

Evaluation of low-pass whole genome sequencing as a cost-effective and improved alternative to arrays

Peng Zhang, Hua Ling, Justin Paschall, Beth Marosy, Jessica Gearhart, Kimberly Doheny

Center for Inherited Disease Research (CIDR), Johns Hopkins Genomics, Institute of Genetic Medicine, The Johns Hopkins School of Medicine

Contact: Peng Zhang (pzhang13@jhmi.edu), Kimberly Doheny (kdoheny@jhmi.edu)

#215

Introduction

We tested using low-pass whole genome sequencing (lpWGS) as a cost-effective, and potentially improved alternative for SNP arrays in genome-wide association studies (GWAS).

We designed and analyzed an experiment with 92 samples to evaluate aspects of practical implementation including: library preparation methodology, impact across ancestry groups, lpWGS coverage levels, DNA source, imputation algorithm, imputation reference and compute resources.

The validation ‘truth’ dataset was created by deep sequencing 92 samples with diverse ancestry background (down sampled to the same coverage of 28X). The lpWGS ‘test’ dataset (“lps”) was generated using a cost-effective and high-throughput library prep (plexWell LP384 from seqWell) for the same set of 92 samples to target 1X coverage, then imputed to the 1000 Genomes (1KG) deep sequencing reference panel (NYGC) using GLIMPSE.

The SNP array ‘test’ dataset (“GSA”) was in silico array data derived by extracting the Illumina Global Screening Array (GSA, 654K variants) genotypes from the deep WGS validation dataset, and then performing imputation with the TOPMed reference panel, to mimic current standard practice for GWAS studies.

Library Prep & Data Generation

Validation (‘truth’) deep WGS dataset: PCR-free library was created with 500-750ng of genomic DNA, sheared , and processed using the Kapa Hyper Prep kit (Roche) for End-Repair, A-Tailing and Ligation of IDT (Integrated DNA Technologies) unique dual-indexed adapters according to the Kapa protocol. The Illumina NovaSeq 6000 platform was used to generate raw sequencing data (2x150). Coverage normalized to 28x, alignment and variant calling was performed using the DRAGEN 3.7.5 pipeline on the Illumina BaseSpace cloud platform.

lpsWGS dataset: Libraries were generated using 10ng of genomic DNA. Samples underwent a transposase-based library prep using the plexWell LP384 kit (SeqWell) which attaches a unique sample (i7) barcode directly into the input DNA (Fig 1). After pooling (92->1), samples undergo a second barcoding step which inserts the secondary (i5) unique pool barcode into each sample within the pool. Dual barcoded, pooled libraries were amplified using 8 cycles followed by a bead-based size selection. Sequencing reads and alignment were done the same as the deep WGS with a targeted coverage of 1x.

GLIMPSE imputation: We used BCFtools (v1.9) to compute the Genotype Likelihoods (GLs) for all bi-allelic sites in the reference panel using the lps reads in CRAM format. GLIMPSE (v1.1.1) then used those GLs as input to impute the genotypes for all the SNVs from to the 1KG deep sequencing reference panel (NYGC).

