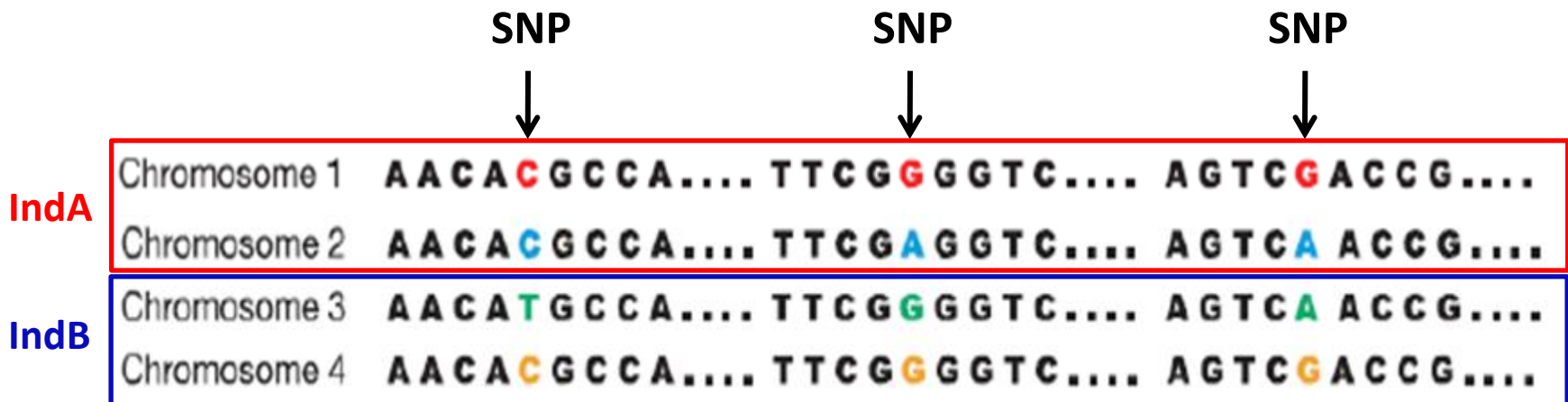


# Grupo de Imputação e Determinação Haplotípica

Wagner C S Magalhães  
Eduardo Martín Tarazona Santos

# SNPs

- SNPs are the most common form of germline variation.
- Occur when a single nucleotide differs between paired chromosomes in an individual, or between individuals in a population.
- **~25 million SNPs in the human genome** with a minor allele frequency of  $>1\%$ .
- Only a (very small) subset have pathogenic potential.



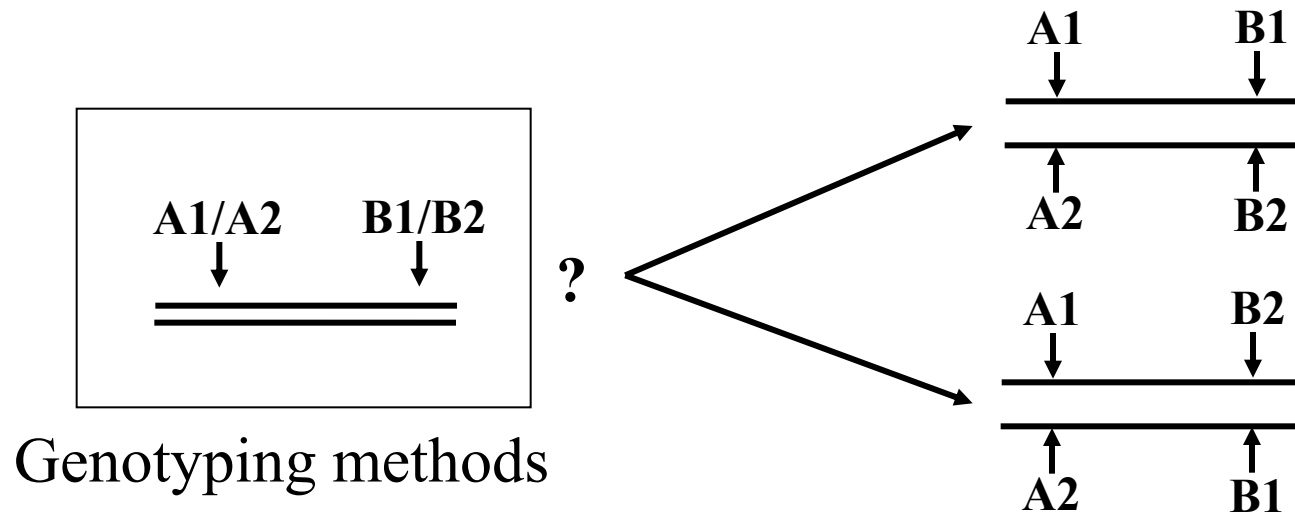
# Haplotype inference

- Problem:
  - Most of experimental methods of sequencing or genotyping, only provide genotype information at each locus: the unordered pair of alleles at the locus.
- A haplotype is a sequence of alleles that are on the same physical chromosome (i.e. that are inherited from the same parent).

