

# Tutorial

## Test Report

Automatically generated by genipe

September 10<sup>th</sup>, 2015

## Contents

<b>1</b>	<b>Background</b>	<b>1</b>
<b>2</b>	<b>Methods</b>	<b>1</b>
<b>3</b>	<b>Results</b>	<b>2</b>
3.1	Cross-validation . . . . .	2
3.2	Completion rate . . . . .	10
3.3	Minor allele frequencies . . . . .	10
<b>4</b>	<b>Conclusions</b>	<b>10</b>
	<b>References</b>	<b>11</b>
	<b>Annex I: Execution time</b>	<b>12</b>

## 1 Background

The aim of this project is to perform genome-wide imputation using the study cohort.

## 2 Methods

The following (cleaned) files provided information about the study cohort dataset for 90 samples and 2,278,357 markers (including 0 markers located on sexual or mitochondrial chromosomes):

- data/hapmap\_CEU\_r23a\_hg19.bed
- data/hapmap\_CEU\_r23a\_hg19.bim
- data/hapmap\_CEU\_r23a\_hg19.fam

IMPUTE2's pre-phasing approach can work with phased haplotypes from SHAPEIT, a highly accurate phasing algorithm that can handle mixtures of unrelated samples, duos or trios. The usage of SHAPEIT is highly recommended to infer haplotypes underlying the study genotypes. The phased haplotypes are then passed to IMPUTE2 for imputation. Although pre-phasing allows for very fast imputation, it leads to a small loss in accuracy since the estimation uncertainty in the study haplotypes is ignored. SHAPEIT version v2.r790 [1] and IMPUTE2 version 2.3.2 [2, 3, 4] were used for this analysis. Binary pedfiles were processed using Plink version v1.07 [5].

To speed up the pre-phasing and imputation steps, the dataset was split by chromosome. The following quality steps were then performed on each chromosome:

1. Ambiguous markers with alleles A/T and C/G, duplicated markers (same position), and markers located on sexual or mitochondrial chromosomes were excluded from the imputation. An initial strand check was also performed using the human reference genome. **In total, 349,533 ambiguous, 0 duplicated and 0 non-autosomal markers were excluded. Also, 338 markers were flipped because of strand issue.**
2. Markers' strand was checked using the SHAPEIT algorithm and IMPUTE2's reference files. **In total, 743 markers had an incorrect strand and were flipped using Plink.**
3. The strand of each marker was checked again using SHAPEIT against IMPUTE2's reference files. **In total, 743 markers were found to still be on the wrong strand, and were hence excluded from the final dataset using Plink.**

**In total, 1,928,081 were used for phasing using SHAPEIT.** IMPUTE2 was then used to impute markers genome-wide using its reference file (filtering out sites where ALL<0.01 or ALL>0.99).

## 3 Results

### 3.1 Cross-validation

According to IMPUTE2's documentation, the cross-validation tables are "based on an internal cross-validation that is performed during each IMPUTE2 run. For this analysis, the program masks the genotypes of one variant at a time in the study data and imputes the masked genotypes by using the remaining study and reference data. The imputed genotypes are then compared with the original genotypes to produce the concordance statistics."

Tables I to XXII show the cross-validation results for the autosomes (chromosomes 1 to 22). Table XXIII shows the cross-validation results across the autosomes.

**Table I:** IMPUTE2's internal cross-validation for chromosome 1. Tables show the percentage of concordance between genotyped calls and imputed calls for 13,280,400 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[≥ 0.0]	100.0	97.9
[0.1 – 0.2]	0	0.0	[≥ 0.1]	100.0	97.9
[0.2 – 0.3]	0	0.0	[≥ 0.2]	100.0	97.9
[0.3 – 0.4]	0	0.0	[≥ 0.3]	100.0	97.9
[0.4 – 0.5]	3,209	34.8	[≥ 0.4]	100.0	97.9
[0.5 – 0.6]	21,669	50.0	[≥ 0.5]	100.0	97.9
[0.6 – 0.7]	21,343	57.7	[≥ 0.6]	99.8	98.0
[0.7 – 0.8]	25,628	65.6	[≥ 0.7]	99.6	98.0
[0.8 – 0.9]	39,797	74.3	[≥ 0.8]	99.5	98.1
[0.9 – 1.0]	13,168,754	98.2	[≥ 0.9]	99.1	98.2

