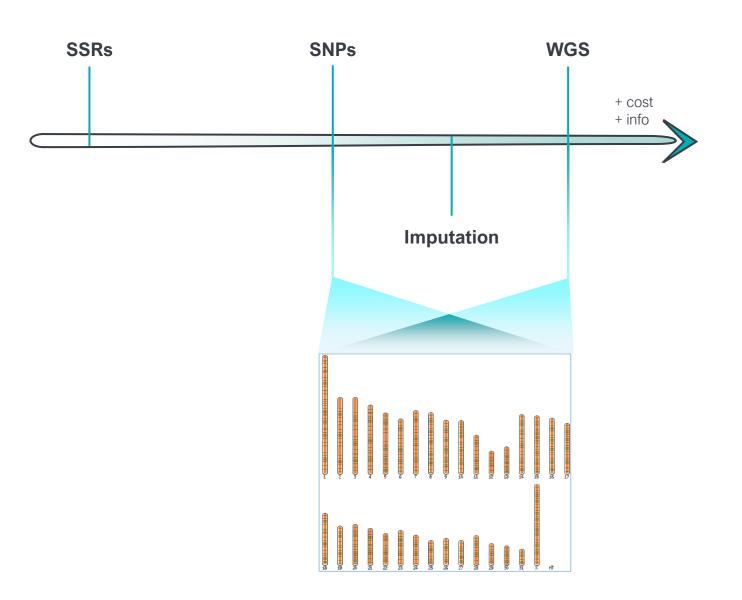




Origin of *Equus caballus*

The domestication began ~5000 years ago in the *Eurasian steppe*, followed by multiple events that led to the definition of different genetic roots. Horse species has played a pivotal role in human evolution, involving both *production* aspects and *recreational* and *sporting* aspects.

- Assessing breeding values
- Establishing management plans
- Safeguarding native breeds
- Improving of existing breeds





WGS Reference panel





SEQUENCES



Illumina paired-end x20



Purosangue Orientale Siciliano



Siciliano



Sanfratellano

MAPPING and VARIANT CALLING

EquCab3.0 genome assembly Best practices recommendations by GATK v4.1.7.0

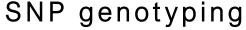
Read group
Mark Illumina adapters
Align and Merge reads
Mark duplicates
BQSR
Haplotype calling
Combine gVCFs
Genotype gVCFs

Merging to a reference panel including 317 horses belonging to 46 worldwide breeds (Reich et al., 2022)



Raw SNP panel







Medium density arrays

Sicilian horses panel

(Criscione et al., 2022)

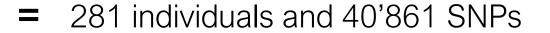
- 65'157 SNPs
- 36 individuals
- 3 Sicilian breeds

Worldwide breeds panel

(Petersen et al., 2013)

- 40'816 SNPs
- 814 individuals
- 38 breeds

	Breed	Code	Size
Endurance / Riding horse	Purosangue Orientale Sic	ORI	9
	Sanfratellano	SAN	13
	Siciliano	SIC	14
Dressage / Show Jumping	German Hanoverian	GER	15
Endurance	Akhal-Teke	AKT	19
	Arabian	ARA	24
Draft	Franches-Montagnes	FRA	19
Light-Draft	Icelandic	ICE	25
	Shetland	SHE	27
Riding horse	Quarter Horse	QUA	40
	Standardbred	STA	40
	Thoroughbred	THO	36



























Imputation





SNP PANEL





- 40'861 SNPs
- 281 individuals
- 12 breeds

IMPUTATION PROCESS

EquCab3.0 genome assembly Based on the horse ReferencePanel Software: Beagle v5.1

Final dataset
Update to Reference assembly
Imputation step per chr
Combine VCFs
Calculation of DR2