# Haplotype inference

- Understanding the interplay of genetic variation and disease,
- Imputation of untyped genetic variation,
- Detecting genotype error,
- Inferring human demographic history,
- Inferring points of recombination,
- Detecting recurrent mutation,
- Signatures of natural selection

The phase problem with double (multiple) heterozygous sites:  
**which alleles are on the same chromosome?**



## Solutions

**-Labor intensive molecular / cellular techniques**

-Long range PCR, Cell fusion

**-Statistical methods (work well, example: Phase algorithm or Expectation-Maximization methods embedded in gen-epi software)**

**-Segregation analysis by typing the parents**

# Haplotype inference

Son      ATCA[G/C][C/G][A/C][T/A]CA[G/C]A.....

True haplotypes      .....ATCAGCATCAGA.....  
                         .....ATCACGCACACA.....

# Number of Haplotypes

- Possible haplotypes = 32

$m$  = loci unphased

$2^m$  haplotypes

Now imagine ----->  $m = 2.5 \text{ M}$  or even  $5 \text{ M}$

 R Console

```
> number_of_Haplotypes = 2^5000000  
> number_of_Haplotypes  
[1] Inf
```

# Haplotype Phasing

- Unrelated individuals can be phased by considering sets of common haplotypes that can explain the observed genotype data. The number of unrelated individuals present in a sample is a crucial factor in determining how well the phase can be estimated: the more individuals, the better the estimation.
- Related individuals, by contrast, can be phased by considering haplotypes that are shared identical-by-descent between individuals within families. This within-family information on identity-by-descent (IBD)



One way to improve the power of GWAS is to infer haplotype phase and use a haplotype-based method for association testing, in addition to applying single-marker association testing methods  
(Browning and Browning 2007).

Group	Feature	BEAGLE	IMPUTE	MACH
Accessibility	Operating system	Java (platform independent)	Linux, Windows, MacOS X, Solaris	Linux, MacOS X
	Licence	Free	Free for academic use	Not clear
	Source code	Not available	Not available	Availability announced
	Documentation	Commendable	Clearly structured	Incomplete
	Authors' response	Quick and detailed	Quick	Quick
	GUI	No	No	No
Input	Genotype format	Discrete; custom format	Probabilities; custom format	QTD T (Linkage)
	Reference format	Custom format	Custom format; prepared HapMap reference available	HapMap format (custom format)
	Conversion utilities	Yes	Yes	No
Processing	Target of imputation	Chromosomes	Chromosomes or segments	Chromosomes
	Memory-saving mode	Yes	No	Yes
	Known checking errors	None	Missing probability and input check	Problematic handling of missing reference
	Runtime [chr. 6]	350 minutes	433 minutes	2781 minutes
	Maximum memory allocation	2 GB	14 GB [< 1 GB with ~ 10 MB segments]	7 GB
	Memory-saving mode	Yes	No	Yes
	Strand orientation	Check	Check + autoflip	Check + autoflip

# Data Imputation

- **Imputation** (Statistics)
  - Is the substitution of some value for missing data. Once all missing values have been imputed, the dataset can then be analyzed using standard techniques for complete data.
- **Imputation** (Genetics)
  - Is the estimation of missing genotype values by using the genotypes at nearby SNPs and the haplotype frequencies seen in other individuals.

**Genotype imputation is now an essential tool in the analysis of genome-wide association scans. This technique allows geneticists to accurately evaluate the evidence for association at genetic markers that are not directly genotyped.**