**Table II:** IMPUTE2's internal cross-validation for chromosome 2. Tables show the percentage of concordance between genotyped calls and imputed calls for 15,643,890 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[≥ 0.0]	100.0	98.7
[0.1 – 0.2]	0	0.0	[≥ 0.1]	100.0	98.7
[0.2 – 0.3]	0	0.0	[≥ 0.2]	100.0	98.7
[0.3 – 0.4]	0	0.0	[≥ 0.3]	100.0	98.7
[0.4 – 0.5]	1,647	37.2	[≥ 0.4]	100.0	98.7
[0.5 – 0.6]	14,990	51.7	[≥ 0.5]	100.0	98.7
[0.6 – 0.7]	15,171	59.6	[≥ 0.6]	99.9	98.8
[0.7 – 0.8]	18,617	68.8	[≥ 0.7]	99.8	98.8
[0.8 – 0.9]	29,503	77.6	[≥ 0.8]	99.7	98.9
[0.9 – 1.0]	15,563,962	98.9	[≥ 0.9]	99.5	98.9

**Table III:** IMPUTE2's internal cross-validation for chromosome 3. Tables show the percentage of concordance between genotyped calls and imputed calls for 11,673,990 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.4
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.4
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.4
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.4
[0.4 – 0.5]	1,982	35.9	[ $\geq$ 0.4]	100.0	98.4
[0.5 – 0.6]	14,363	52.1	[ $\geq$ 0.5]	100.0	98.4
[0.6 – 0.7]	14,536	60.2	[ $\geq$ 0.6]	99.9	98.5
[0.7 – 0.8]	17,386	68.9	[ $\geq$ 0.7]	99.7	98.5
[0.8 – 0.9]	27,969	77.7	[ $\geq$ 0.8]	99.6	98.6
[0.9 – 1.0]	11,597,754	98.6	[ $\geq$ 0.9]	99.3	98.6

**Table IV:** IMPUTE2's internal cross-validation for chromosome 4. Tables show the percentage of concordance between genotyped calls and imputed calls for 10,945,350 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.2
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.2
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.2
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.2
[0.4 – 0.5]	1,796	36.5	[ $\geq$ 0.4]	100.0	98.2
[0.5 – 0.6]	14,037	51.9	[ $\geq$ 0.5]	100.0	98.2
[0.6 – 0.7]	14,334	59.0	[ $\geq$ 0.6]	99.8	98.3
[0.7 – 0.8]	17,259	66.7	[ $\geq$ 0.7]	99.7	98.3
[0.8 – 0.9]	26,882	76.3	[ $\geq$ 0.8]	99.6	98.4
[0.9 – 1.0]	10,871,042	98.4	[ $\geq$ 0.9]	99.3	98.4

**Table V:** IMPUTE2's internal cross-validation for chromosome 5. Tables show the percentage of concordance between genotyped calls and imputed calls for 10,952,820 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.6
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.6
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.6
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.6
[0.4 – 0.5]	1,420	37.0	[ $\geq$ 0.4]	100.0	98.6
[0.5 – 0.6]	11,532	51.5	[ $\geq$ 0.5]	100.0	98.6
[0.6 – 0.7]	11,343	60.0	[ $\geq$ 0.6]	99.9	98.7
[0.7 – 0.8]	13,985	68.2	[ $\geq$ 0.7]	99.8	98.7
[0.8 – 0.9]	21,961	76.6	[ $\geq$ 0.8]	99.7	98.7
[0.9 – 1.0]	10,892,579	98.8	[ $\geq$ 0.9]	99.5	98.8

**Table VI:** IMPUTE2's internal cross-validation for chromosome 6. Tables show the percentage of concordance between genotyped calls and imputed calls for 11,962,800 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.7
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.7
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.7
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.7
[0.4 – 0.5]	1,284	36.1	[ $\geq$ 0.4]	100.0	98.7
[0.5 – 0.6]	11,177	50.7	[ $\geq$ 0.5]	100.0	98.7
[0.6 – 0.7]	11,000	60.7	[ $\geq$ 0.6]	99.9	98.7
[0.7 – 0.8]	13,491	67.8	[ $\geq$ 0.7]	99.8	98.8
[0.8 – 0.9]	21,129	76.4	[ $\geq$ 0.8]	99.7	98.8
[0.9 – 1.0]	11,904,719	98.9	[ $\geq$ 0.9]	99.5	98.9