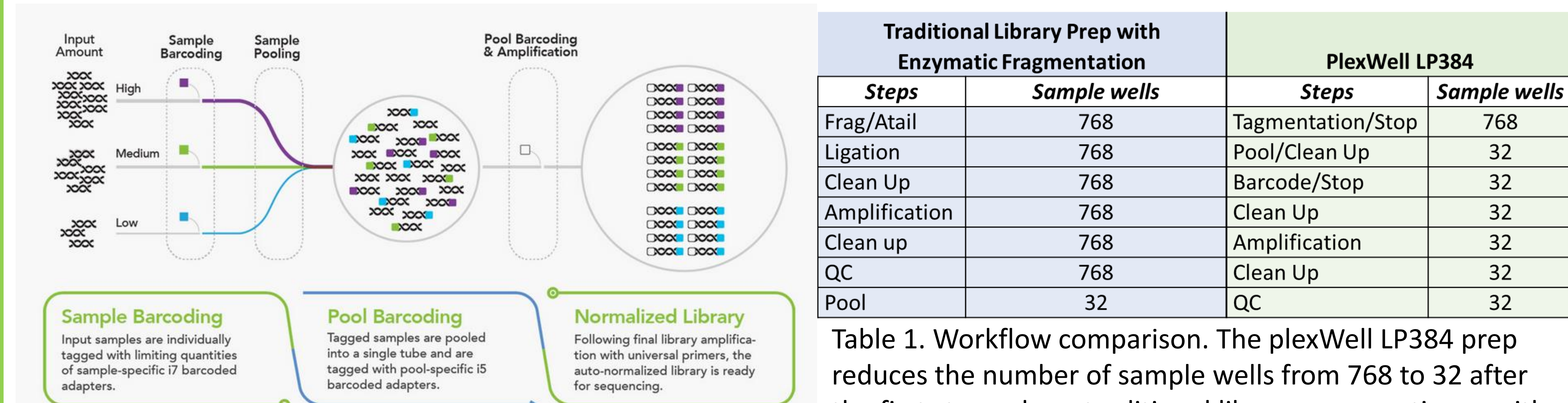
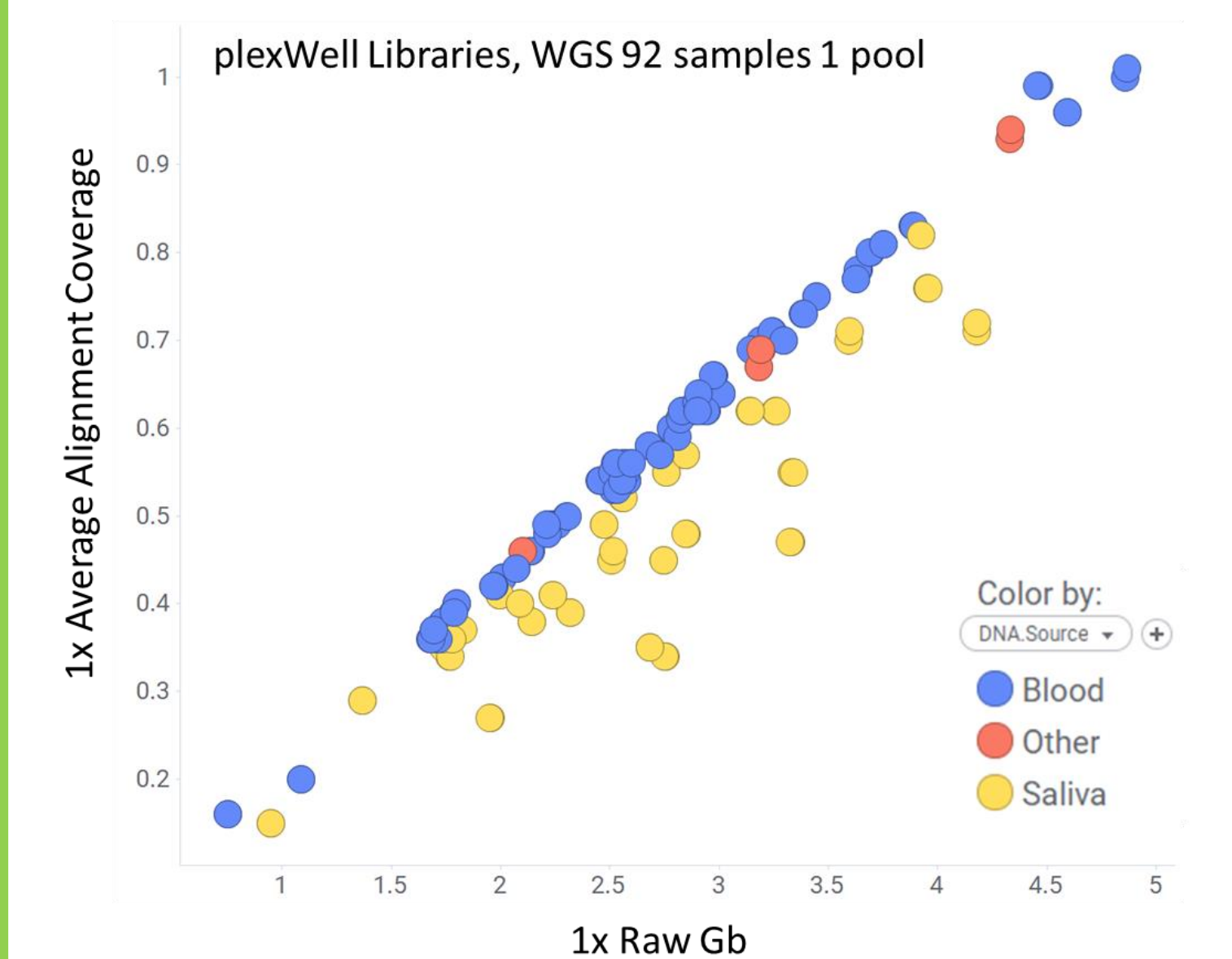