Imputation



SNP PANEL

Medium density arrays



- 40'861 SNPs
- 281 individuals
- 12 breeds

IMPUTATION PROCESS

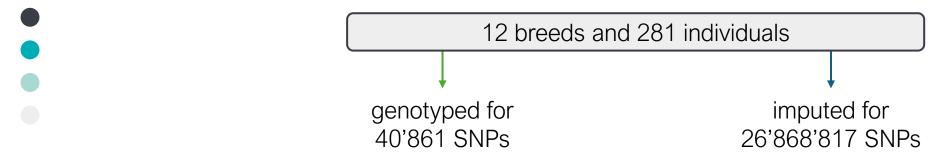
EquCab3.0 genome assembly Based on the horse ReferencePanel Software: Beagle v5.1

Final dataset
Update to Reference assembly
Imputation step per chr
Combine VCFs
Calculation of DR2



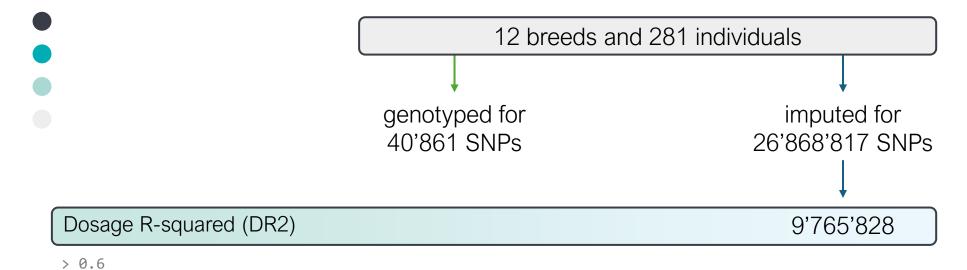
26'868'817 SNPs



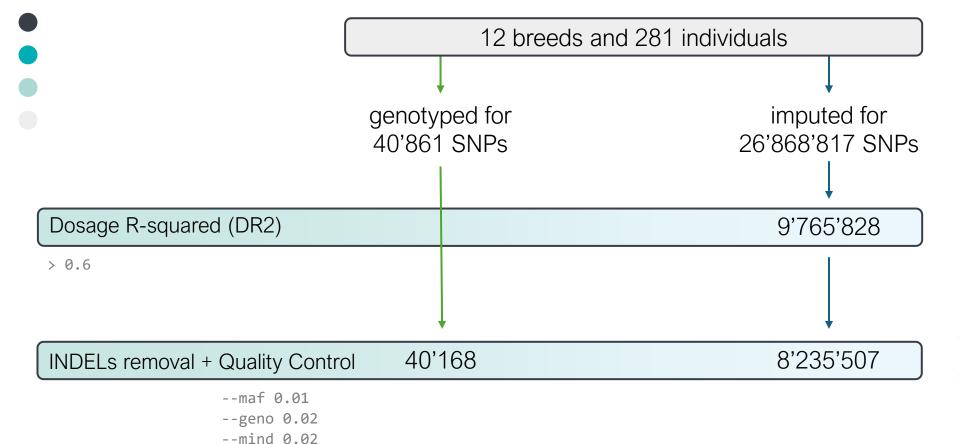


Are imputed data always more informative?