**Table VII:** IMPUTE2's internal cross-validation for chromosome 7. Tables show the percentage of concordance between genotyped calls and imputed calls for 9,180,270 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.6
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.6
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.6
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.6
[0.4 – 0.5]	1,508	34.1	[ $\geq$ 0.4]	100.0	98.6
[0.5 – 0.6]	11,874	52.4	[ $\geq$ 0.5]	100.0	98.6
[0.6 – 0.7]	11,664	60.1	[ $\geq$ 0.6]	99.9	98.6
[0.7 – 0.8]	14,063	68.1	[ $\geq$ 0.7]	99.7	98.7
[0.8 – 0.9]	21,801	77.2	[ $\geq$ 0.8]	99.6	98.7
[0.9 – 1.0]	9,119,360	98.8	[ $\geq$ 0.9]	99.3	98.8

**Table VIII:** IMPUTE2's internal cross-validation for chromosome 8. Tables show the percentage of concordance between genotyped calls and imputed calls for 10,412,010 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.9
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.9
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.9
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.9
[0.4 – 0.5]	854	36.8	[ $\geq$ 0.4]	100.0	98.9
[0.5 – 0.6]	8,676	53.0	[ $\geq$ 0.5]	100.0	98.9
[0.6 – 0.7]	8,608	60.7	[ $\geq$ 0.6]	99.9	98.9
[0.7 – 0.8]	10,518	69.9	[ $\geq$ 0.7]	99.8	98.9
[0.8 – 0.9]	16,815	78.0	[ $\geq$ 0.8]	99.7	99.0
[0.9 – 1.0]	10,366,539	99.0	[ $\geq$ 0.9]	99.6	99.0

**Table IX:** IMPUTE2's internal cross-validation for chromosome 9. Tables show the percentage of concordance between genotyped calls and imputed calls for 8,442,990 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.5
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.5
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.5
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.5
[0.4 – 0.5]	993	36.7	[ $\geq$ 0.4]	100.0	98.5
[0.5 – 0.6]	8,640	53.2	[ $\geq$ 0.5]	100.0	98.5
[0.6 – 0.7]	9,018	60.8	[ $\geq$ 0.6]	99.9	98.6
[0.7 – 0.8]	10,970	68.6	[ $\geq$ 0.7]	99.8	98.6
[0.8 – 0.9]	17,372	77.8	[ $\geq$ 0.8]	99.6	98.7
[0.9 – 1.0]	8,395,997	98.7	[ $\geq$ 0.9]	99.4	98.7

**Table X:** IMPUTE2's internal cross-validation for chromosome 10. Tables show the percentage of concordance between genotyped calls and imputed calls for 8,925,210 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.5
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.5
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.5
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.5
[0.4 – 0.5]	1,396	35.4	[ $\geq$ 0.4]	100.0	98.5
[0.5 – 0.6]	11,085	52.7	[ $\geq$ 0.5]	100.0	98.6
[0.6 – 0.7]	11,153	58.8	[ $\geq$ 0.6]	99.8	98.6
[0.7 – 0.8]	13,620	68.0	[ $\geq$ 0.7]	99.7	98.6
[0.8 – 0.9]	21,169	76.7	[ $\geq$ 0.8]	99.6	98.7
[0.9 – 1.0]	8,866,787	98.7	[ $\geq$ 0.9]	99.3	98.7

**Table XI:** IMPUTE2's internal cross-validation for chromosome 11. Tables show the percentage of concordance between genotyped calls and imputed calls for 8,593,020 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.6
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.6
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.6
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.6
[0.4 – 0.5]	1,206	35.7	[ $\geq$ 0.4]	100.0	98.6
[0.5 – 0.6]	9,953	51.1	[ $\geq$ 0.5]	100.0	98.6
[0.6 – 0.7]	9,990	60.2	[ $\geq$ 0.6]	99.9	98.7
[0.7 – 0.8]	11,853	68.9	[ $\geq$ 0.7]	99.8	98.7
[0.8 – 0.9]	18,350	77.3	[ $\geq$ 0.8]	99.6	98.7
[0.9 – 1.0]	8,541,668	98.8	[ $\geq$ 0.9]	99.4	98.8