Fig 1. plexWell LP384 transposase-based library prep chemistry



Genetic Ancestry Determination: We determined five genetic ancestry groups by performing principal component analysis (PCA) on the 92 samples (GSA) along with an external set of 1201 unrelated HapMap 3 samples genotyped on the Illumina 1M array. Results: 9 African, 6 Admix_CEU_African, 69 European, 7 Admix_CEU_Asian and 1 Asian.

The plexWell LP384 prep workflow saves time and cost:

Time/Labor saving: 1 day versus 3 days, automation required for processing 768 samples only for 1st step versus all steps.

Cost saving: reduced 24-fold of cost by collapsing down to 32 sample wells early without the need to purchase additional equipment specifically for lp384 kit.

Fig 2. The plexWell chemistry uses a reagent limiting initial transposition step which normalizes the sample input into pooling. Here we show plots the 1x average alignment coverage (y-axis) vs the 1x Raw Gb (x-axis) for each of the 92 samples processed into 1 pool. Input sample DNA was derived from both blood and saliva, as saliva is a common DNA source used in array studies.

Contamination Detection

To determine if we could detect sample mixtures, we derived in silico mixtures from the lps datasets. We selected six sample pairs to mix, with the same or similar genetic ancestry and similar level of coverage (within 0.1x). Sample mixtures were detectable at a 5-10% range (dependent on coverage), similar to detection limits for arrays.

0.5X			1X		
Mixture	Contamination prediction percent	MT avg cov	Mixture	Contamination prediction percent	MT avg cov
50%	42.90%	67	50%	48.10%	138
30%	29.40%	68	30%	32.40%	137
15%	16.50%	71	15%	17.20%	138
10%	10.30%	69	10%	12.10%	138
7%	None	NA	5%		135
5%	None	NA	2%	None	NA
2%	None	NA			

Main References:

- GLIMPSE: Rubinacci Simone et al. (2021) *Nature Genetics*, 153:120-126
- HaploCheck: Weissensteiner H. et al. (2021) *Genome Research*, 31: 309-316
- GWASTools: Gogarten SM et al. (2012) *Bioinformatics*, 28(24), 3329-3331.
- Summary metrics FPR etc: Kishikawa T et al. (2019) *Scientific Report*, 9:1784

Analyses & Results

Table 3. Sensitivity and other metrics from different filters compared to validation (average #SNVs: 4,043,555) dataset.

Type	#SNVs (m)	Sensitivity	FPR	FNR	NTPR
lps.glimpse.all	61,715,567 (100%)	0.897	0.140%	0.157%	96.751%
lps.glimpse.GP	61,294,115 (99.32%)	0.85	0.055%	0.070%	98.626%
lps.glimpse.INFO	59,218,836 (95.95%)	0.825	0.116%	0.067%	97.036%
lps.glimpse.GP.INFO	59,051,021 (95.68%)	0.793	0.050%	0.035%	98.662%
lps.glimpse.HI	59,791,338 (96.88%)	0.852	0.083%	0.054%	98.150%
GSA.all	292,323,483 (100%)	0.841	0.092%	0.022%	92.303%
GSA.R2	16,018,114 (5.48%)	0.813	1.632%	0.301%	92.244%