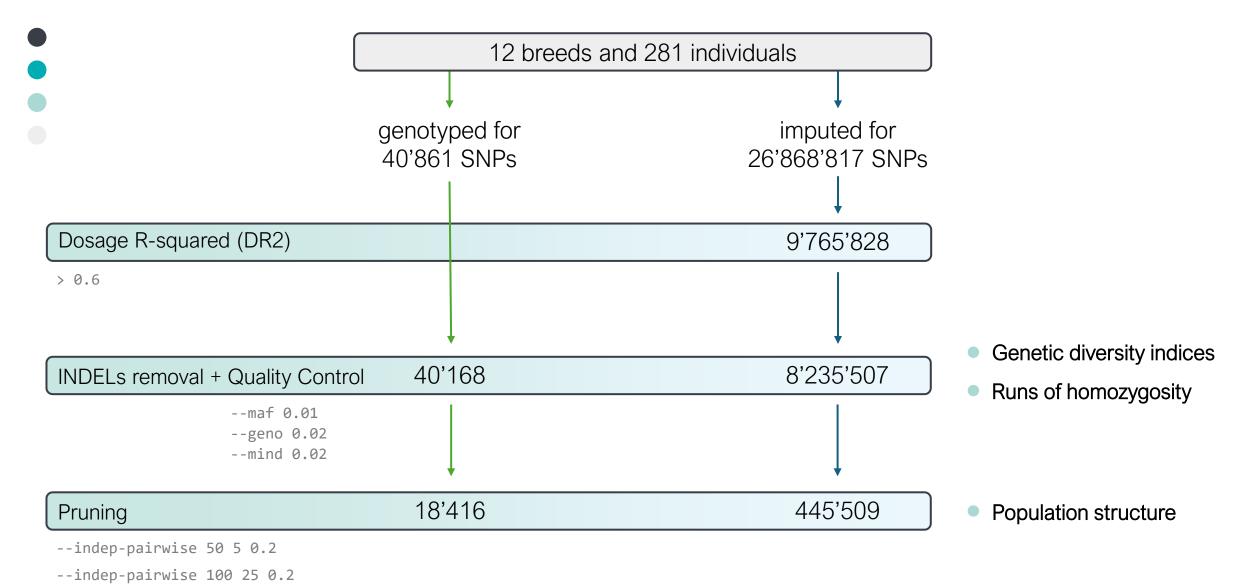








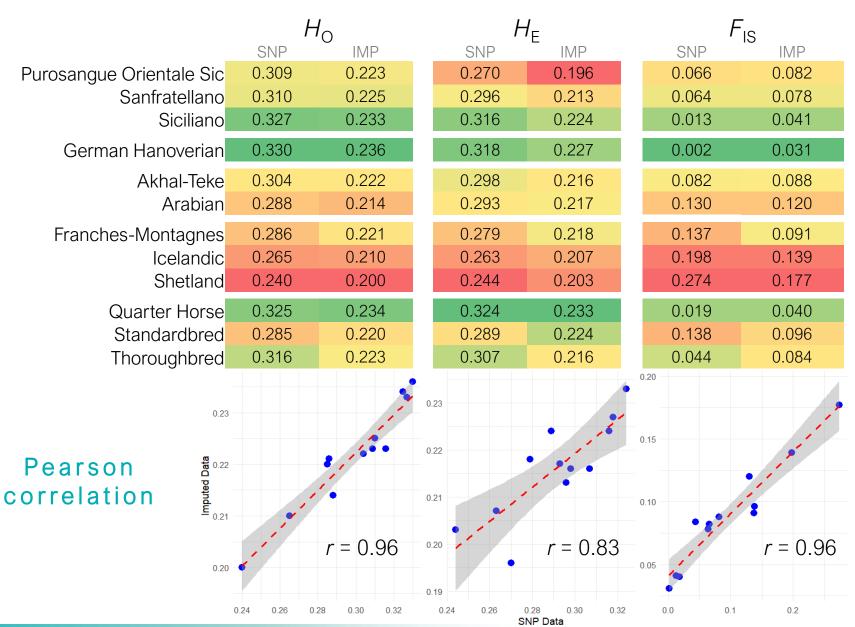
- Genetic diversity indices
- Runs of homozygosity