**Table XII:** IMPUTE2's internal cross-validation for chromosome 12. Tables show the percentage of concordance between genotyped calls and imputed calls for 8,039,970 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.5
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.5
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.5
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.5
[0.4 – 0.5]	1,472	39.2	[ $\geq$ 0.4]	100.0	98.5
[0.5 – 0.6]	10,582	52.5	[ $\geq$ 0.5]	100.0	98.5
[0.6 – 0.7]	10,490	60.3	[ $\geq$ 0.6]	99.8	98.6
[0.7 – 0.8]	13,152	68.4	[ $\geq$ 0.7]	99.7	98.6
[0.8 – 0.9]	20,700	77.4	[ $\geq$ 0.8]	99.5	98.7
[0.9 – 1.0]	7,983,574	98.8	[ $\geq$ 0.9]	99.3	98.8

**Table XIII:** IMPUTE2's internal cross-validation for chromosome 13. Tables show the percentage of concordance between genotyped calls and imputed calls for 6,720,480 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.6
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.6
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.6
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.6
[0.4 – 0.5]	834	36.2	[ $\geq$ 0.4]	100.0	98.6
[0.5 – 0.6]	7,146	52.9	[ $\geq$ 0.5]	100.0	98.7
[0.6 – 0.7]	7,763	59.3	[ $\geq$ 0.6]	99.9	98.7
[0.7 – 0.8]	9,244	68.0	[ $\geq$ 0.7]	99.8	98.8
[0.8 – 0.9]	14,193	76.8	[ $\geq$ 0.8]	99.6	98.8
[0.9 – 1.0]	6,681,300	98.9	[ $\geq$ 0.9]	99.4	98.9

**Table XIV:** IMPUTE2's internal cross-validation for chromosome 14. Tables show the percentage of concordance between genotyped calls and imputed calls for 5,804,370 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.8
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.8
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.8
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.8
[0.4 – 0.5]	634	36.8	[ $\geq$ 0.4]	100.0	98.8
[0.5 – 0.6]	5,703	53.0	[ $\geq$ 0.5]	100.0	98.8
[0.6 – 0.7]	6,057	60.9	[ $\geq$ 0.6]	99.9	98.9
[0.7 – 0.8]	7,047	70.4	[ $\geq$ 0.7]	99.8	98.9
[0.8 – 0.9]	11,394	78.3	[ $\geq$ 0.8]	99.7	99.0
[0.9 – 1.0]	5,773,535	99.0	[ $\geq$ 0.9]	99.5	99.0

**Table XV:** IMPUTE2's internal cross-validation for chromosome 15. Tables show the percentage of concordance between genotyped calls and imputed calls for 4,791,060 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.6
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.6
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.6
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.6
[0.4 – 0.5]	760	42.2	[ $\geq$ 0.4]	100.0	98.6
[0.5 – 0.6]	6,577	53.3	[ $\geq$ 0.5]	100.0	98.6
[0.6 – 0.7]	6,810	59.5	[ $\geq$ 0.6]	99.8	98.7
[0.7 – 0.8]	8,083	68.6	[ $\geq$ 0.7]	99.7	98.7
[0.8 – 0.9]	13,339	78.2	[ $\geq$ 0.8]	99.5	98.8
[0.9 – 1.0]	4,755,491	98.9	[ $\geq$ 0.9]	99.3	98.9

**Table XVI:** IMPUTE2's internal cross-validation for chromosome 16. Tables show the percentage of concordance between genotyped calls and imputed calls for 4,533,930 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.1
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.1
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.1
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.1
[0.4 – 0.5]	1,199	37.5	[ $\geq$ 0.4]	100.0	98.1
[0.5 – 0.6]	9,679	50.8	[ $\geq$ 0.5]	100.0	98.1
[0.6 – 0.7]	9,499	57.8	[ $\geq$ 0.6]	99.8	98.2
[0.7 – 0.8]	11,543	66.7	[ $\geq$ 0.7]	99.5	98.3
[0.8 – 0.9]	18,260	75.4	[ $\geq$ 0.8]	99.3	98.4
[0.9 – 1.0]	4,483,750	98.5	[ $\geq$ 0.9]	98.9	98.5

**Table XVII:** IMPUTE2's internal cross-validation for chromosome 17. Tables show the percentage of concordance between genotyped calls and imputed calls for 3,821,760 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.1
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.1
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.1
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.1
[0.4 – 0.5]	1,118	37.9	[ $\geq$ 0.4]	100.0	98.1
[0.5 – 0.6]	8,438	51.7	[ $\geq$ 0.5]	100.0	98.1
[0.6 – 0.7]	8,751	59.2	[ $\geq$ 0.6]	99.8	98.2
[0.7 – 0.8]	10,220	68.1	[ $\geq$ 0.7]	99.5	98.3
[0.8 – 0.9]	15,895	75.1	[ $\geq$ 0.8]	99.2	98.4
[0.9 – 1.0]	3,777,338	98.4	[ $\geq$ 0.9]	98.8	98.4

