

Population genetic and admixture analyses

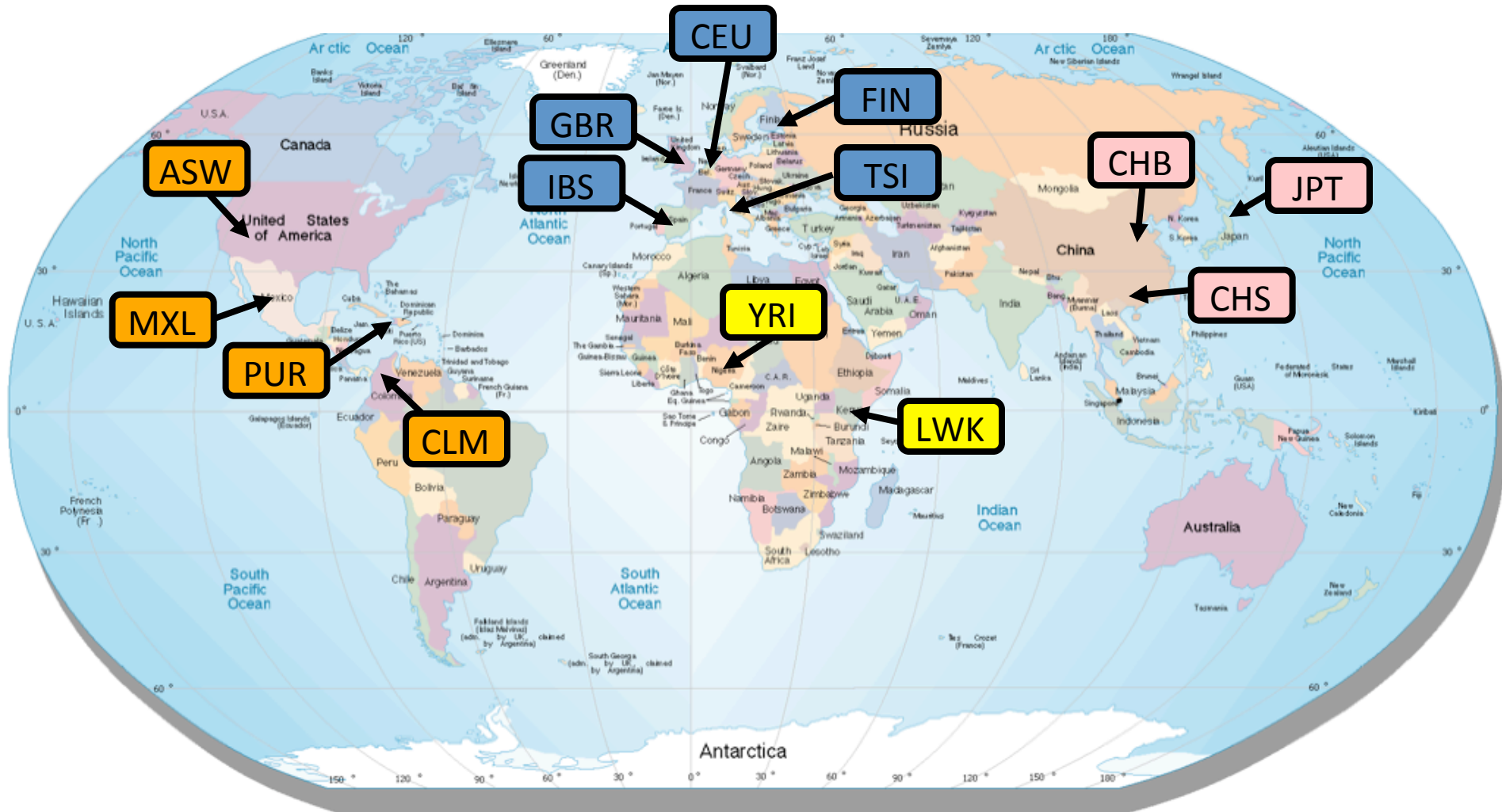
Eimear Kenny

Assistant Professor

Institutes of Personalized Medicine and Multiscale Biology at
Mount Sinai School of Medicine

Representing and presenting slides from the 1000 Genomes
Project Consortium

Phase1: 1092 samples from 14 populations

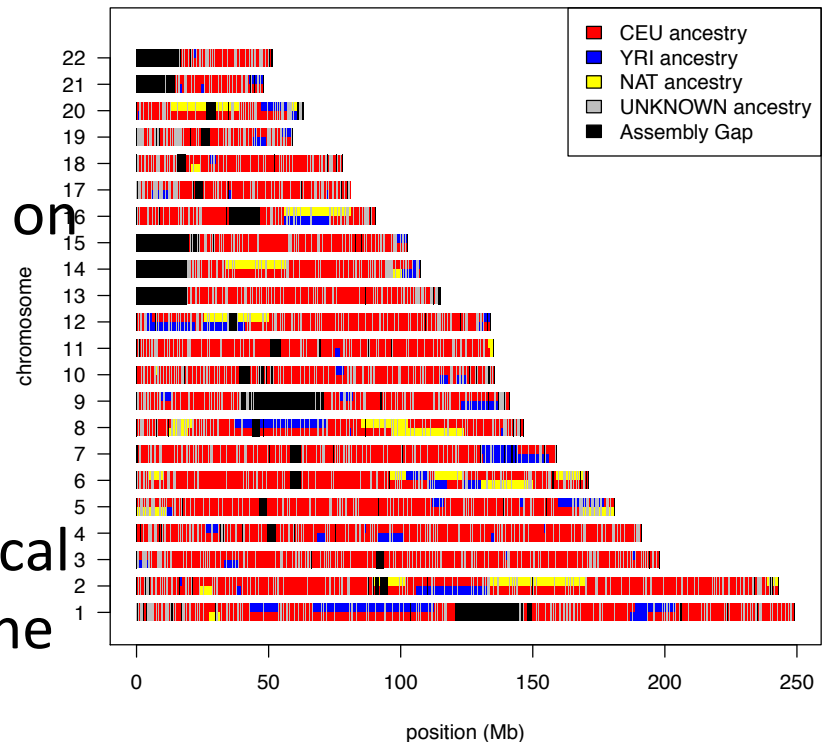


Motivation for Local Ancestry Inference

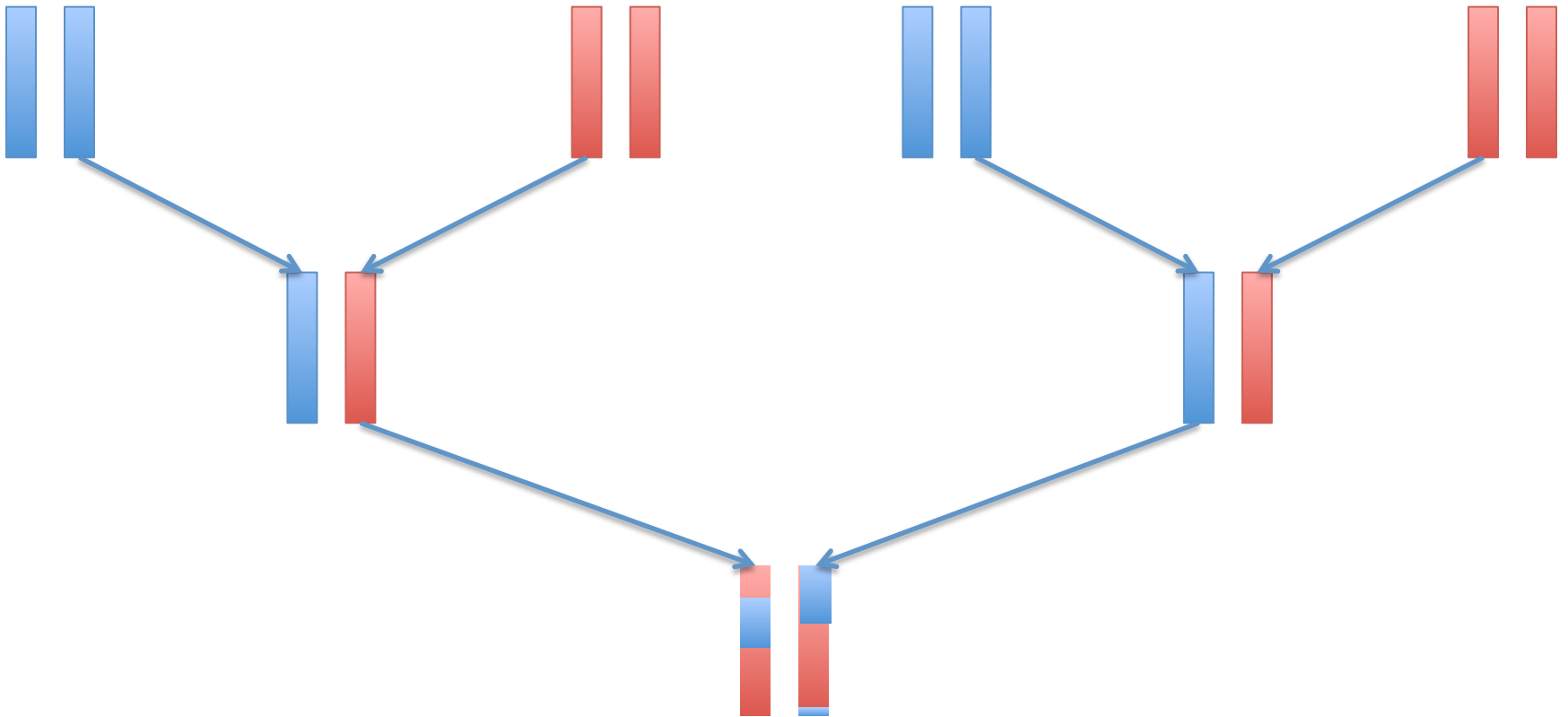
Using 1000 Genomes admixed populations to:

- study impact of recent admixture on patterns of genomic variation
- understand pre-historical demography
- inform variant discovery for medical genetics and personalized medicine

Ancestry Karyogram (Viterbi) for HG00551

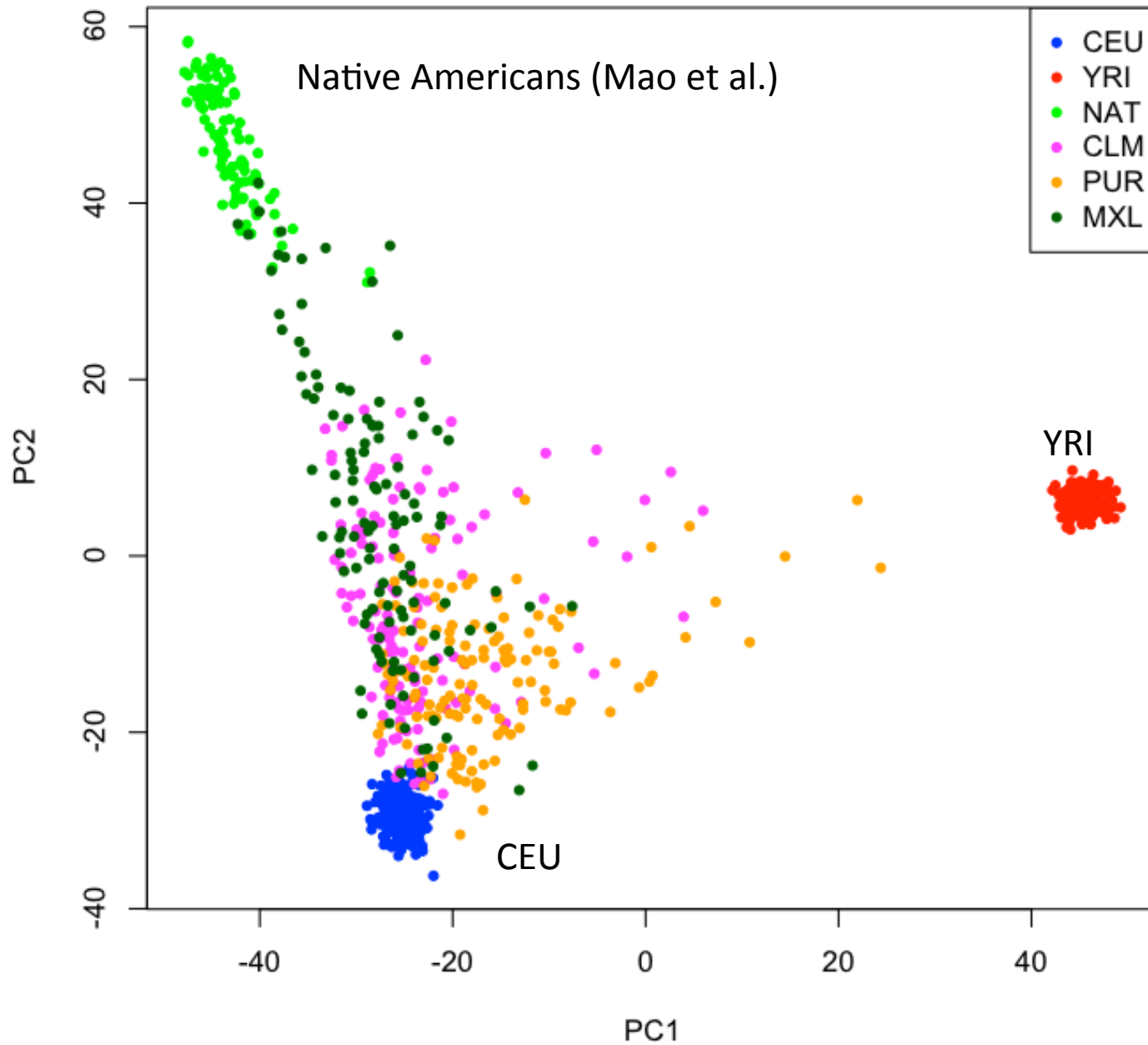


Admixture

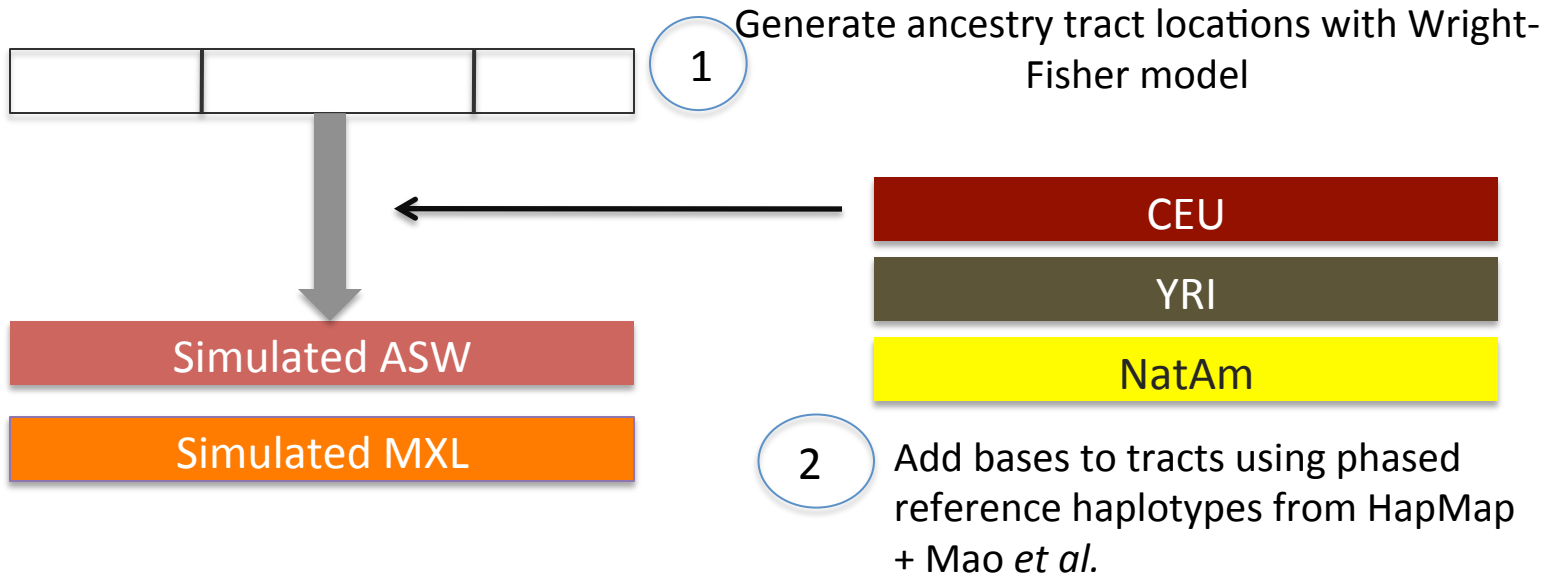


Recombination creates mosaic of ancestry

Phase 1 AMR Samples



Accuracy assessment using simulated haplotypes