Results and Discussion •





H_O = Observed heterozygosity

H_E = Expected heterozygosity

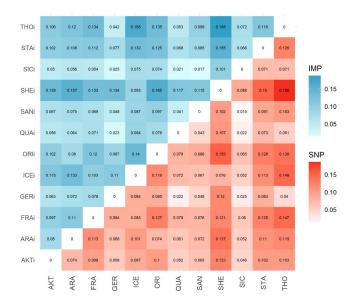
 F_{IS} = Inbreeding coefficient



Results and Discussion •



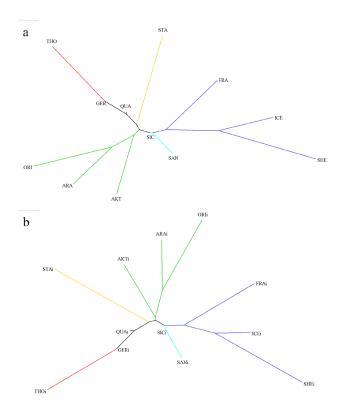
Pairwise F_{ST} distances



r: 0.962

p-value: 0.001

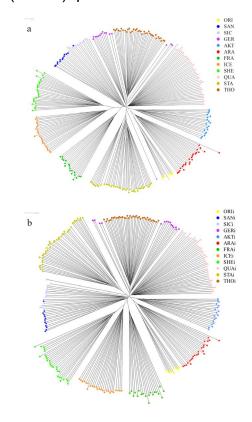
Neighbor-Joining tree on Reynolds' genetic distances