Results shown are average across 92 samples and all autosomal chromosomes (chr1-22). Filter.type: for lps, we applied GP≥0.9, INFO≥0.3, HI was defined as INFO≥0.3 for MAF<0.05 and GP≥0.9 for MAF≥0.05. For GSA, R2≥0.3. Sensitivity here is the non-reference sensitivity. **Note:** Validation data have no reference homozygote (0/0)

Sensitivity = (e+f+h+i)/(b+e+f+h+i+k)
FPR = FalsePosRate = (d+h)/m
FNR = FalseNegRate = (b+f)/m
NTPR = Non-referenceTruePositiveRate = (e+i)/(d+e+f+h+i)
m is the #SNVs in test datasets
(a+b+d+e+f+h+i)

		Validation WGS(28X)		
Lps or GSA	0/0	NA	0/1	1/1
	a	b		
	0/1	d	e	f
	1/1	h	i	
No call	NA	k		

Fig 3. Concordance (squared Pearson correlation) from different filters compared to the validation dataset.

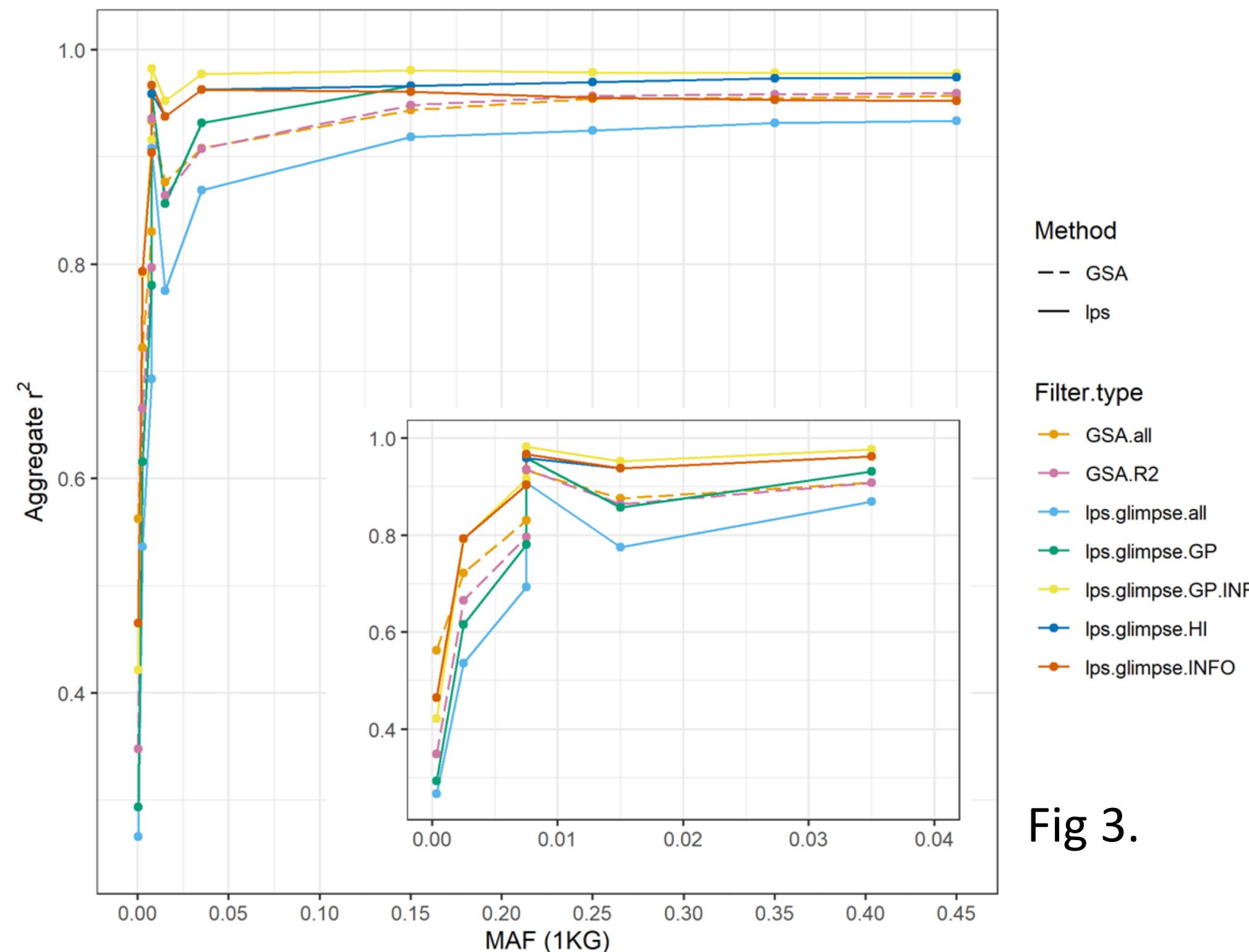


Fig 3.

Fig 4. Sensitivity by genetically defined ancestry for lps and GSA.