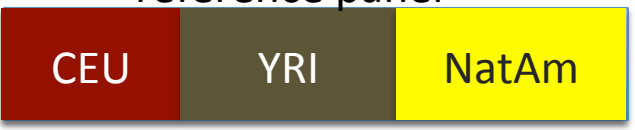




Simulated ASW

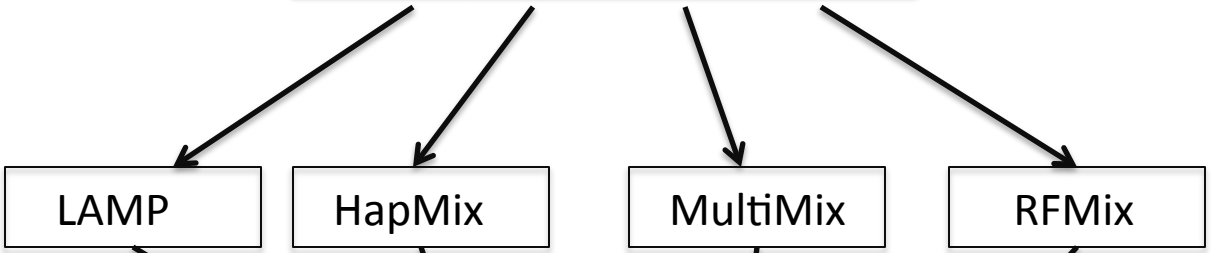
Simulated MXL

reference panel



3

Provide reference panel from other HapMap samples + Mao *et al*



4

Each group runs local ancestry estimation

Per site diploid calls

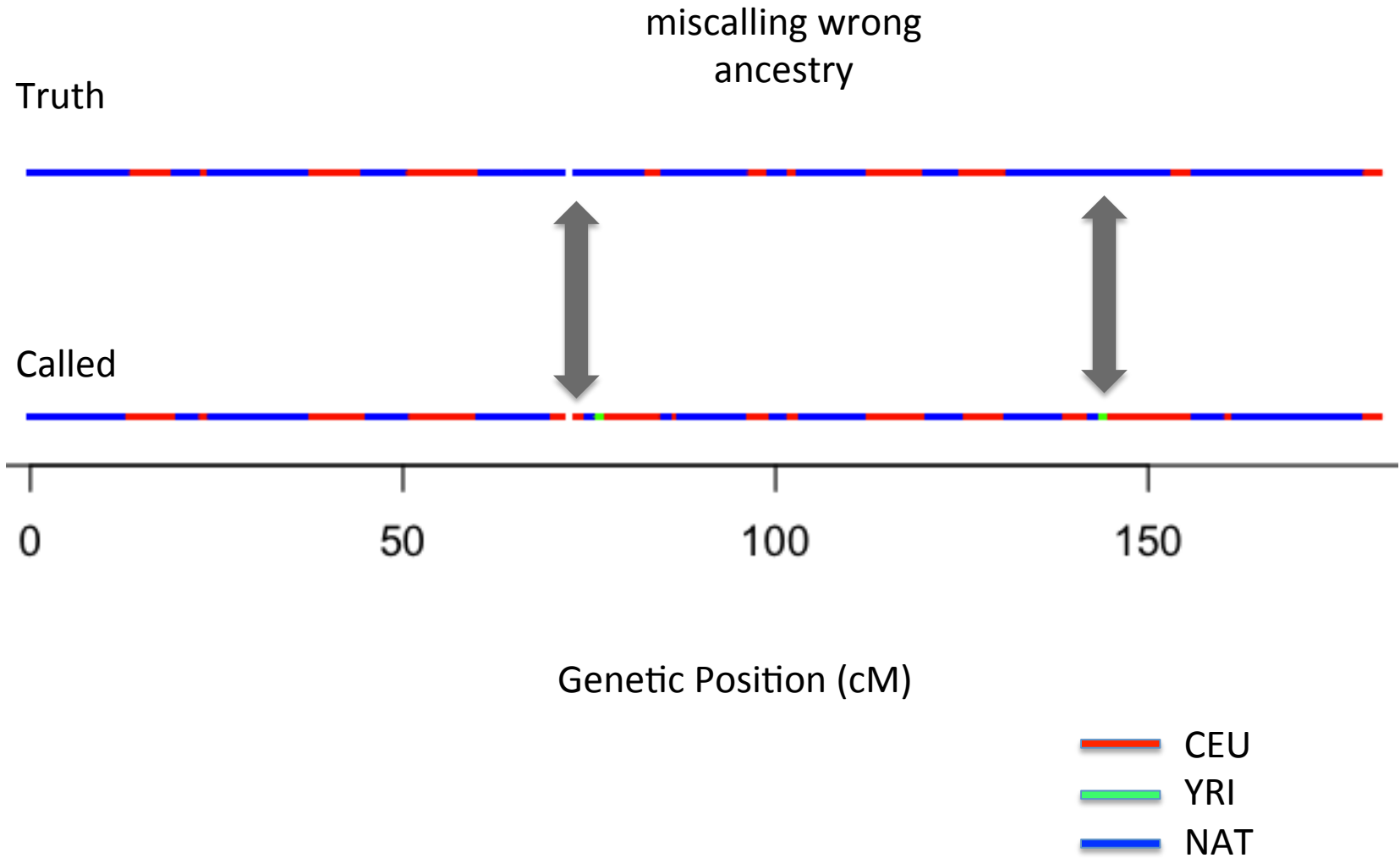
5

3- or 6-way diploid call per site

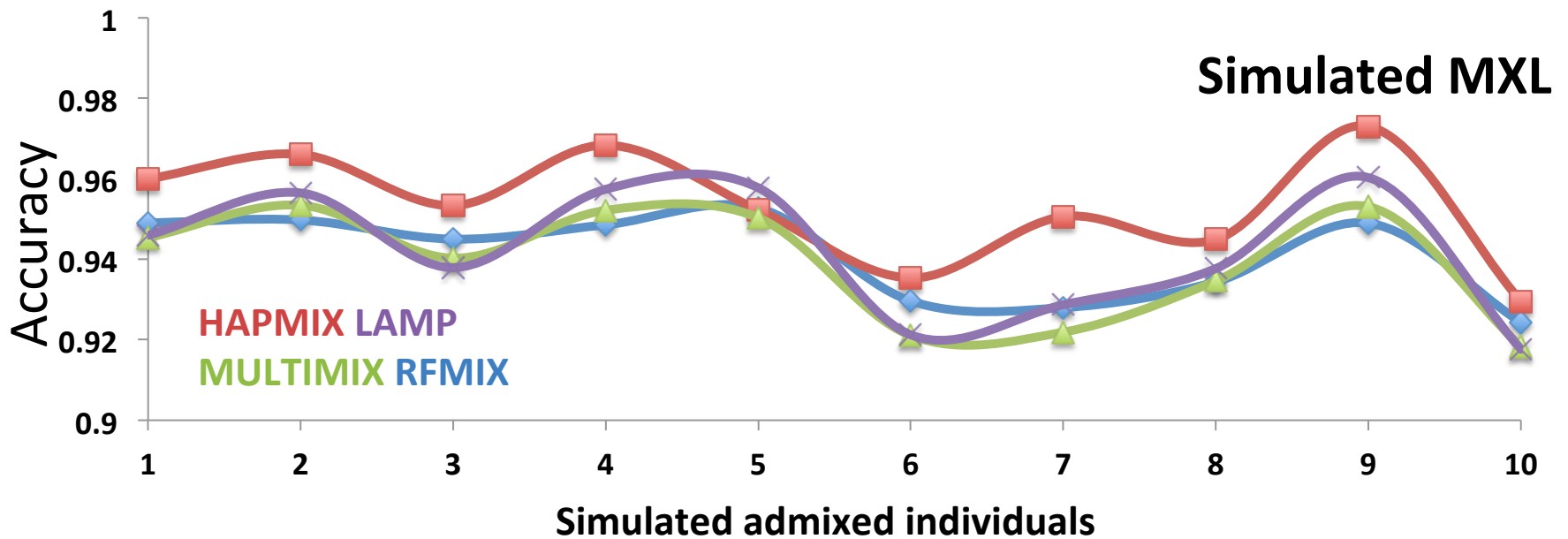
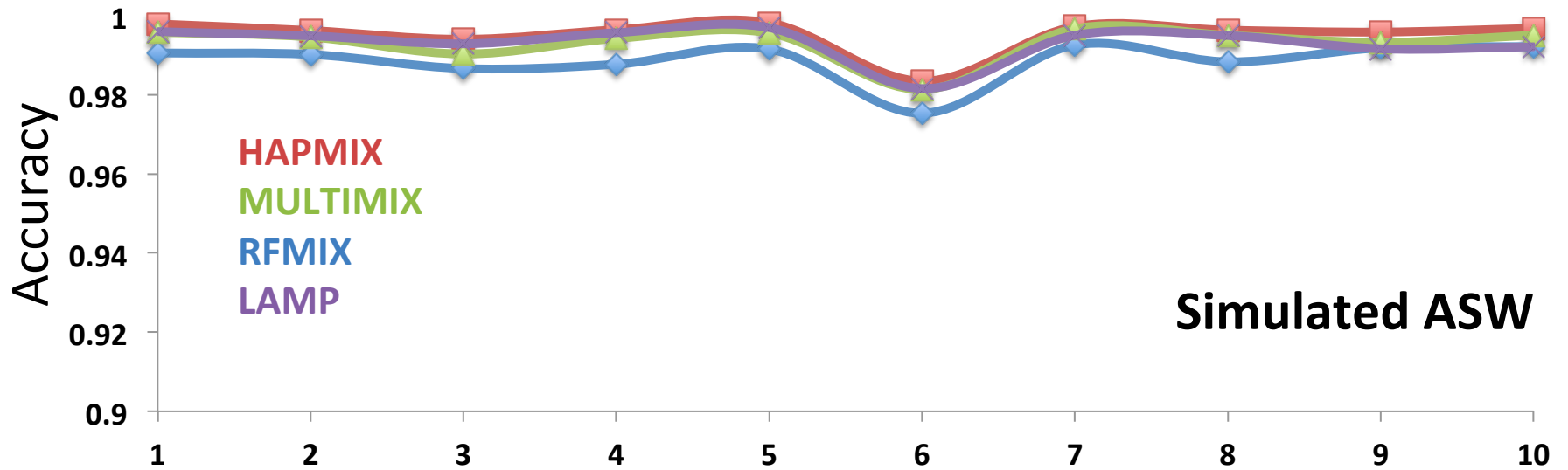