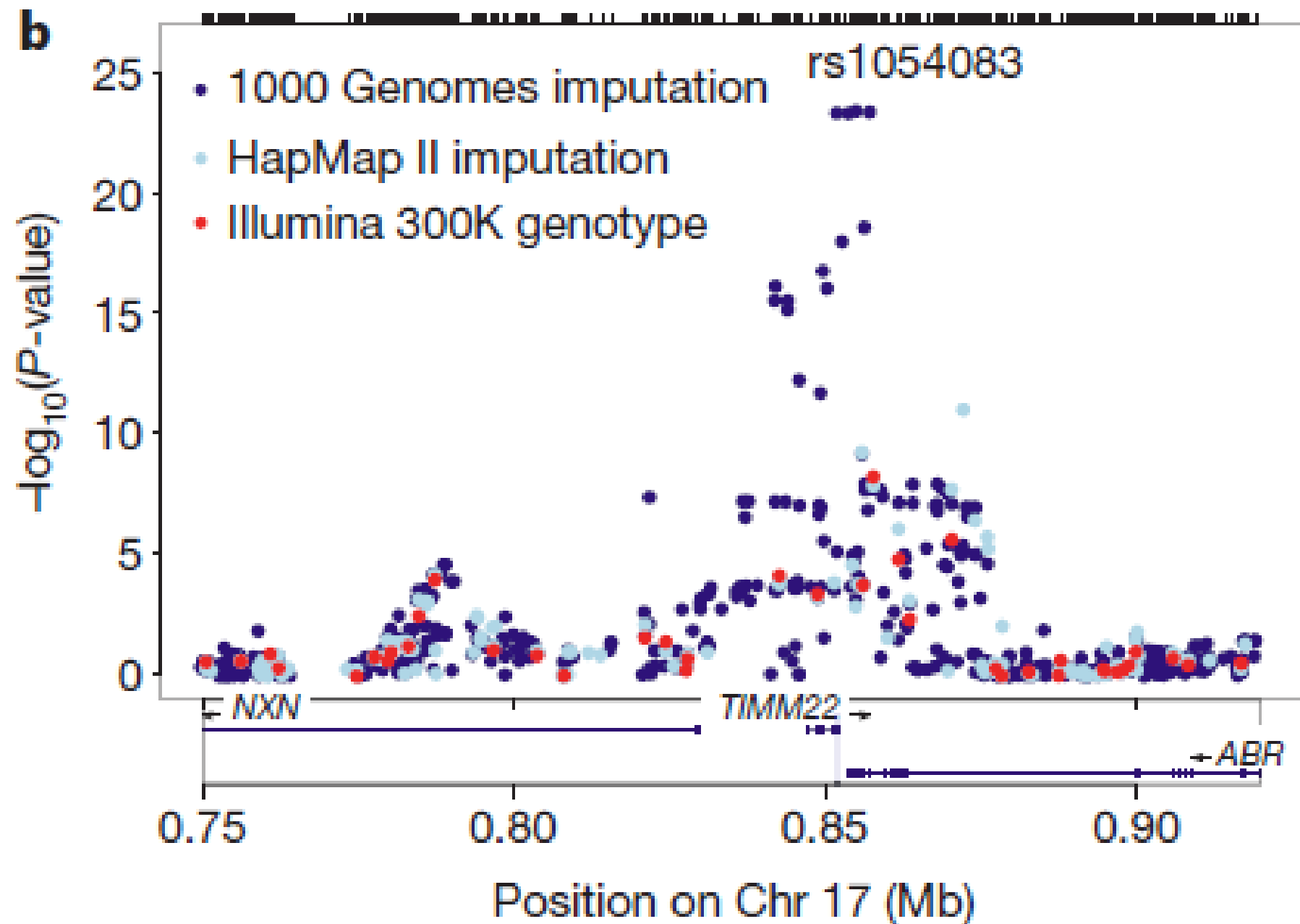
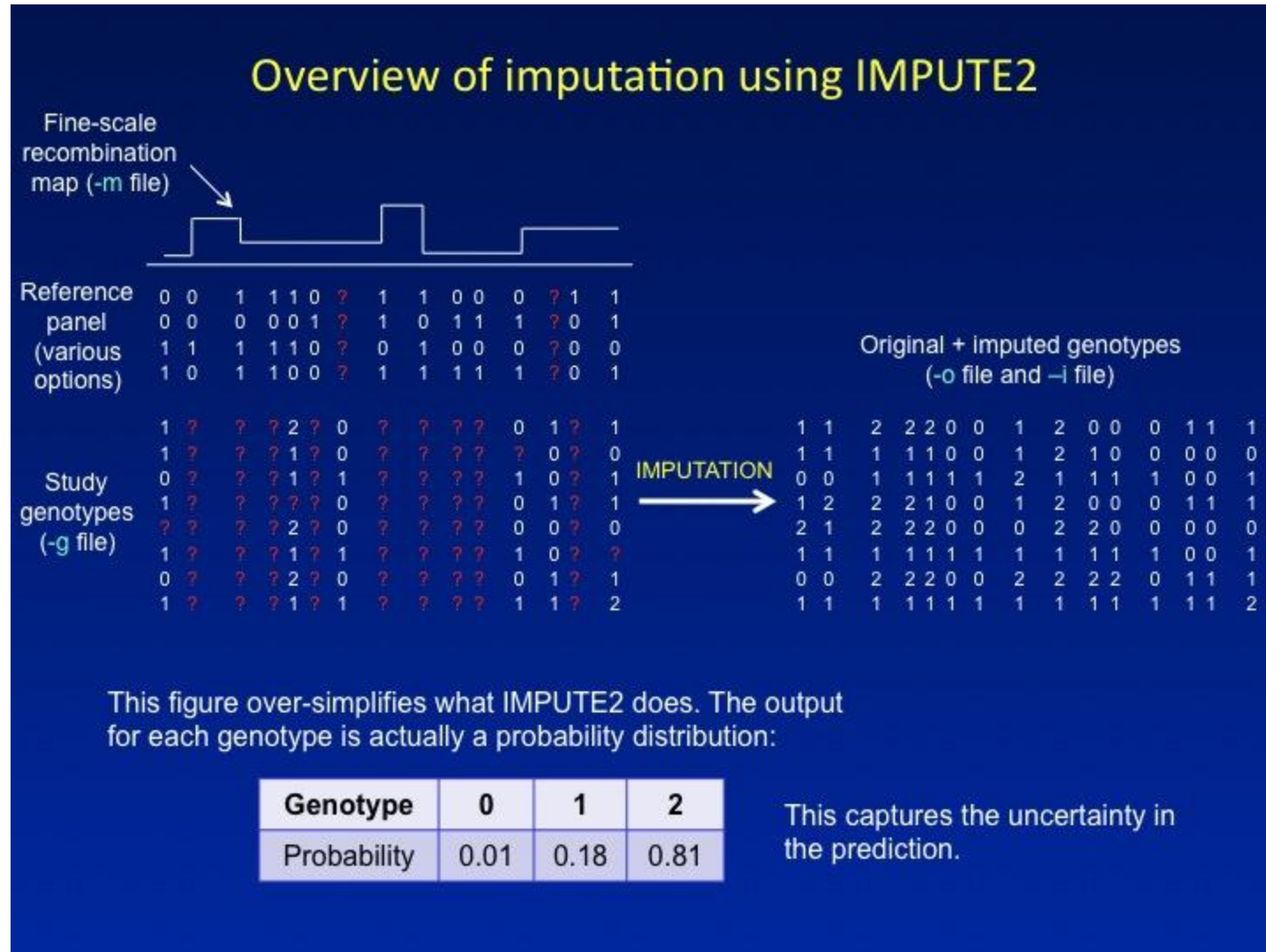
**b**