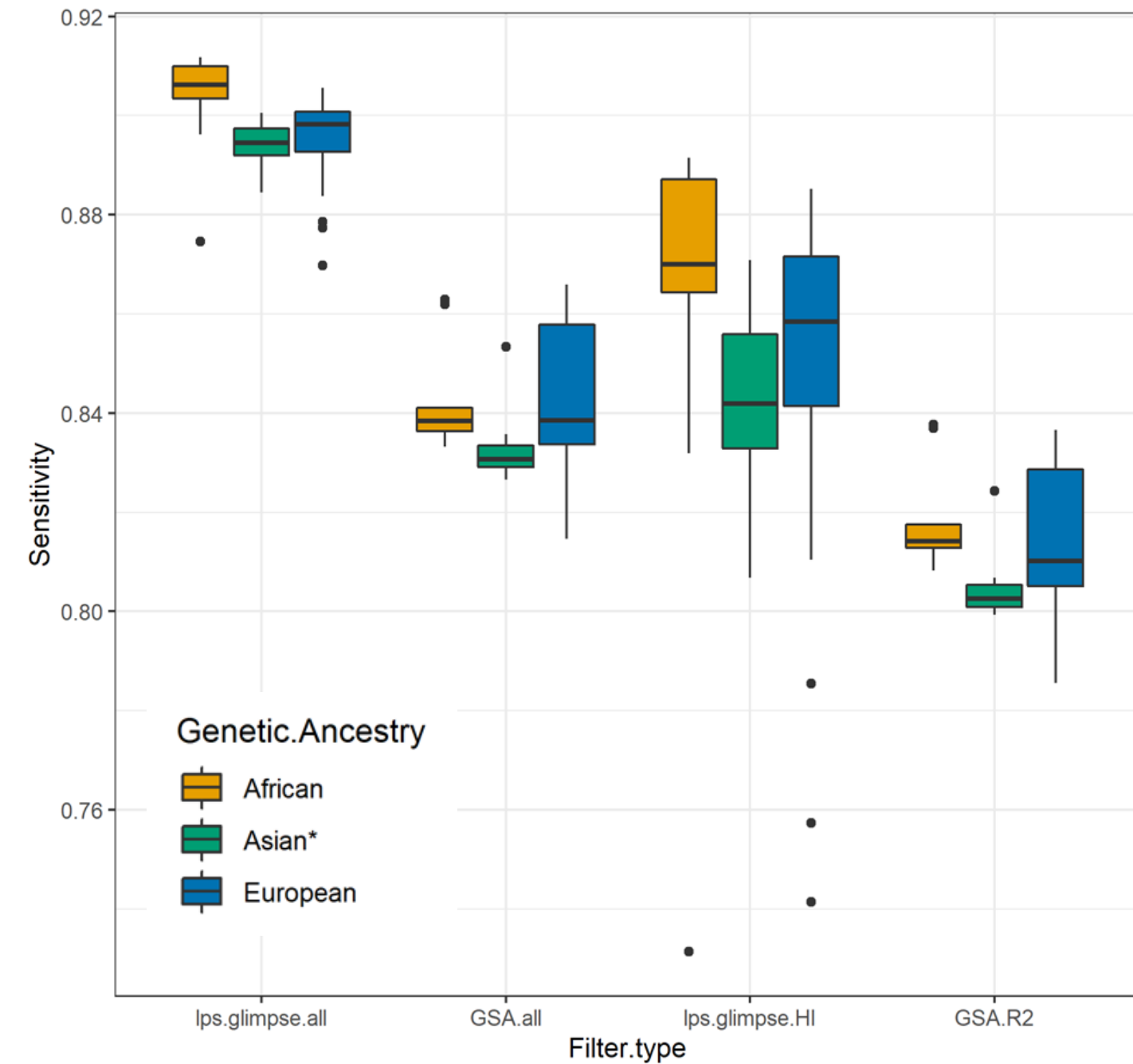