**Table XVIII:** IMPUTE2's internal cross-validation for chromosome 18. Tables show the percentage of concordance between genotyped calls and imputed calls for 5,635,350 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.8
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.8
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.8
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.8
[0.4 – 0.5]	647	41.3	[ $\geq$ 0.4]	100.0	98.8
[0.5 – 0.6]	6,216	50.9	[ $\geq$ 0.5]	100.0	98.8
[0.6 – 0.7]	6,119	60.3	[ $\geq$ 0.6]	99.9	98.8
[0.7 – 0.8]	7,384	69.1	[ $\geq$ 0.7]	99.8	98.9
[0.8 – 0.9]	11,481	78.6	[ $\geq$ 0.8]	99.6	98.9
[0.9 – 1.0]	5,603,503	99.0	[ $\geq$ 0.9]	99.4	99.0

**Table XIX:** IMPUTE2's internal cross-validation for chromosome 19. Tables show the percentage of concordance between genotyped calls and imputed calls for 2,419,650 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	97.5
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	97.5
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	97.5
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	97.5
[0.4 – 0.5]	930	38.2	[ $\geq$ 0.4]	100.0	97.5
[0.5 – 0.6]	7,224	51.2	[ $\geq$ 0.5]	100.0	97.6
[0.6 – 0.7]	7,407	57.4	[ $\geq$ 0.6]	99.7	97.7
[0.7 – 0.8]	8,647	67.4	[ $\geq$ 0.7]	99.4	97.8
[0.8 – 0.9]	13,254	76.3	[ $\geq$ 0.8]	99.0	97.9
[0.9 – 1.0]	2,382,188	98.1	[ $\geq$ 0.9]	98.5	98.1

**Table XX:** IMPUTE2's internal cross-validation for chromosome 20. Tables show the percentage of concordance between genotyped calls and imputed calls for 4,379,490 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[ $\geq$ 0.0]	100.0	98.6
[0.1 – 0.2]	0	0.0	[ $\geq$ 0.1]	100.0	98.6
[0.2 – 0.3]	0	0.0	[ $\geq$ 0.2]	100.0	98.6
[0.3 – 0.4]	0	0.0	[ $\geq$ 0.3]	100.0	98.6
[0.4 – 0.5]	620	33.1	[ $\geq$ 0.4]	100.0	98.6
[0.5 – 0.6]	5,656	52.6	[ $\geq$ 0.5]	100.0	98.7
[0.6 – 0.7]	5,446	60.5	[ $\geq$ 0.6]	99.9	98.7
[0.7 – 0.8]	6,758	67.0	[ $\geq$ 0.7]	99.7	98.8
[0.8 – 0.9]	10,544	77.2	[ $\geq$ 0.8]	99.6	98.8
[0.9 – 1.0]	4,350,466	98.9	[ $\geq$ 0.9]	99.3	98.9



**Table XXI:** IMPUTE2's internal cross-validation for chromosome 21. Tables show the percentage of concordance between genotyped calls and imputed calls for 2,423,520 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[≥ 0.0]	100.0	98.2
[0.1 – 0.2]	0	0.0	[≥ 0.1]	100.0	98.2
[0.2 – 0.3]	0	0.0	[≥ 0.2]	100.0	98.2
[0.3 – 0.4]	0	0.0	[≥ 0.3]	100.0	98.2
[0.4 – 0.5]	513	37.6	[≥ 0.4]	100.0	98.2
[0.5 – 0.6]	3,805	50.3	[≥ 0.5]	100.0	98.3
[0.6 – 0.7]	3,807	59.1	[≥ 0.6]	99.8	98.4
[0.7 – 0.8]	4,627	67.3	[≥ 0.7]	99.7	98.4
[0.8 – 0.9]	7,241	77.0	[≥ 0.8]	99.5	98.5
[0.9 – 1.0]	2,403,527	98.5	[≥ 0.9]	99.2	98.5