r: 0.961

p-value: 0.001

Allele Sharing Distances (ASD) pairwise distances



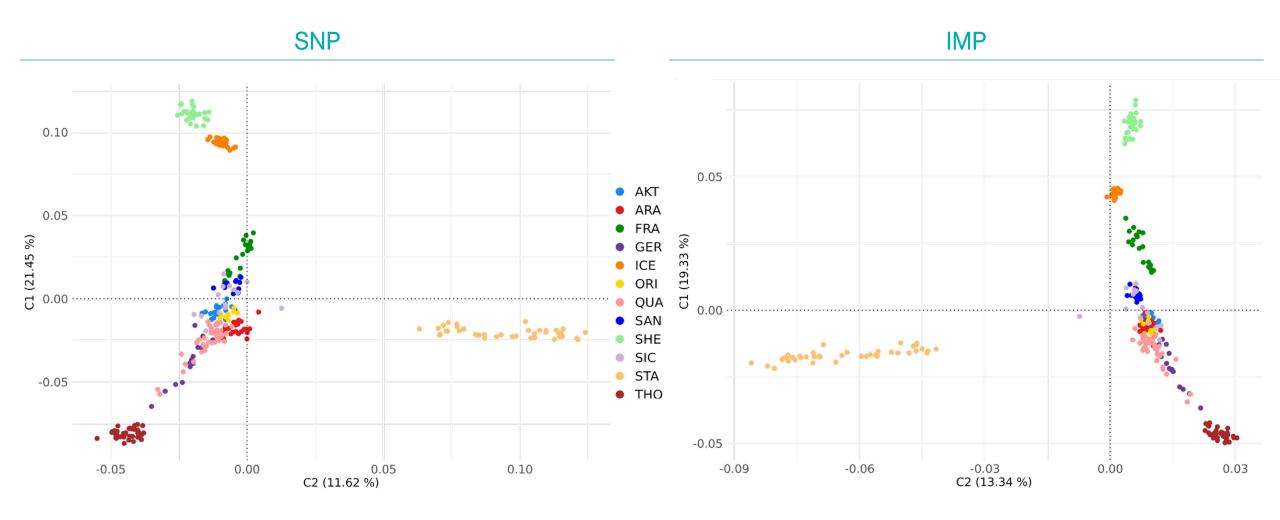
r: 0.813

p-value: 0.001



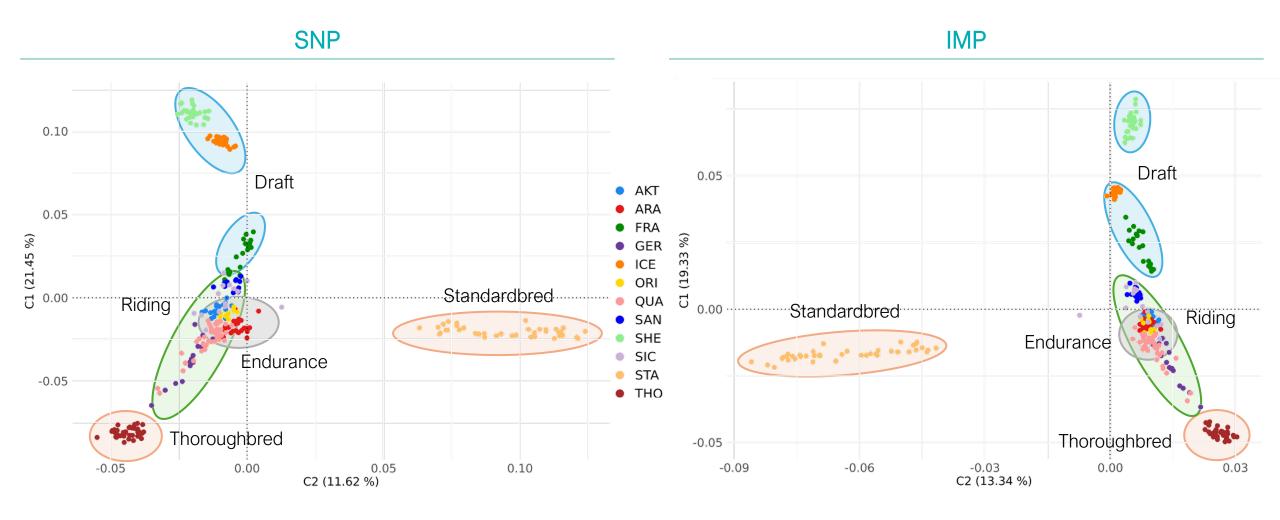


Results and Discussion • • • •



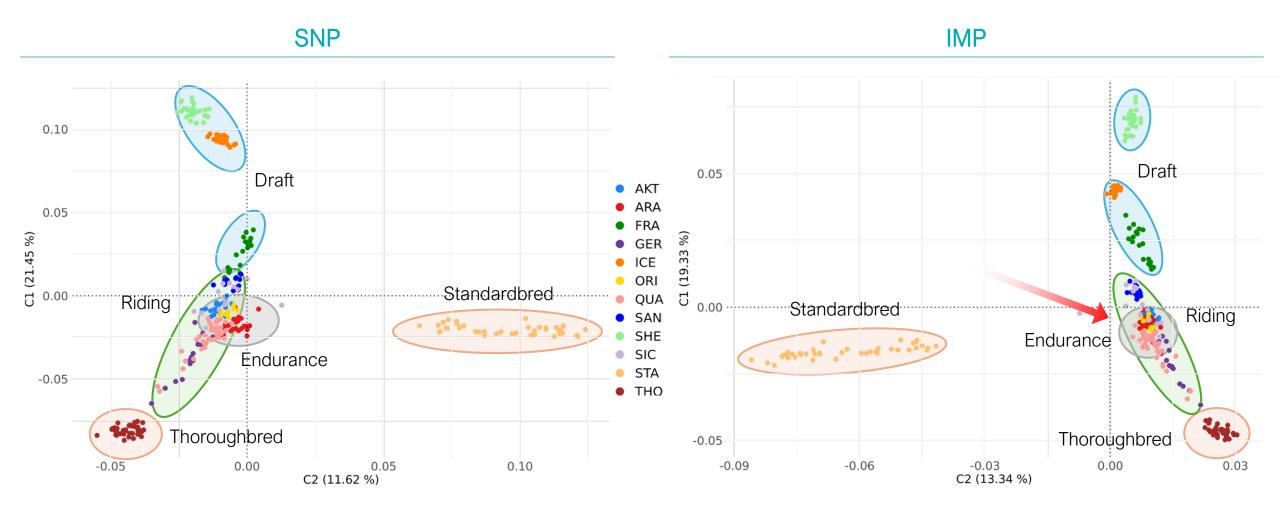


Results and Discussion • • • •





Results and Discussion • • • •





Runs of homozygosity



SNP dataset after quality control

Sliding Windows

Threshold 0.05
windowSize 25
minSNP 20
ROHet FALSE
maxOppRun 0
maxMissRun 0
maxOppWindow 0
maxMissWindow 0
maxGap 1'000'000
minLengthBps 100'000
minDensity 1/100