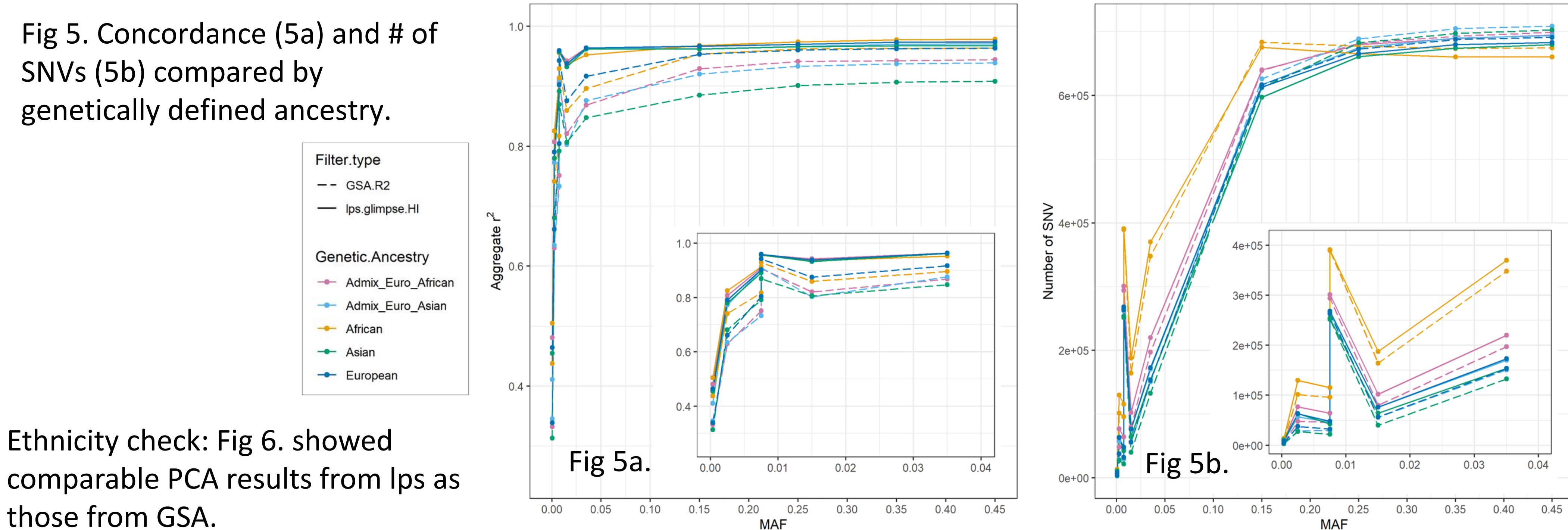