Compare ancestry calls

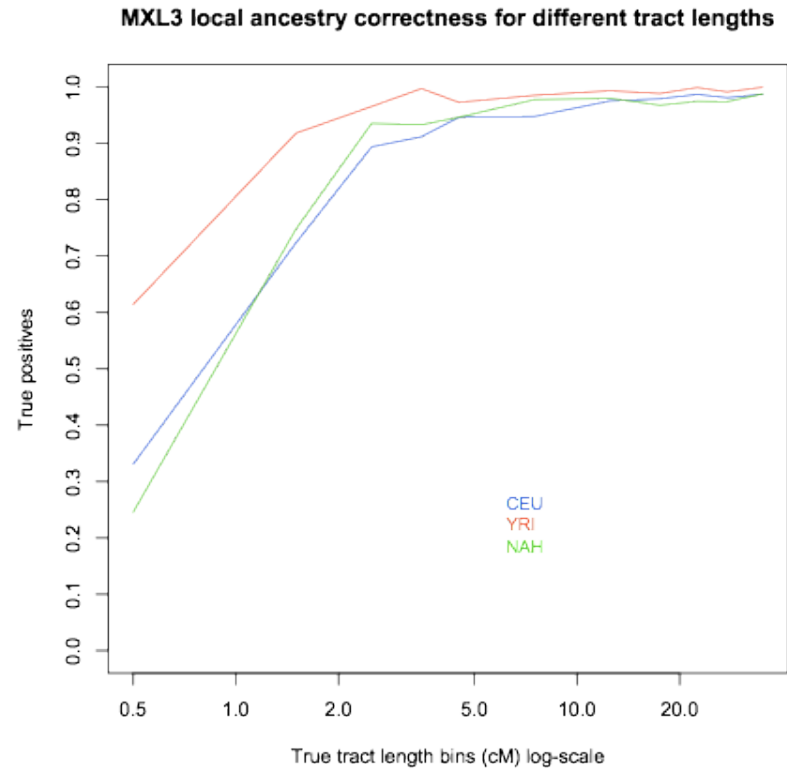
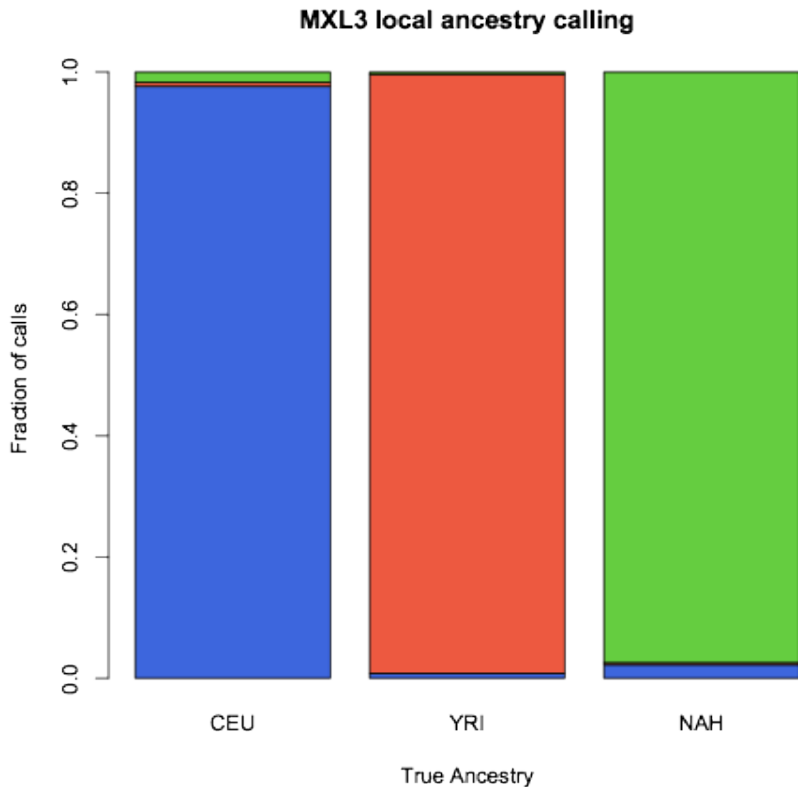
Example miscall



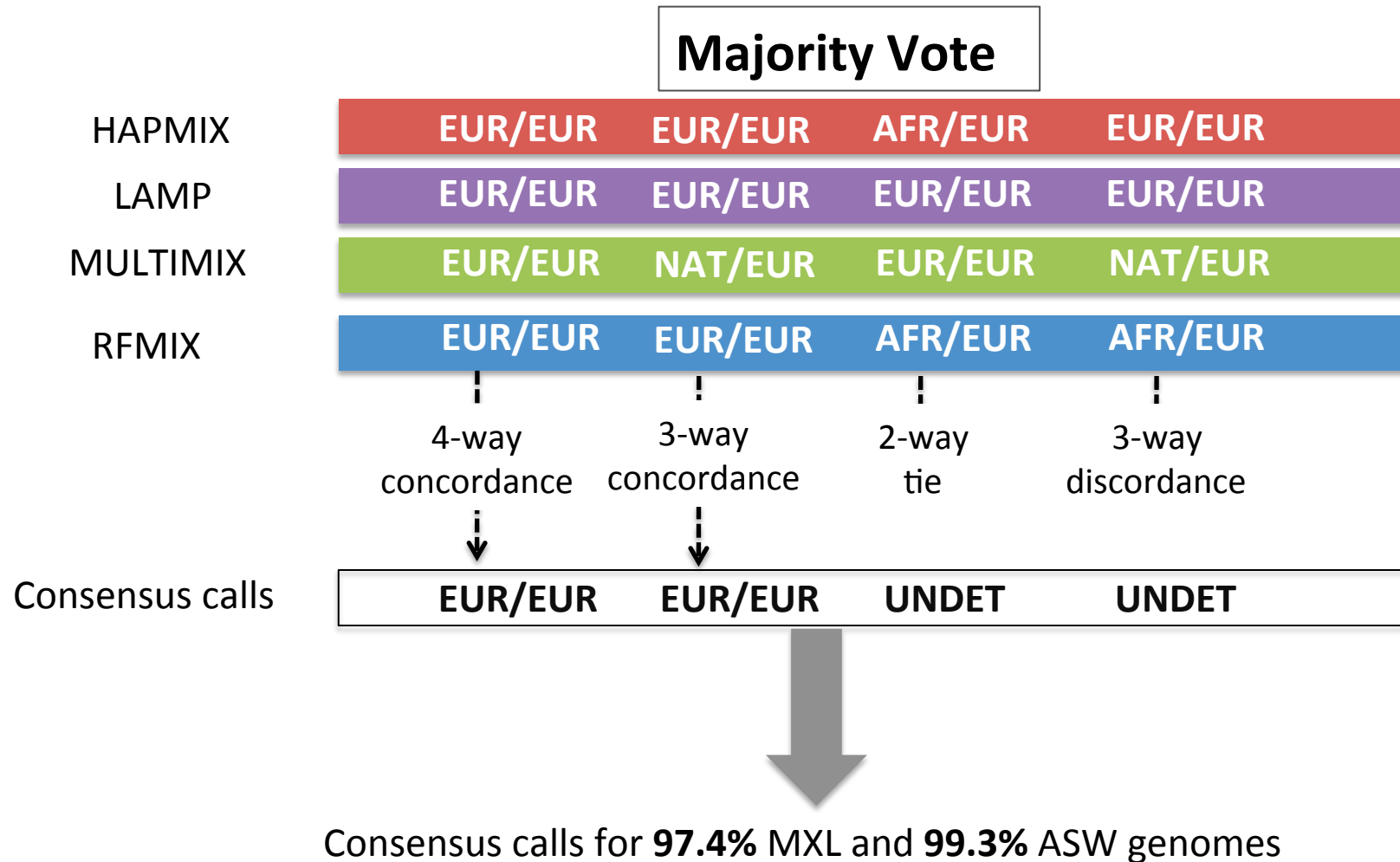
Results from the simulated data: some methods perform better than others