**Table XXII:** IMPUTE2's internal cross-validation for chromosome 22. Tables show the percentage of concordance between genotyped calls and imputed calls for 2,343,690 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[≥ 0.0]	100.0	98.2
[0.1 – 0.2]	0	0.0	[≥ 0.1]	100.0	98.2
[0.2 – 0.3]	0	0.0	[≥ 0.2]	100.0	98.2
[0.3 – 0.4]	0	0.0	[≥ 0.3]	100.0	98.2
[0.4 – 0.5]	464	42.3	[≥ 0.4]	100.0	98.2
[0.5 – 0.6]	4,164	52.7	[≥ 0.5]	100.0	98.2
[0.6 – 0.7]	4,458	60.8	[≥ 0.6]	99.8	98.3
[0.7 – 0.8]	5,248	67.6	[≥ 0.7]	99.6	98.4
[0.8 – 0.9]	8,330	77.6	[≥ 0.8]	99.4	98.5
[0.9 – 1.0]	2,321,026	98.5	[≥ 0.9]	99.0	98.5

**Table XXIII:** IMPUTE2's internal cross-validation across the genome. Tables show the percentage of concordance between genotyped calls and imputed calls for 170,926,020 genotypes.

Interval	Nb Geno	Concordance (%)	Interval	Called (%)	Concordance (%)
[0.0 – 0.1]	0	0.0	[≥ 0.0]	100.0	98.5
[0.1 – 0.2]	0	0.0	[≥ 0.1]	100.0	98.5
[0.2 – 0.3]	0	0.0	[≥ 0.2]	100.0	98.5
[0.3 – 0.4]	0	0.0	[≥ 0.3]	100.0	98.5
[0.4 – 0.5]	26,486	36.7	[≥ 0.4]	100.0	98.5
[0.5 – 0.6]	213,186	51.8	[≥ 0.5]	100.0	98.5
[0.6 – 0.7]	214,767	59.5	[≥ 0.6]	99.9	98.6
[0.7 – 0.8]	259,343	68.0	[≥ 0.7]	99.7	98.6
[0.8 – 0.9]	407,379	76.8	[≥ 0.8]	99.6	98.7
[0.9 – 1.0]	169,804,859	98.7	[≥ 0.9]	99.3	98.7

## 3.2 Completion rate

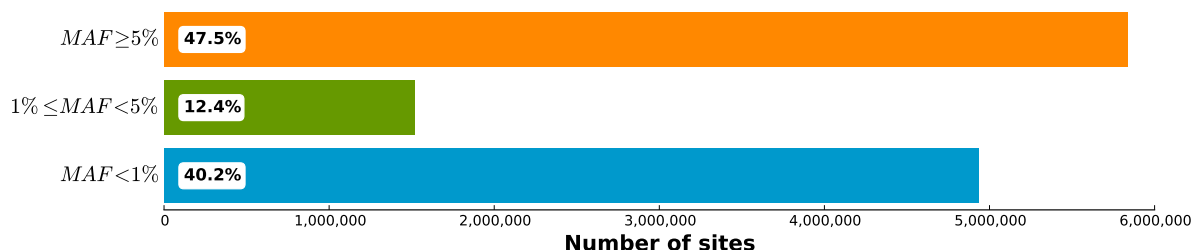
To evaluate the completion rate, we first used a probability threshold of  $\geq 90.0\%$ , which means that a genotype must have one of the three allele combination (AA, AB or BB) probabilities higher or equal to 90.0% to be considered as a *good call*.

For the 13,771,150 imputed variants, an average completion rate of 98.9% was obtained. When removing variants with an information value under 0.00, and a completion rate under 98.0%, 12,286,990 (89.2%) markers were left, with an average completion rate of 100.0%, meaning that there is a mean of 0.0 missing genotypes (for 90 samples) for each markers.

A total of 1,928,081 variants were previously genotyped, 406,033 (21.1%) of which had a call rate lower than 100% (*i.e.* 406,033 missing genotypes). A total of 406,033 (100.0%) missing genotypes were imputed with high quality (*i.e.* 1,928,081 markers now have a call rate of 100%).

## 3.3 Minor allele frequencies

Out of the 12,286,990 imputed variants with a completion rate  $\geq 98.0\%$ , there were 7,353,523 (59.8%) variants with a minor allele frequency (MAF)  $\geq 1\%$ , 5,834,575 (47.5%) variants with a MAF  $\geq 5\%$ , and 6,452,415 (52.5%) variants with a MAF  $< 5\%$ . Figure 1 shows the proportions of ultra rare (MAF  $< 1\%$ ), rare ( $1\% \leq \text{MAF} < 5\%$ ) and common (MAF  $\geq 5\%$ ) variants.



**Figure 1:** Proportions of minor allele frequencies for imputed sites with a completion rate of 98.0% or more at a probability of 90.0% or more.

