# Split populations (after subject QC)

N=485,595
M=16,703,829
PLINK hardcall
After subject QC

- M same for all population samples
- Population samples other than EUR may contain SNPs with MAF~0
- Additional SNP filtering is typically performed within each project (to ensure appropriate MAF threshold for project/subsample)

**EUR**
N=460,527
**(387,614 unrelated)**

**SAS**
N=10,427
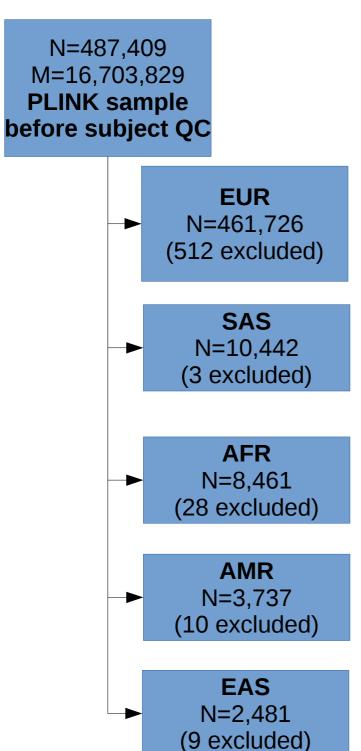(9,653 unrelated)

**AFR**
N=8,439
(7,840 unrelated)

**AMR**
N=3,726
(3,527 unrelated)

**EAS**
N=2,476
(2,415 unrelated)

# Ancestry assignment procedure

N=487,409
M=16,703,829
**PLINK sample
before subject QC**

**EUR**
N=461,726
(512 excluded)

**SAS**
N=10,442
(3 excluded)

**AFR**
N=8,461
(28 excluded)

**AMR**
N=3,737
(10 excluded)

**EAS**
N=2,481
(9 excluded)

**Assign 1kg ancestry**
- restrict to ukb-1kgph3 overlapping SNPs (650,232 SNPs)
- apply filters ( 593,693 SNPs)
    - maf_ukb < 0.001
    - hwe_ukb < 0.000001
    - remove longrange LD regions (Price et al., 2008)
    - remove non-HRC SNPs
    - merge ukb and 1kg and restrict to nonmissing SNPs
-prune SNPs pairwise
    - window = 1500kb
    - step = 150snps
    - r2 = 0.1
(145,692 SNPs)
- compute 30 PCs in 1kg
- project ukb subjects on these PCs
- assign ancestry to closest 1kg subpopulation
- exclude subjects with Mahalanobis distance > 6 S.D. from subpop average