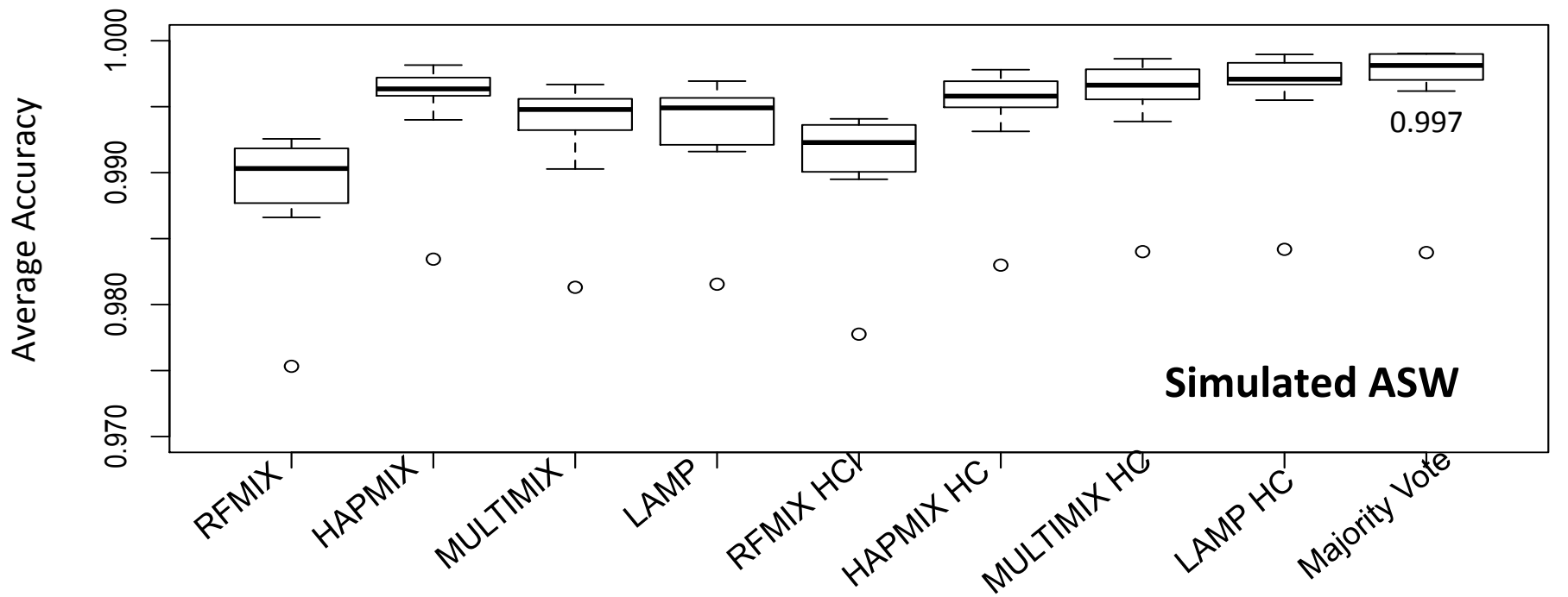
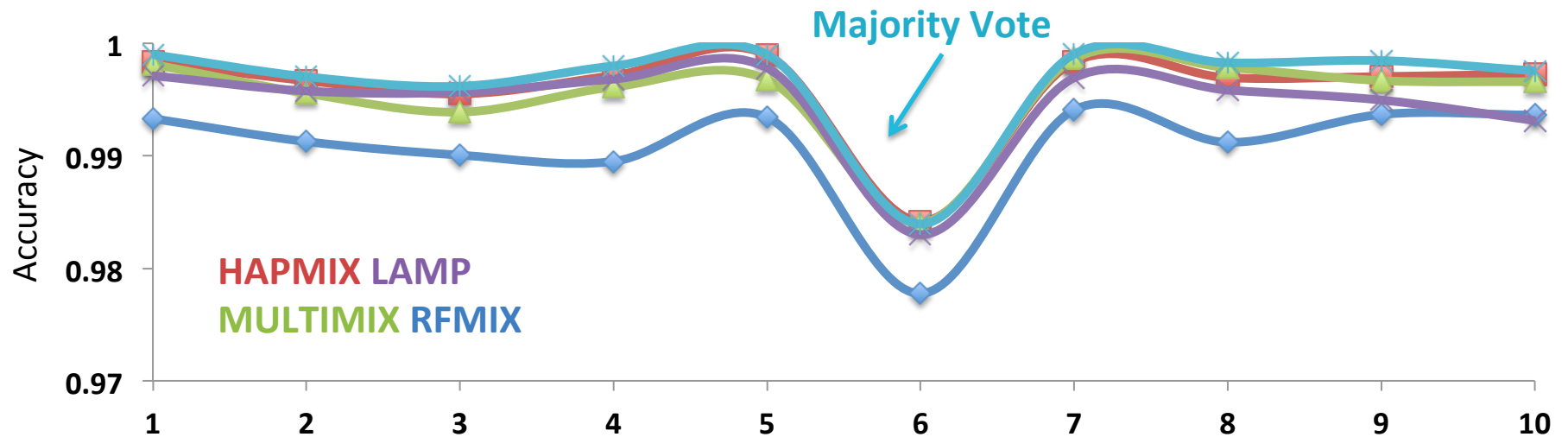
HapMix MXL ($K = 3$) Results



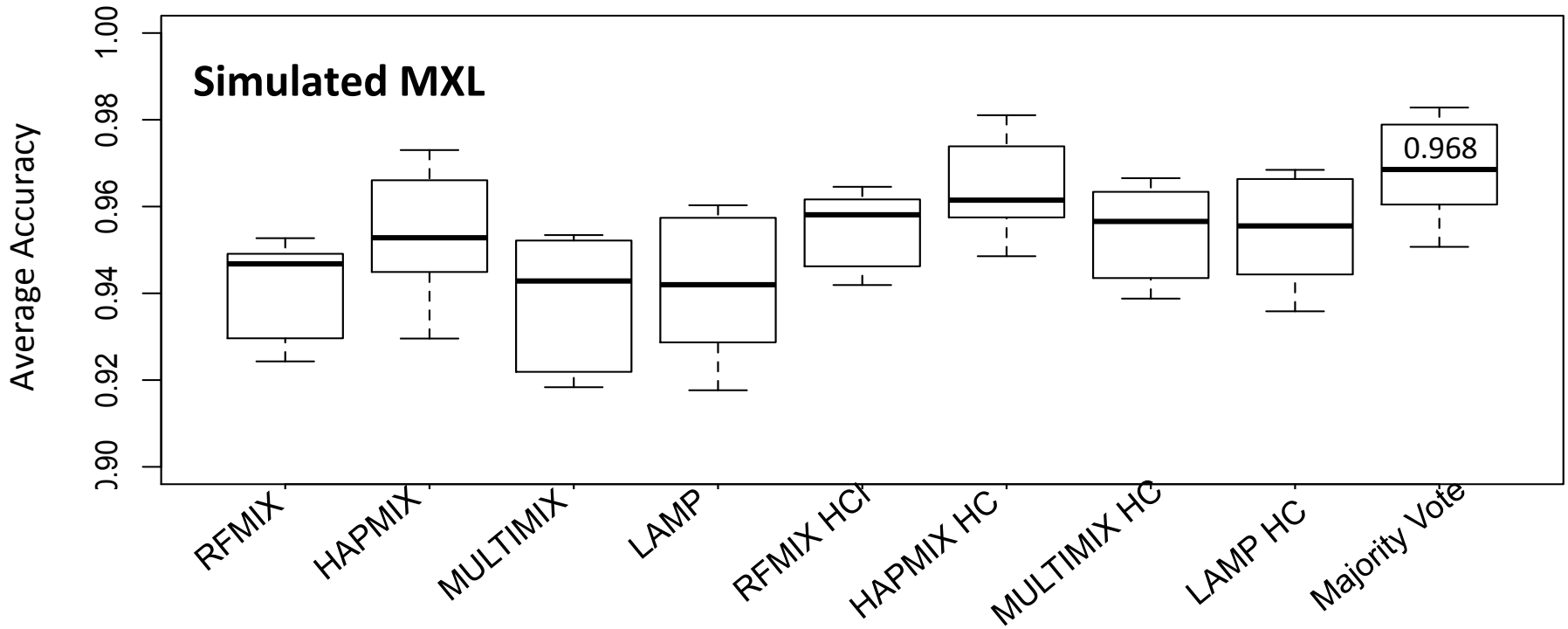
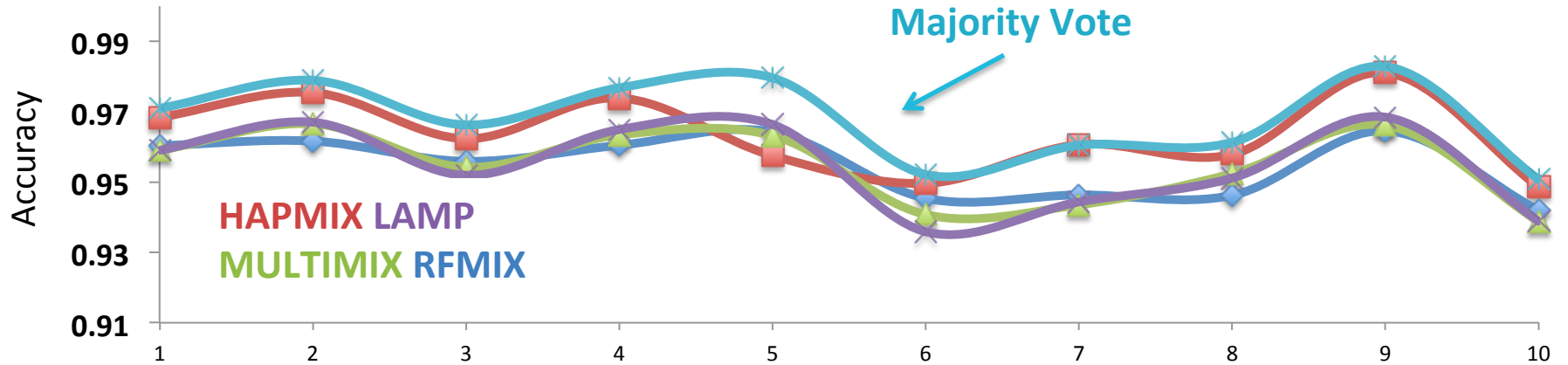
Strategy to obtain highest confidence calls for downstream analysis