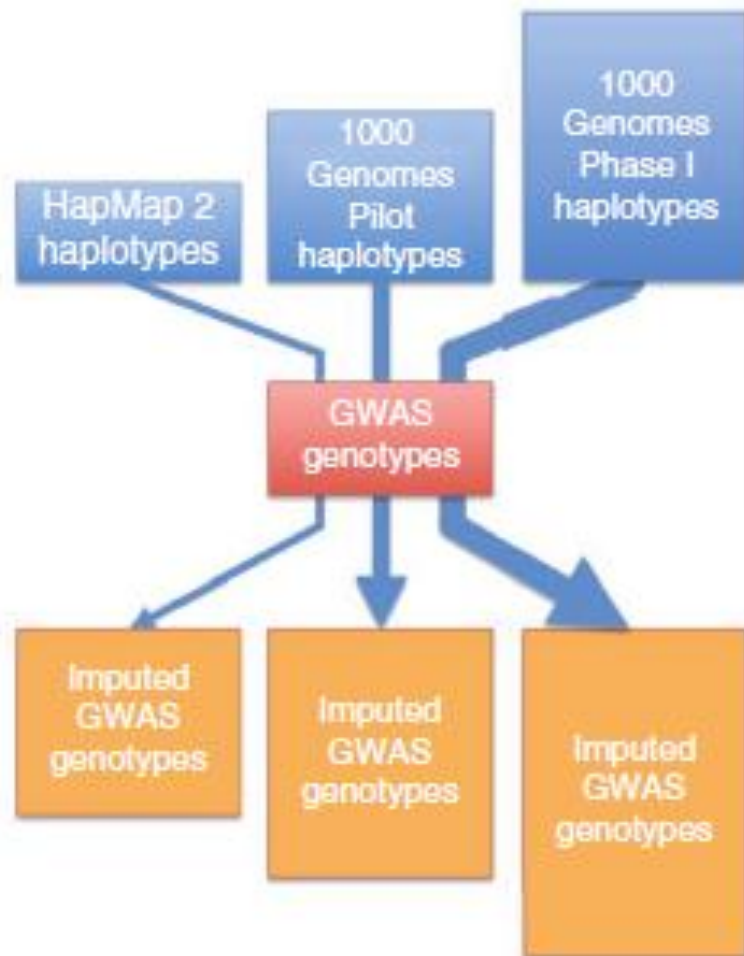
Figure 1. Schematic drawing of imputation Scenario A.