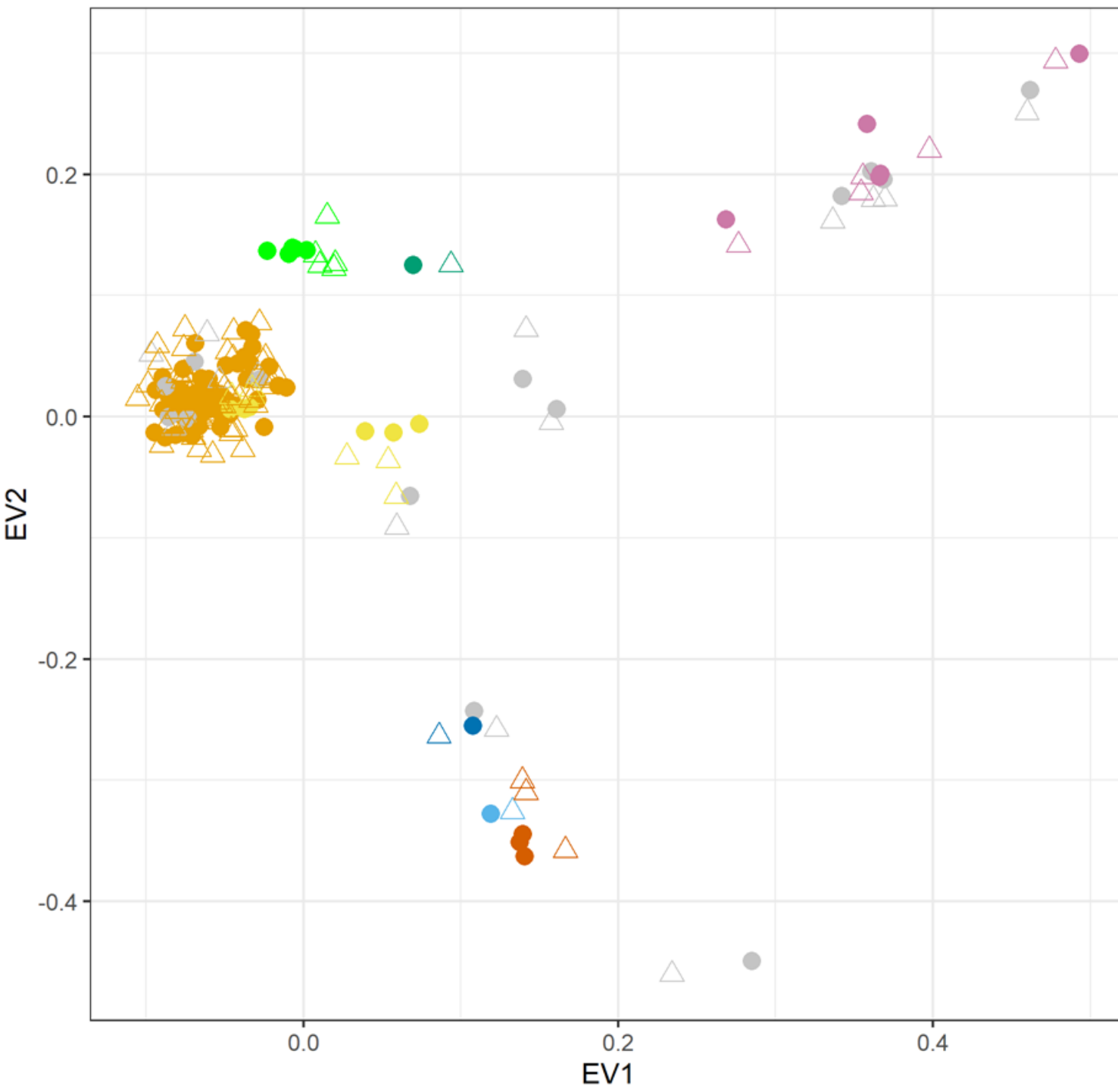
Fig 4.

Fig 5. Concordance (5a) and # of SNVs (5b) compared by genetically defined ancestry.