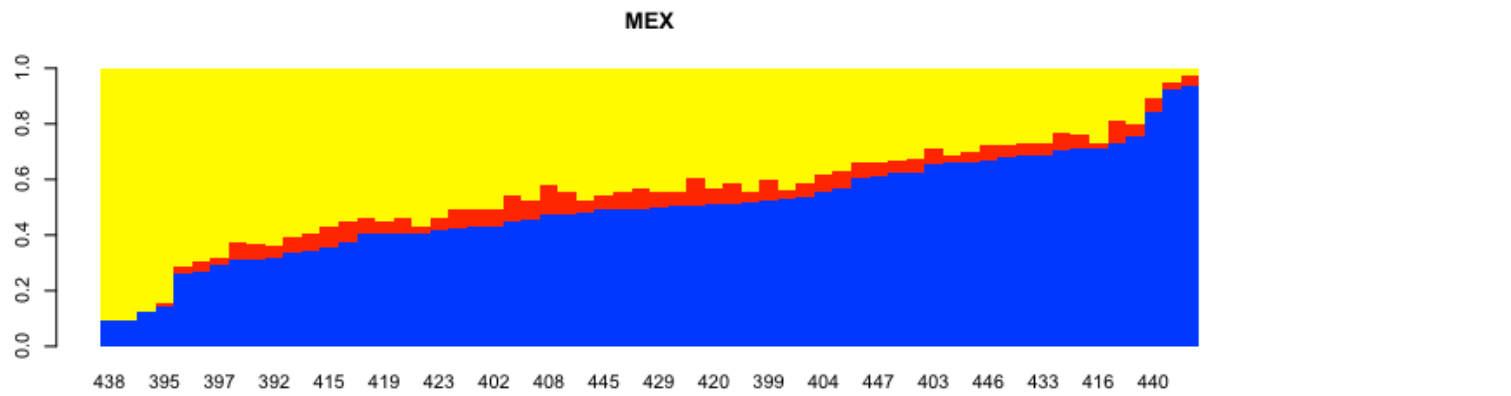
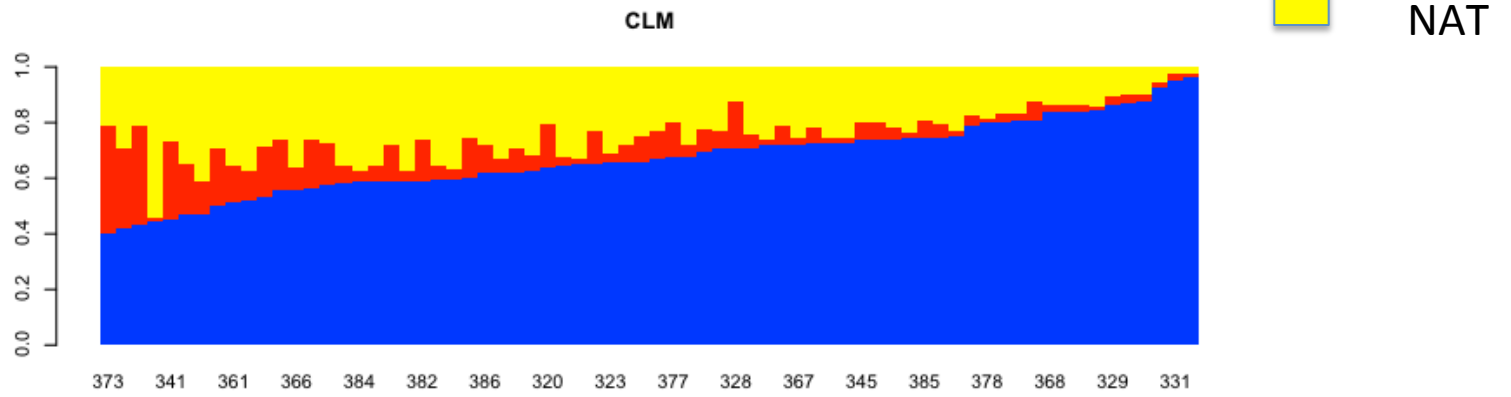
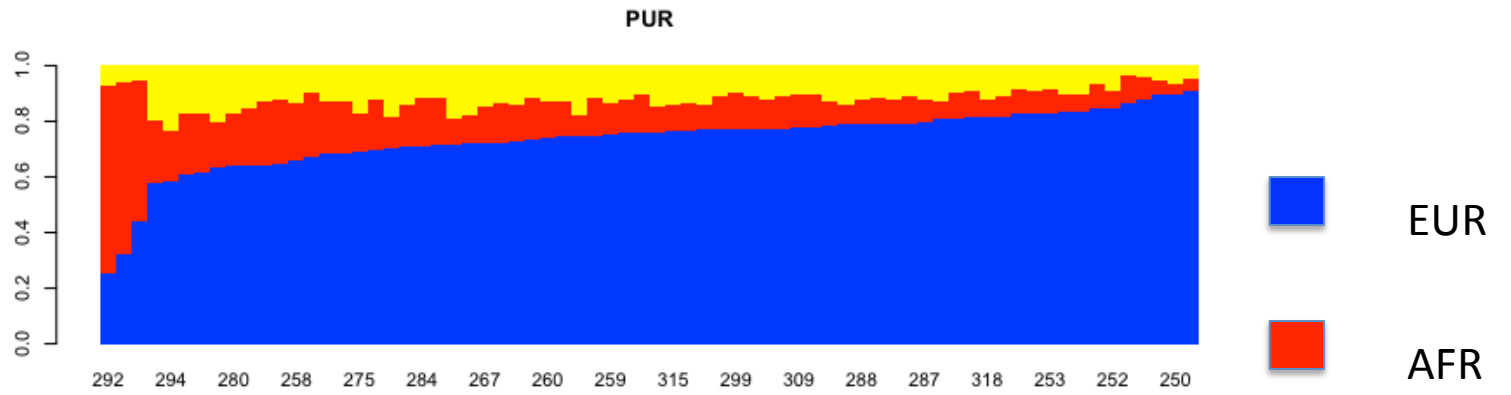
Marginal improvement for ASW