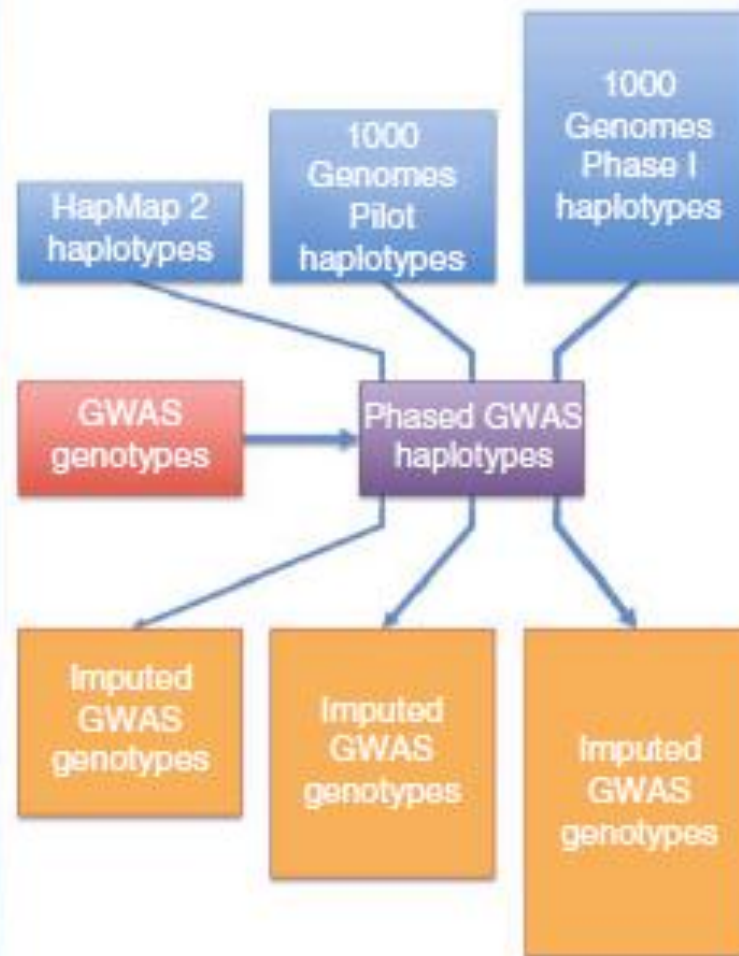
# PRE-IMPUTATION FILTERING OF STUDY GENOTYPES

- Remove low-quality variants and individuals
- Take care of their positions on the chromosome
- Strand alignment between study and reference data
- Choose reference panels that match the ancestry of their study samples