For a total of 23'026 ROHs

IMP dataset after quality control

Sliding Windows

```
Threshold 0.05
windowSize 500
minSNP 50
ROHet FALSE
maxOppRun 0
maxMissRun 0
maxOppWindow 0
maxMissWindow 0
maxGap 1'000'000
minLengthBps 100'000
minDensity 1/50
```

For a total of 117'040 ROHs



Runs of homozygosity



SNP dataset after quality control

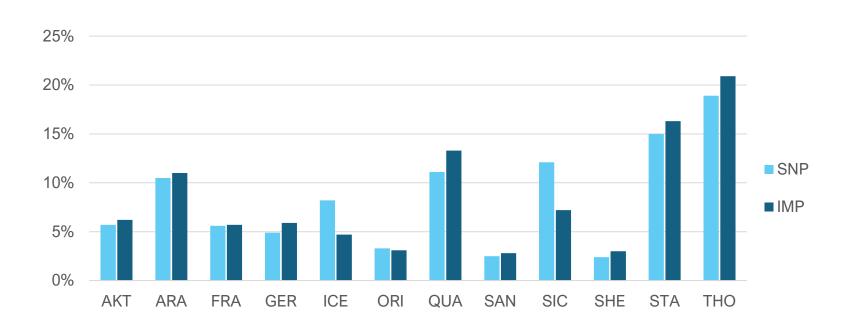
Sliding Windows

For a total of 23'026 ROHs

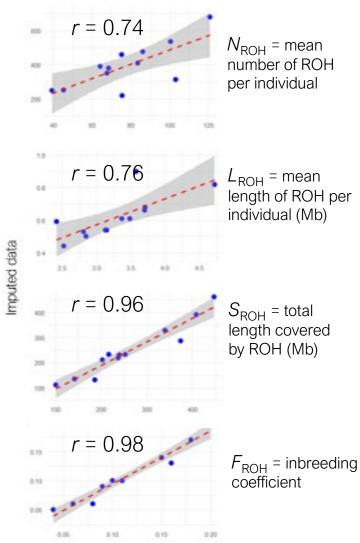
IMP dataset after quality control

Sliding Windows

For a total of 117'040 ROHs



Pearson correlation



SNP data





SNP dataset after quality control

Sliding Windows

Threshold 0.05
windowSize 25
minSNP 20
ROHet FALSE
maxOppRun 0
maxMissRun 0
maxOppWindow 0
maxMissWindow 0
maxGap 1'000'000
minLengthBps 100'000
minDensity 1/100

For a total of 23'026 ROHs

IMP dataset after quality control

Sliding Windows

Threshold 0.05
windowSize 500
minSNP 50
ROHet FALSE
maxOppRun 0
maxMissRun 0
maxOppWindow 0
maxMissWindow 0
maxGap 1'000'000
minLengthBps 100'000
minDensity 1/50

For a total of 117'040 ROHs

WGS dataset after quality control

Sliding Windows

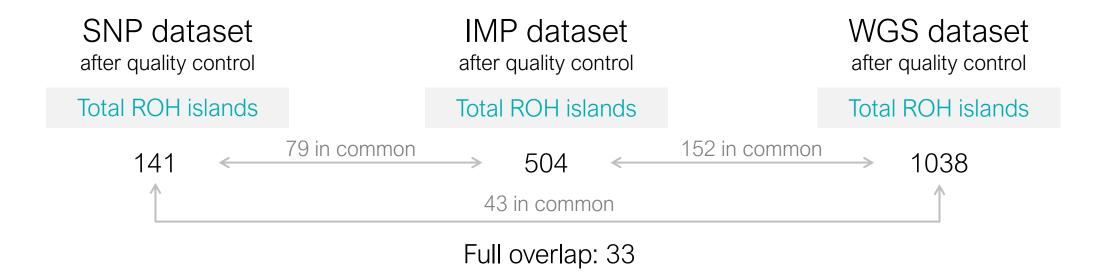
Threshold 0.05
windowSize 500
minSNP 50
ROHet FALSE
maxOppRun 0
maxMissRun 0
maxOppWindow 0
maxMissWindow 0
maxGap 1'000'000
minLengthBps 100'000
minDensity 1/50