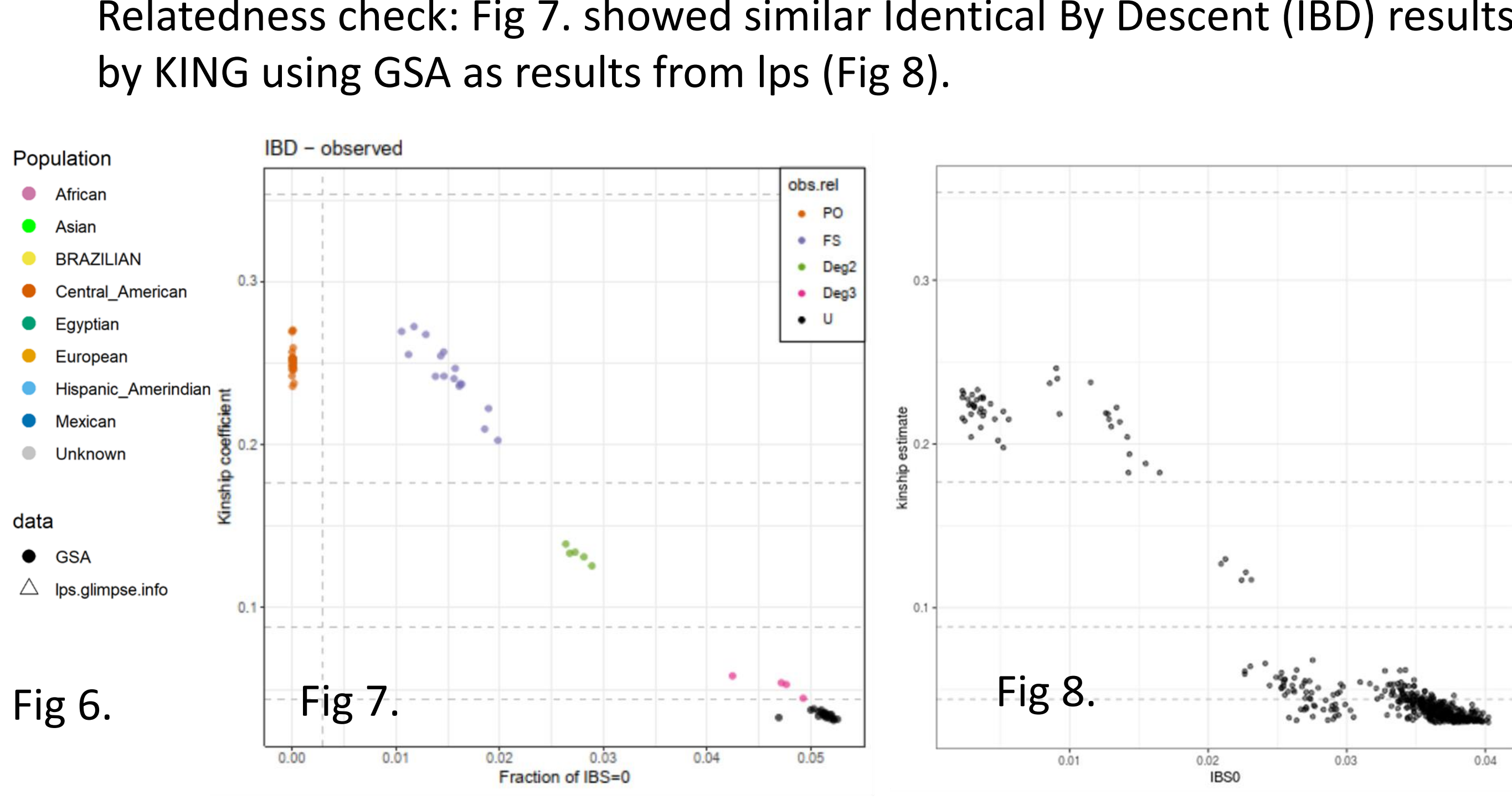


Ethnicity check: Fig 6. showed comparable PCA results from lps as those from GSA.

Population refers to self-reported race.



Relatedness check: Fig 7. showed similar Identical By Descent (IBD) results by KING using GSA as results from lps (Fig 8).



Summary & Discussions

- We demonstrated the feasibility of using lpWGS as a promising alternative for SNP array with currently available resources.
- The plexWell LP384 library prep method provides the dramatically increased throughput and reduced costs required to replace array for extremely large studies.
- Using results from deep coverage WGS (28X) as the validation dataset, we show that lpWGS (1X) can achieve higher sensitivity and concordance compared to GSA based imputation.
- The gain in sensitivity and concordance of lpWGS is more prominent for non-European ancestry or less common variants (e.g. MAF < 5%).
- Using lps data, we can detect sample mixtures, and perform analysis for ethnicity (PCA) and relatedness (IBD) with similar results as using array data.
- Future directions for lpWGS include comparing chr X, imputation of INDELs, and detections of chromosomal anomalies.