## Traditional imputation



## Pre-phasing imputation

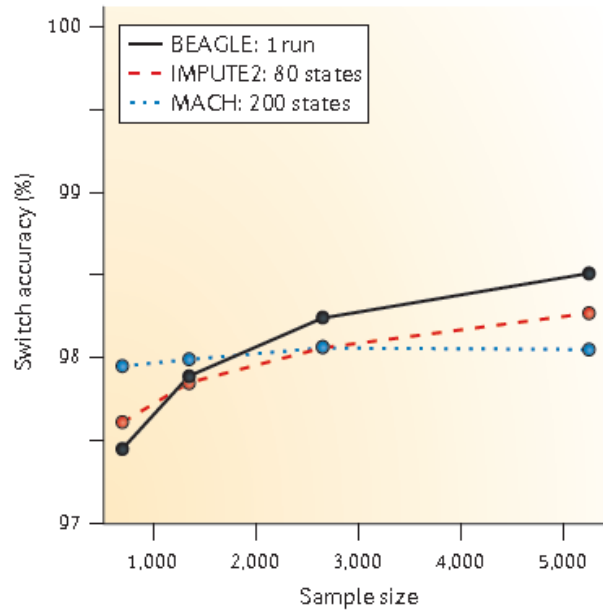
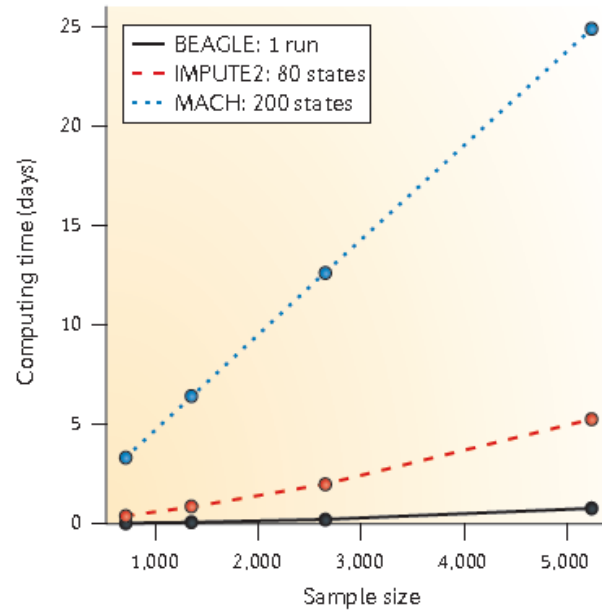
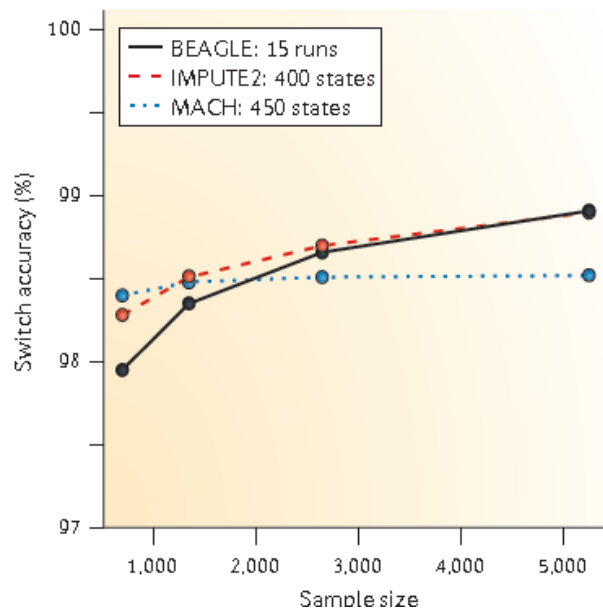
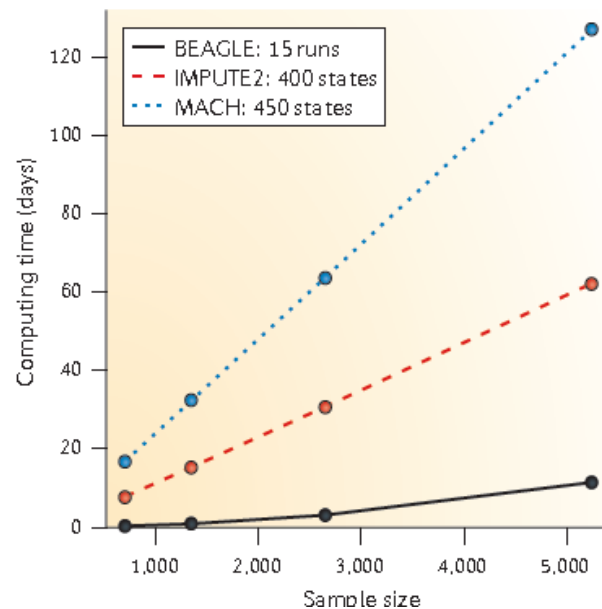


Link to download page	NCBI build	Haplotype release date	Release status
<a href="#">1000 Genomes Phase I integrated variant set</a>	b37	Mar 2012	Includes chrX; updated 19 Apr 2012
<a href="#">1000 Genomes Phase I (interim)</a>	b37	Jun 2011	Includes chrX; updated 19 Apr 2012
<a href="#">1000 Genomes (2010 interim)</a>	b37	Dec 2010	
<a href="#">1000 Genomes Pilot + HapMap 3</a>	b36	Jun 2010 / Feb 2009	
<a href="#">1000 Genomes Pilot</a>	b36	Jun 2010	
<a href="#">HapMap 3 (release #2)</a>	b36	Feb 2009	Includes chrX
<a href="#">HapMap 2 (release #24)</a>	b36	Oct 2008	
<a href="#">HapMap 2 (release #22)</a>	b36	Jan 2008	
<a href="#">HapMap 2 (release #21)</a>	b35	Jul 2006	



**Table 2 Recommended choices of HapMap reference panel haplotypes for imputing genotypes in Human Genome Diversity Panel different samples**