For a total of 137'959 ROHs

Highly homozygous genomic regions (ROH islands) were identified by transforming the SNP-within-ROH incidences per population into standard normal z-scores and thus calculating the p-value; only the top 0.5% of SNPs were considered to constitute an island

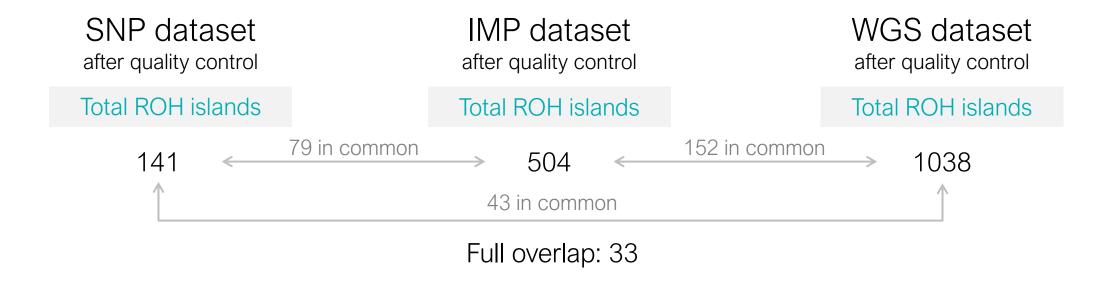








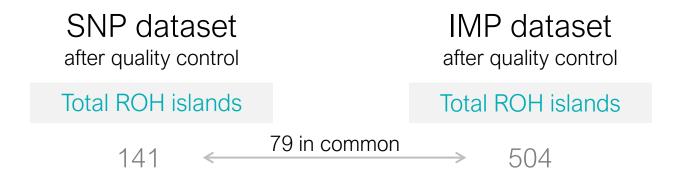


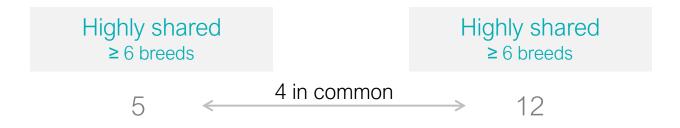






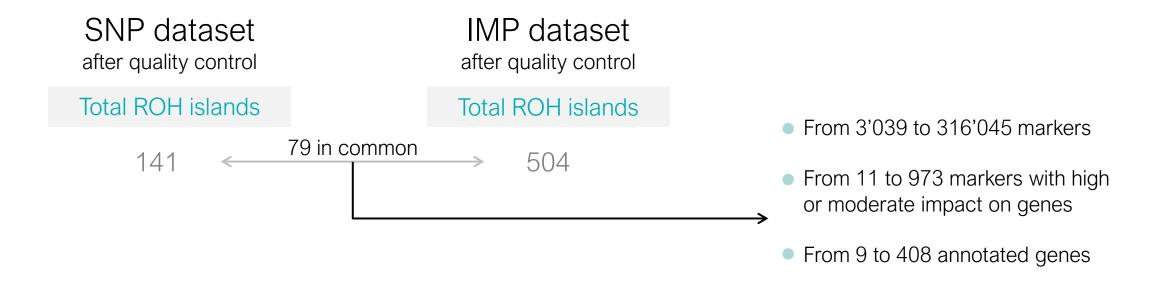
















Conclusion



Are imputed data always more informative?

Reliability of imputed dataset

2 Application for predicting selection signatures

3 Limitation of the method: accuracy



Contacts

Any questions?



giorgio.chessari@phd.unict.it



giorgio.chessari@agr.uni-goettingen.de