## 4 Conclusions

Statistical analyses will be performed with the genome-wide imputed dataset, which include 12,286,990 imputed variants (done with an information threshold of  $\geq 0.00$ , and a completion rate of  $\geq 98.0\%$  at an imputation probability threshold of  $\geq 90.0\%$ ). This total includes 1,928,081 previously genotyped variants.

All files were generated in the **genipe** directory and were separated by chromosomes (**genipe/chr\*** directories). The final (merged) results (generated by IMPUTE2) are located in the **genipe/chr\*/final\_impute2** directories. All the output files are described below.

- **chr\*.imputed.alleles:** description of the reference and alternative allele at each site.
- **chr\*.imputed.completion\_rates:** number of missing values and completion rate for all site (using a probability threshold  $\geq 90.0\%$ ).
- **chr\*.imputed.good\_sites:** list of sites which pass the information threshold ( $\geq 0.00$ ) and the completion rate threshold ( $\geq 98.0\%$ ) using the probability threshold  $\geq 90.0\%$ .
- **chr\*.imputed.impute2:** imputation results (merged from all segments).
- **chr\*.imputed.impute2\_info:** the IMPUTE2 marker-wise information file (merged from all segments).
- **chr\*.imputed.imputed\_sites:** list of imputed sites (excluding sites that were previously genotyped in the study cohort).
- **chr\*.imputed.log:** log file of the merging procedure.

- `chr*.imputed.maf`: minor allele frequency (along with minor allele identification) for all sites using the probability threshold  $\geq 90.0\%$ .
- `chr*.imputed.map`: a map file describing the genomic location of all sites.
- `chr*.imputed.sample`: the sample file generated by the phasing step.

## References

- [1] Delaneau O, Zagury JF, Marchini J: **Improved whole-chromosome phasing for disease and population genetic studies**. *Nature methods* 2013, **10**:5–6. [DOI:10.1038/nmeth.2307].
- [2] Howie BN, Donnelly P, Marchini J: **A flexible and accurate genotype imputation method for the next generation of genome-wide association studies**. *PLoS genetics* 2009, **5**(6):e1000529. [DOI:10.1371/journal.pgen.1000529].
- [3] Howie B, Marchini J, Stephens M: **Genotype imputation with thousands of genomes**. *G3: Genes, Genomes, Genetics* 2011, **1**(6):457–470. [DOI:10.1534/g3.111.001198].
- [4] Howie B, Fuchsberger C, Stephens M, Marchini J, Abecasis GR: **Fast and accurate genotype imputation in genome-wide association studies through pre-phasing**. *Nature Genetics* 2012, **44**(8):955–959. [DOI:10.1038/ng.2354].
- [5] Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J, Sklar P, De Bakker PI, Daly MJ, et al.: **PLINK: a tool set for whole-genome association and population-based linkage analyses**. *The American Journal of Human Genetics* 2007, **81**(3):559–575. [DOI:10.1086/519795].

## Annex I: Execution times

The following tables show the execution time required by all the different tasks. All tasks are split by chromosomes. Execution times for imputation for each chromosome are means of individual segment times. Computing all genotyped markers' missing rate took 24 seconds.

**Table XXIV:** Execution time for the 'plink\_exclude\_chr\*' tasks.

Chrom	Time	Chrom	Time
1	00:00:13	12	00:00:12
2	00:00:13	13	00:00:12
3	00:00:14	14	00:00:12
4	00:00:13	15	00:00:11
5	00:00:13	16	00:00:11
6	00:00:12	17	00:00:09
7	00:00:13	18	00:00:09
8	00:00:13	19	00:00:10
9	00:00:12	20	00:00:09
10	00:00:11	21	00:00:09
11	00:00:12	22	00:00:08

**Table XXV:** Execution time for the 'shapeit\_check\_chr\*\_1' tasks.

Chrom	Time	Chrom	Time
1	00:00:32	12	00:00:19
2	00:00:33	13	00:00:14
3	00:00:28	14	00:00:11
4	00:00:30	15	00:00:09
5	00:00:26	16	00:00:12
6	00:00:24	17	00:00:08
7	00:00:25	18	00:00:09
8	00:00:27	19	00:00:07
9	00:00:17	20	00:00:06
10	00:00:18	21	00:00:05
11	00:00:17	22	00:00:04