These reference panel haplotypes...	...are best for imputing genotypes in these Human Genome Diversity Panel samples
CEU	<b>Europe:</b> Orcadian, Basque, French, Italian, Sardinian
	<b>Middle East:</b> Druze
CHB + JPT	<b>East Asia:</b> Han, Han-Nchina, Dai, Lahu, Miao, Oroqen, She, Tujia, Tu, Xibo, Yi, Mongola, <sup>a</sup> Naxi, Japanese
YRI	<b>Africa:</b> Bantu, Yoruba, San, Mandenka, MbutiPygmy, BiakaPygmy
Combined (CEU, CHB, JPT, YRI)	<b>Europe:</b> Adygei, Russian, Tuscan
	<b>Middle East:</b> Mozabite, Bedouin, Palestinian
	<b>Asian:</b> Balochi, Brahui, Makrani, Sindhi, Pathan, Burusho, Hazara, Uygur, Kalash
	<b>East Asia:</b> Daur, Hezhen, Mongola, <sup>*</sup> Cambodian, Yakut
	<b>Oceania:</b> Melanesian, Papuan
	<b>Americas:</b> Colombian, Karitiana, Surui, Maya, Pima

**a Accuracy: standard settings****b Computing times: standard settings****c Accuracy: high-accuracy settings****d Computing times: high-accuracy settings**

**Table 1. Running times and memory requirements for various algorithms in Scenario A.**

<b>Method</b>	<b>Avg. running time (min)</b>	<b>Avg. required RAM (MB)</b>
BEAGLE	56	3100
BEAGLE (50iter)	392	3200
fastPHASE (K = 20)	397	8
fastPHASE (K = 30)	855	16
IMPUTE v1	43	1000
IMPUTE v2 (k = 40)	270	155
IMPUTE v2 (k = 80)	505	180
MACH	105	80

Running times are in minutes (min) and RAM requirements are in megabytes (MB). Each entry in the table is an average across four runs on different 7.5 Mb regions of chromosome 10. Each analysis included a reference panel of 120 chromosomes (CEU HapMap) and a study sample of 1,377 individuals genotyped on the Affymetrix 500 K SNP chip.

doi:10.1371/journal.pgen.1000529.t001

Howie BN, Donnelly P, Marchini J (2009) A Flexible and Accurate Genotype Imputation Method for the Next Generation of Genome-Wide Association Studies. *PLoS Genet* 5(6): e1000529. doi:10.1371/journal.pgen.1000529

<http://www.plosgenetics.org/article/info:doi/10.1371/journal.pgen.1000529>

## The Effect of Reference Panels and Software Tools on Genotype Imputation

[Kwangsik Nho](#), PhD,<sup>1,2</sup> [Li Shen](#), PhD,<sup>2,3</sup> [Sungeun Kim](#), PhD,<sup>2,3</sup> [Shanker Swaminathan](#), BTech,<sup>2,4</sup> [Shannon L. Risacher](#), BS,<sup>2</sup> [Andrew J. Saykin](#), PsyD,<sup>2,3,4</sup> and the Alzheimer's Disease Neuroimaging Initiative (ADNI)

[Author information ►](#) [Copyright and License information ►](#)



# Genotype imputation for Latinos using the HapMap and 1000 Genomes Project reference panels

**Xiaoyi Gao<sup>1,2\*</sup>, Talin Haritunians<sup>3</sup>, Paul Marjoram<sup>2</sup>, Roberta Mckean-Cowdin<sup>2</sup>, Mina Torres<sup>1</sup>, Kent D. Taylor<sup>3</sup>, Jerome I. Rotter<sup>3</sup>, William J. Gauderman<sup>2</sup> and Rohit Varma<sup>1,2</sup>**

<sup>1</sup> Department of Ophthalmology, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA

<sup>2</sup> Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA

<sup>3</sup> Medical Genetics Institute, Cedars-Sinai Medical Center, Los Angeles, CA, USA

## Edited by:

Rongling Wu, Pennsylvania State University, USA

## Reviewed by:

Ashok Ragavendran, Purdue University, USA

Wei Hou, University of Florida, USA

## \*Correspondence:

Xiaoyi Gao, Departments of Ophthalmology and Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA 90033, USA.  
e-mail: xiaoyiga@usc.edu

Genotype imputation is a vital tool in genome-wide association studies (GWAS) and meta-analyses of multiple GWAS results. Imputation enables researchers to increase genomic coverage and to pool data generated using different genotyping platforms. HapMap samples are often employed as the reference panel. More recently, the 1000 Genomes Project resource is becoming the primary source for reference panels. Multiple GWAS and meta-analyses are targeting Latinos, the most populous, and fastest growing minority group in the US. However, genotype imputation resources for Latinos are rather limited compared to individuals of European ancestry at present, largely because of the lack of good reference data. One choice of reference panel for Latinos is one derived from the population of Mexican individuals in Los Angeles contained in the HapMap Phase 3 project and the 1000 Genomes Project. However, a detailed evaluation of the quality of the imputed genotypes derived from the public reference panels has not yet been reported. Using simulation studies, the Illumina OmniExpress GWAS data from the Los Angeles Latino Eye Study and the MACH software package, we evaluated the accuracy of genotype imputation in Latinos. Our results show that the 1000 Genomes Project AMR + CEU + YRI reference panel provides the highest imputation accuracy for Latinos, and that also including Asian samples in the panel can reduce imputation accuracy. We also provide the imputation accuracy for each autosomal chromosome using the 1000 Genomes Project panel for Latinos. Our results serve as a guide to future imputation based analysis in Latinos.