For downstream analysis select high confidence regions



Sequence diversity in admixed TGP samples

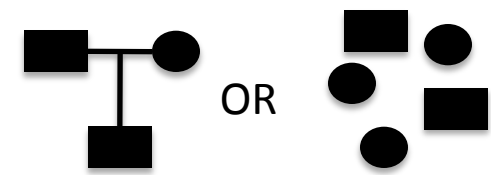


Inferred local ancestry calling in admixed samples



Select Affy6.0 sites
←

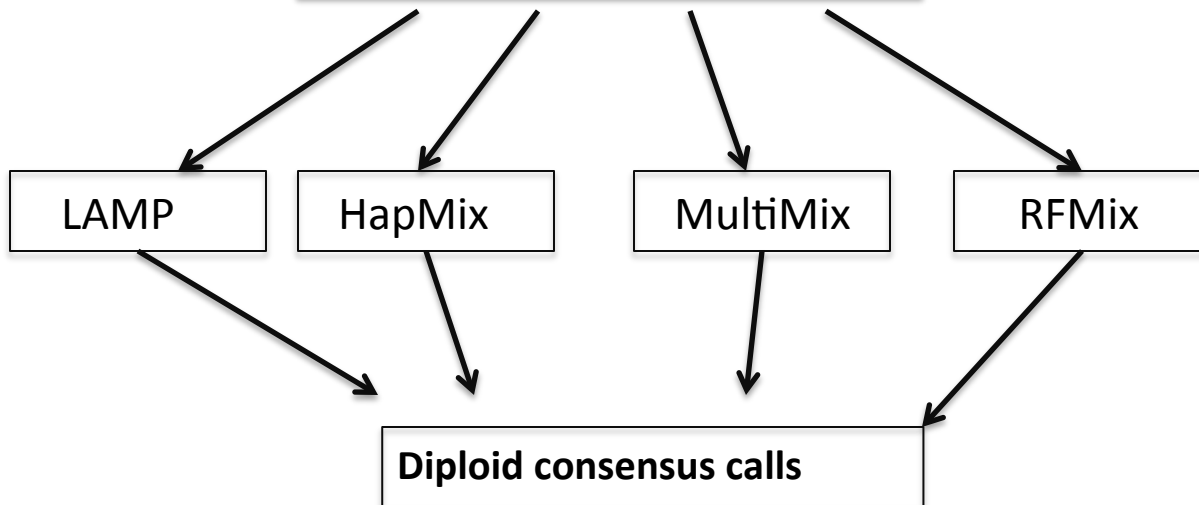
Phase Omni + Sequence data



Used Omni data for trio phasing

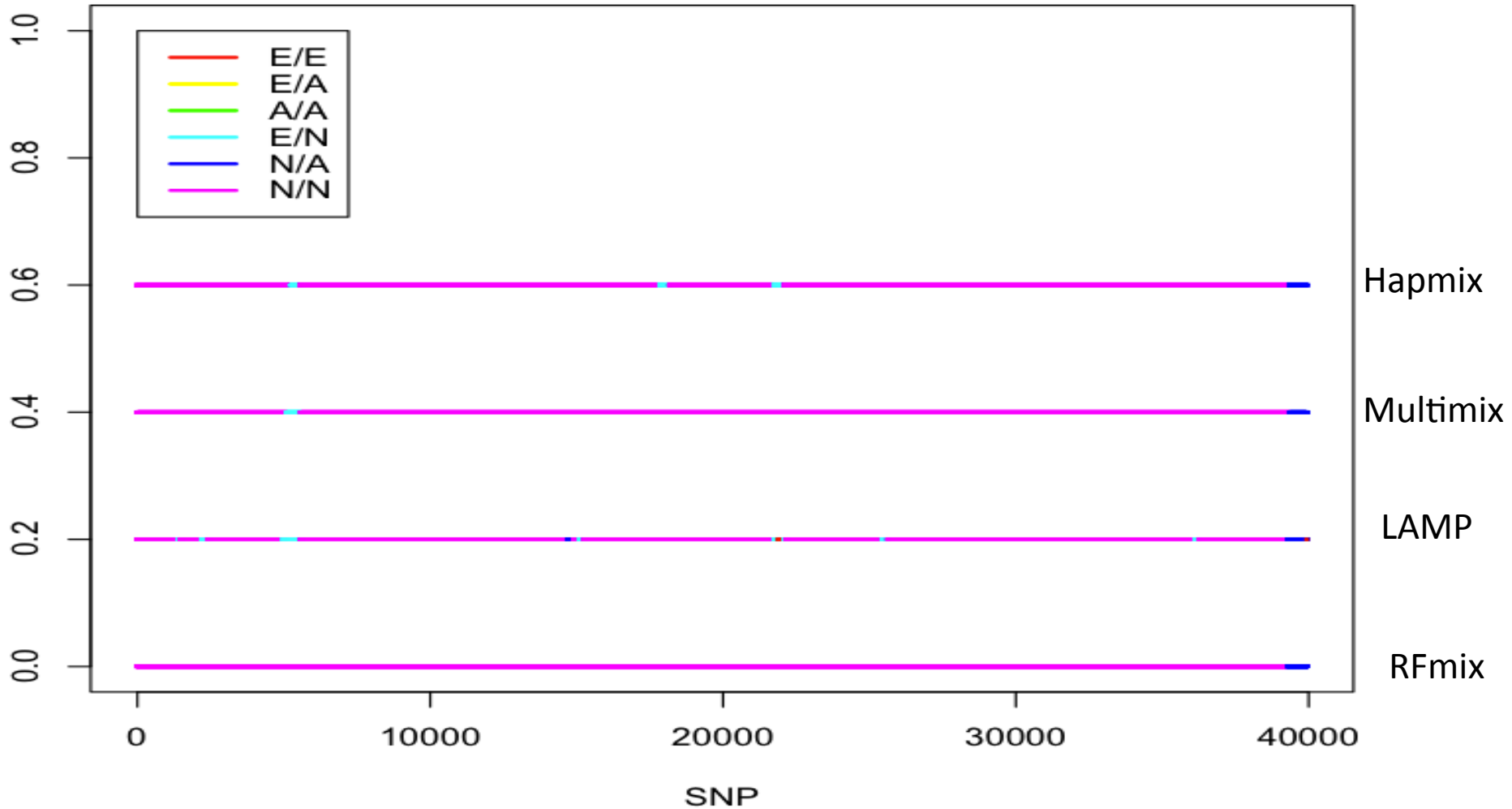
Phased unrelated (10 CLM + 1 MXL)

NatAm Affy 6.0 (Mao *et al*) CEU + YRI TGP



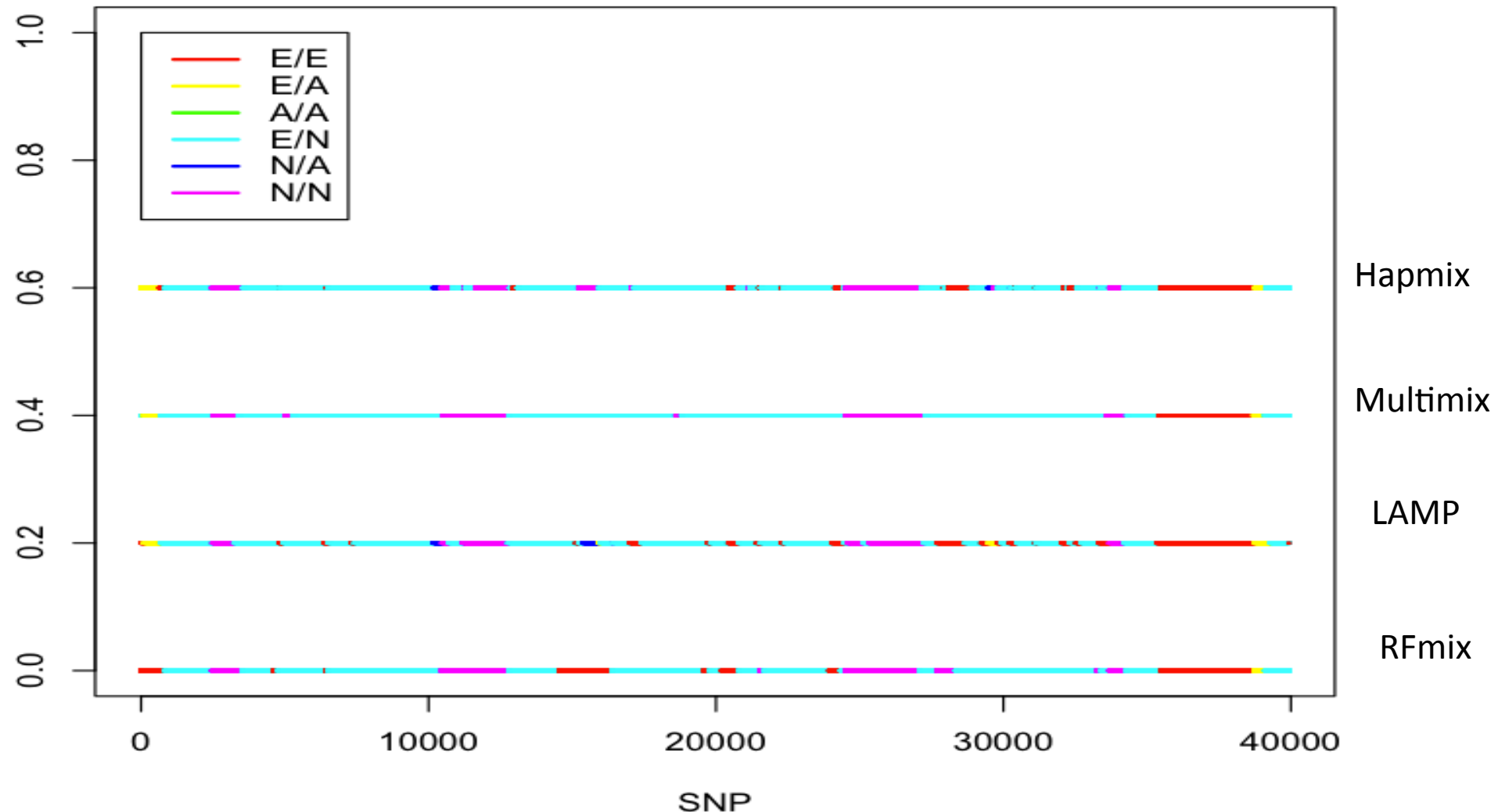
*Marchini (Shapeit) approach used for phasing all samples together

Good agreement



Not so good agreement

(Probably best to “mask” and not make call...)