**Table XXVI:** Execution time for the 'plink\_flip\_chr\*' tasks.

Chrom	Time	Chrom	Time
1	00:00:02	12	00:00:01
2	00:00:02	13	00:00:01
3	00:00:02	14	00:00:01
4	00:00:01	15	00:00:01
5	00:00:02	16	00:00:01
6	00:00:02	17	00:00:00
7	00:00:01	18	00:00:01
8	00:00:01	19	00:00:00
9	00:00:01	20	00:00:00
10	00:00:01	21	00:00:00
11	00:00:01	22	00:00:00

**Table XXVII:** Execution time for the 'shapeit\_check\_chr\*\_2' tasks.

Chrom	Time	Chrom	Time
1	00:00:28	12	00:00:15
2	00:00:36	13	00:00:11
3	00:00:27	14	00:00:10
4	00:00:28	15	00:00:09
5	00:00:24	16	00:00:09
6	00:00:20	17	00:00:08
7	00:00:24	18	00:00:08
8	00:00:27	19	00:00:07
9	00:00:14	20	00:00:05
10	00:00:17	21	00:00:03
11	00:00:14	22	00:00:03

**Table XXVIII:** Execution time for the 'plink\_final\_exclude\_chr\*' tasks.

Chrom	Time	Chrom	Time
1	00:00:02	12	00:00:01
2	00:00:02	13	00:00:01
3	00:00:02	14	00:00:01
4	00:00:01	15	00:00:01
5	00:00:02	16	00:00:01
6	00:00:01	17	00:00:01
7	00:00:01	18	00:00:01
8	00:00:02	19	00:00:00
9	00:00:01	20	00:00:01
10	00:00:01	21	00:00:00
11	00:00:01	22	00:00:00

**Table XXIX:** Execution time for the 'shapeit\_phase\_chr\*' tasks.

Chrom	Time	Chrom	Time
1	01:54:18	12	01:03:09
2	02:04:32	13	00:50:57
3	01:33:17	14	00:45:18
4	01:30:26	15	00:38:14
5	01:27:43	16	00:38:52
6	01:31:26	17	00:33:02
7	01:14:22	18	00:36:21
8	01:18:56	19	00:21:00
9	01:07:18	20	00:29:59
10	01:09:38	21	00:17:01
11	01:06:16	22	00:15:41

**Table XXX:** Execution time for the 'impute2\_chr\*' tasks.

Chrom	Nb Seg.	Mean T.	Max T.	Chrom	Nb Seg.	Mean T.	Max T.
1	50	00:02:29	00:04:17	12	27	00:01:58	00:02:58
2	49	00:02:43	00:04:25	13	24	00:01:29	00:02:47
3	40	00:02:31	00:03:51	14	22	00:01:31	00:02:36
4	39	00:02:24	00:03:37	15	21	00:01:18	00:02:22
5	37	00:02:14	00:03:31	16	19	00:01:36	00:02:38
6	35	00:02:22	00:03:46	17	17	00:01:29	00:02:12
7	32	00:02:09	00:03:25	18	16	00:01:40	00:02:13
8	30	00:02:17	00:03:19	19	12	00:01:24	00:01:48
9	29	00:01:44	00:02:50	20	13	00:01:30	00:02:17
10	28	00:02:05	00:03:33	21	10	00:01:04	00:01:49
11	28	00:02:01	00:03:06	22	11	00:00:58	00:01:57

**Table XXXI:** Execution time for the 'merge\_impute2\_chr\*' tasks.

Chrom	Time	Chrom	Time
1	00:05:03	12	00:03:20
2	00:05:23	13	00:02:17
3	00:04:35	14	00:02:19
4	00:04:48	15	00:02:07
5	00:03:59	16	00:02:12
6	00:03:57	17	00:01:50
7	00:04:00	18	00:01:45
8	00:03:27	19	00:01:32
9	00:03:02	20	00:01:26
10	00:03:20	21	00:00:53
11	00:03:15	22	00:00:49

**Table XXXII:** Execution time for the 'bgzip\_chr\*' tasks.

Chrom	Time	Chrom	Time
1	00:00:48	12	00:00:09
2	00:00:43	13	00:00:11
3	00:00:47	14	00:00:05
4	00:00:41	15	00:00:06
5	00:00:35	16	00:00:06
6	00:00:42	17	00:00:05
7	00:00:36	18	00:00:05
8	00:00:29	19	00:00:05
9	00:00:25	20	00:00:03
10	00:00:10	21	00:00:03
11	00:00:21	22	00:00:03