**Keywords:** genotype imputation, Latino, HapMap Project, 1000 Genomes Project

**Table 1 | Phased haplotypes downloaded from the MACH website.**

Population	Code	HapMap phase 3	1000 Genomes project	
		Number of haplotypes	Number of haplotypes	Group code
Mexican ancestry in Los Angeles, California	MEX	104	132	AMR
Colombian in Medellin, Colombia	CLM		120	
Puerto Rican in Puerto Rico	PUR		110	
CEPH in Utah residents	CEU	234	174	EUR
Tuscans in Italy	TSI	176	196	
Finnish individuals from Finland	FIN		186	
British individuals from England and Scotland	GBR		178	
Iberian populations in Spain	IBS		28	
Yoruba in Ibadan, Nigeria	YRI	230	176	AFR
African ancestry individuals from Southwest, US	ASW		122	
Luhya in Webuye, Kenya	LWK		194	
Han Chinese in Beijing, China	CHB	168	194	ASN
Japanese in Tokyo, Japan	JPT	172	178	
Han Chinese South, China	CHS		200	
Total		1084	2188	

*The population labels were obtained from the HapMap and the 1000 Genomes Project websites.*

# NIH Public Access

## Author Manuscript

*Nat Genet.* Author manuscript; available in PMC 2012 February 9.

Published in final edited form as:

*Nat Genet.* ; 44(1): 6–7. doi:10.1038/ng.1044.

## Improved Imputation of Common and Uncommon Single Nucleotide Polymorphisms (SNPs) with a New Reference Set

Zhaoming Wang<sup>1,2</sup>, Kevin B. Jacobs<sup>1,2</sup>, Meredith Yeager<sup>1,2</sup>, Amy Hutchinson<sup>1,2</sup>, Joshua Sampson<sup>2</sup>, Nilanjan Chatterjee<sup>2</sup>, Demetrius Albanes<sup>2</sup>, Sonja I. Berndt<sup>2</sup>, Charles C. Chung<sup>2</sup>, W. Ryan Diver<sup>3</sup>, Susan M. Gapstur<sup>3</sup>, Lauren R. Teras<sup>3</sup>, Christopher A. Haiman<sup>4</sup>, Brian E. Henderson<sup>4</sup>, Daniel Stram<sup>4</sup>, Xiang Deng<sup>1,2</sup>, Ann W. Hsing<sup>2</sup>, Jarmo Virtamo<sup>5</sup>, Michael A. Eberle<sup>6</sup>, Jennifer L. Stone<sup>6</sup>, Mark P. Purdue<sup>2</sup>, Phil Taylor<sup>2</sup>, Margaret Tucker<sup>2</sup>, and Stephen J. Chanock<sup>2</sup>

Group	Populations					Illumina Array				
	European	American	African	American	African	Asian	Hap660	Hap1	Omni1	Omni2.5
ATBC		246						✓	✓	✓
CPSII		227						✓	✓	✓
PLCO		255						✓	✓	✓
PLCO				98				✓		✓
SHNX						74	✓			✓
HapMap										
CEU		116								✓
CHB						44				✓
JPT						44				✓
TSI		86								✓
YRI					59					✓
Total		930		98	59	162				





# Literatura

- **J. Marchini**, B. Howie, S. Myers, G. McVean and P. Donnelly (2007) *A new multipoint method for genome-wide association studies via imputation of genotypes*. Nature Genetics 39 : 906-913
- B Servin and **M Stephens**. Imputation-based analysis of association studies: candidate regions and quantitative traits. *PLoS Genet*, 3(7), Jul 2007.
- **Browning, S. R. and B. L. Browning**, 2007. Rapid and accurate haplotype phasing and missing data inference for whole genome association studies using localized haplotype clustering. American Journal of Human Genetics, 81:1084-1097
- **Browning, S. R. and B. L. Browning**, 2011. Haplotype phasing: existing methods and new developments. Nature Reviews Genetics, 12: 703-714.