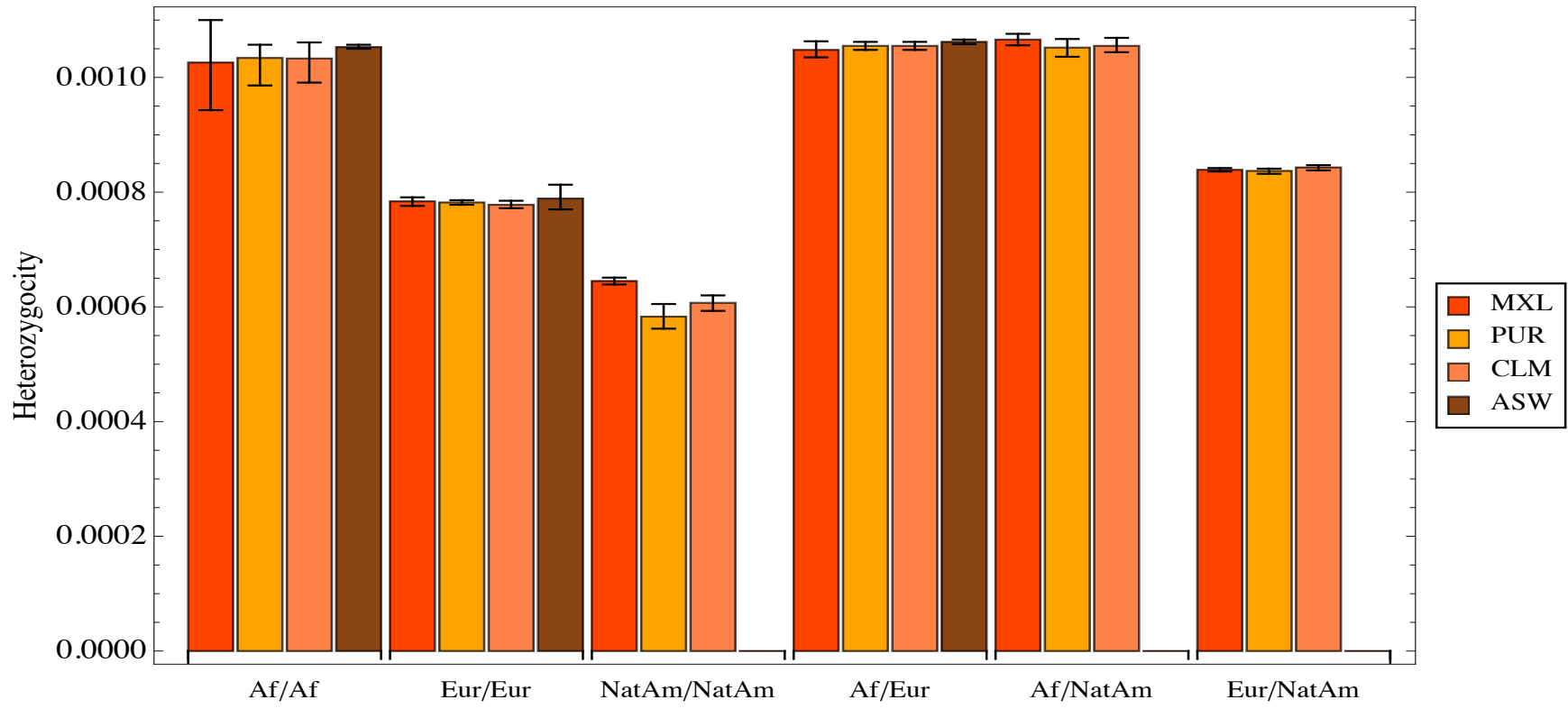


Summary

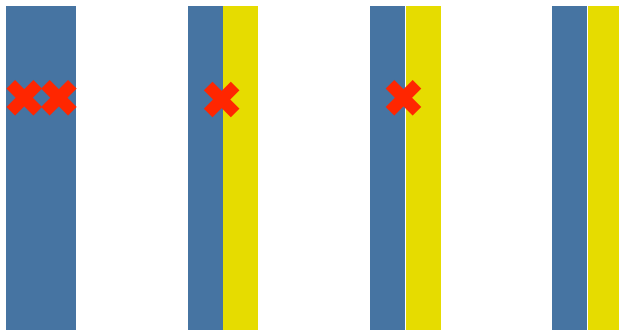
- Admixture deconvolution produced for all Phase 1 samples and pushed to DCC
 - Merged calling improved accuracy
 - Estimated 99%+ accuracy for ASW and 96-97% for MXL, CLM, PUR samples
- Local ancestry calls are available as bed files per individual on the DCC

Some cool applications

Heterozygosity per population



Inferring “ancestral” allele frequencies

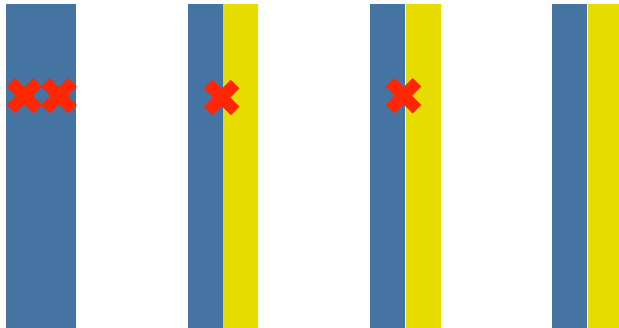


What are the allele frequencies in each population?

-Bayesian answer

$$P(\mathbf{f}|D) = \frac{P(D|\mathbf{f})P(\mathbf{f})}{\int d\mathbf{f}'P(D|\mathbf{f}')P(\mathbf{f}')}$$

Inferring “ancestral” allele frequencies



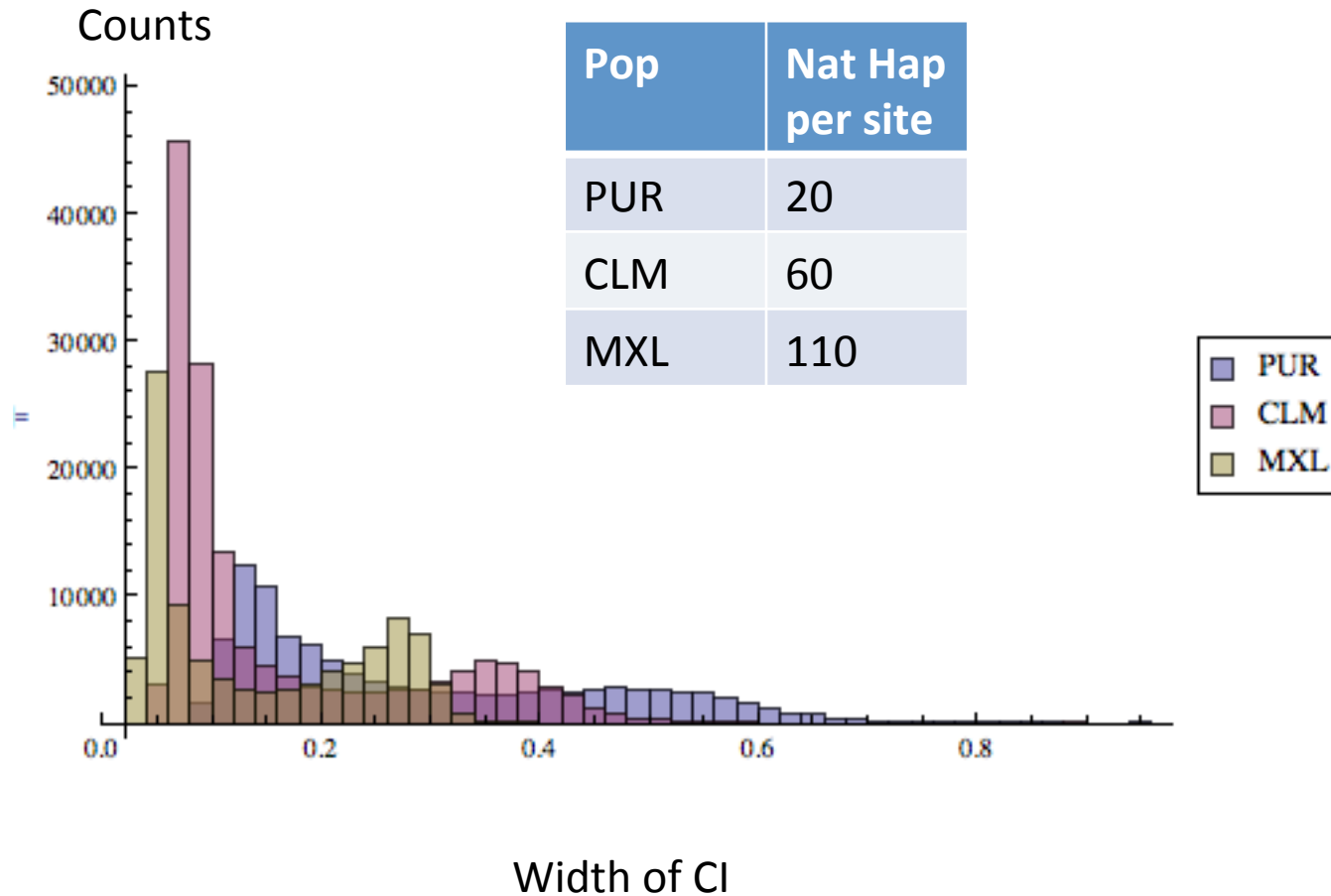
What are the allele frequencies in each population?

-Bayesian answer

$$P(\mathbf{f}|D) = \frac{P(D|\mathbf{f})P(\mathbf{f})}{\int d\mathbf{f}'P(D|\mathbf{f}')P(\mathbf{f}')}$$

Calculate across all sites

Confidence interval width varies by population



Summary

- SNP diversity, novelty rate, and Non-Syn/Syn ratios mirror demography
- Recover allele frequencies in ancestral populations
- Samples and pipelines are useful beyond TGP

Admixture Working Group

- Bustamante lab (Eimear Kenny, Simon Gravel, Fouad Zakharia, Brian Maples)
- Marchini lab (Claire Churchouse)
- Halperin lab (Yael Baran)
- Myers lab (Anjali GuptaHinch)
- Burchard (Chris Gignoux)
- Abigail Bigham and Mark Shriver (Mao et al. Affy 6.0 data)

Credits



More information at www.1000genomes.org