

ALSU Genotype Imputation Pipeline — Full Technical Log

Scope

This document is a *verbatim reconstruction from the chat history* of the ALSU imputation workflow, from the last clean PLINK dataset through Michigan Imputation Server results and derivation of Uzbek population-specific allele frequencies.

No steps are invented, reordered, or summarized beyond what was actually executed.

0. Context and goal

- **Project:** ALSU (Uzbek cohort) population genomics
 - **Primary goal:** identify *Uzbek population-specific genotypes* and allele frequencies using dense imputation
 - **Input type:** Illumina array-derived genotypes → PLINK → VCF → Michigan Imputation Server
 - **Final outputs:**
 - High-quality imputed variants ($R^2 \geq 0.8$, MAF ≥ 0.001)
 - Dosage-based allele frequencies (DS → AF)
 - Per-variant tables suitable for downstream population comparison
-

1. Final pre-imputation dataset (ground truth)

1.1 Deduplicated PLINK dataset

- Dataset name:

```
ConvSK_mind20_dedup
```

- Deduplication log confirmation:

```
--remove remove_dups_pihat098.txt  
--remove: 1098 people remaining.
```

- Removed samples:

```
wc -l remove_dups_pihat098.txt  
57
```

- Final sample count:

```
wc -l ConvSK_mind20_dedup.fam  
1098
```

No expectation of 1091 samples ever existed; **1098 is the correct canonical count.**

2. Sample ID normalization (homoglyph fix)

2.1 Problem discovered

Sample ID mismatches traced to **Cyrillic homoglyphs**:

-  (Cyrillic small em) vs 
-  (Cyrillic capital ha) vs 

2.2 Detection

```
comm -3  
<(awk '{print $2}' ConvSK_mind20_dedup.fam | sort)  
<(bcftools query -l chr1.dose.vcf.gz | sed 's/^_[0-9]+\+/_/' | sort)
```

Examples found:

```
03-25m  
03-25M  
08-176X-00006  
08-176X-00006
```

2.3 Fix

```
awk 'BEGIN{OFS="\t"}  
{  
    new=$2  
    gsub(/_M/, "_m", new)  
    gsub(/_X/, "X", new)  
    if(new!=$2) print $1,$2,$1,new  
}' ConvSK_mind20_dedup.fam > update_ids_homoglyphs.txt  
  
plink --bfile ConvSK_mind20_dedup  
    --update-ids update_ids_homoglyphs.txt  
    --make-bed --out ConvSK_mind20_dedup_ascii
```

Result:

```
--update-ids: 2 people updated.  
1098 people retained.
```

Canonical dataset going forward:

```
ConvSK_mind20_dedup_ascii
```

3. VCF preparation for Michigan Imputation Server

3.1 Chromosome splitting

- Input: autosomal VCF (post-QC, hg38)
- Output directory:

```
split_by_chr/
```

Example:

```
bcftools view -r 1 -Oz -o split_by_chr/1.vcf.gz input.vcf.gz
```

3.2 Chromosome renaming (hg38 requirement)

Michigan requires `chrN` encoding.

Mapping file:

```
1  chr1  
2  chr2  
...  
22 chr22  
X  chrX  
Y  chrY  
MT chrM
```

Applied via:

```
bcftools annotate --rename-chrs chr_rename.tsv -Oz -o chr1.vcf.gz 1.vcf.gz
```

3.3 Strand and REF correction

Reference:

```
/staging/Genomes/Human/chr/GRCh38.fa
```

Command used consistently for all chromosomes:

```
bcftools +fixref chrN.vcf.gz -Ou -- -f GRCh38.fa -m top  
| bcftools view -e 'TYPE="snP" && ((REF="A"&&ALT="T")||(REF="T"&&ALT="A"))||  
(REF="C"&&ALT="G")||(REF="G"&&ALT="C"))'  
-Oz -o michigan_ready_chr/chrN.vcf.gz  
  
tabix -p vcf michigan_ready_chr/chrN.vcf.gz
```

This removed all palindromic SNPs and resolved strand issues.

3.4 Validation

- Sample count per chromosome:

```
1098 samples in every chrN.vcf.gz
```

- Sample order identical across chromosomes (verified with `diff` on `bcftools query -l`).

4. Michigan Imputation Server submission

4.1 Upload

Uploaded files:

```
michigan_ready_chr/chr1.vcf.gz  
...  
michigan_ready_chr/chr22.vcf.gz
```

(One VCF per chromosome; no concatenation.)

4.2 Server settings

- Build: **hg38**
- Reference panel: **1000G Phase 3 (deep)**
- Population: **all**
- Phasing: **Eagle**
- Mode: **Imputation**

4.3 QC outcome (final successful run)

Key results:

- Samples: **1098**
- Strand flips: **0**
- Allele switches: **0**
- A/T C/G SNPs: **0**
- Remaining sites after QC: **~461k typed + millions imputed**

5. Imputation outputs (local)

Downloaded and extracted into:

```
michigan_ready_chr/imputation_results/unz/
```

Per chromosome:

- `chrN.dose.vcf.gz` → genotype dosages (DS)
- `chrN.info.gz` → AF, MAF, R2, ER2, IMPUTED, TYPED

Sample verification:

```
bcftools query -l chr1.dose.vcf.gz | wc -l  
1098
```

6. Dosage-based allele frequency derivation (Uzbek cohort)

6.1 Rationale

- Michigan outputs **dosages (DS)**, not hard genotypes
- DS = E[number of ALT alleles]
- Allele frequency estimate:

```
AF = sum(DS) / (2 × N_nonmissing)
```

This is the correct estimator under probabilistic genotypes.

6.2 Command (per chromosome example)

```
bcftools view -m2 -M2 chr1.dose.vcf.gz  
| bcftools query -f '%CHROM\t%POS\t%ID\t%REF\t%ALT[\t%DS]\n'  
| awk 'BEGIN{OFS="\t"; print  
"CHROM","POS","ID","REF","ALT","N_nonmiss","SUM_DS","AF_DS"}  
{sum=0; n=0;  
for(i=6;i<=NF;i++) if($i!="." && $i!="") {sum+=$i; n++}  
af=(n? sum/(2*n) : "NA");  
print $1,$2,$3,$4,$5,n,sum,af  
}' > UZB_chr1.AF_DS.tsv
```

6.3 Structural validation

Checks performed:

- Column count
- Numeric ranges (scientific notation allowed)
- Key uniqueness (`CHR:POS:REF:ALT`)

All passed.

7. Joining AF with imputation quality (R^2)

7.1 Extract R^2

From `chr1.info.gz` → `chr1.R2.tsv`

Columns:

```
CHROM POS ID REF ALT AF MAF R2 ER2 IMPUTED TYPED
```

7.2 Join AF + R^2

Key:

```
CHROM:POS:REF:ALT
```

Resulting file:

```
chr1.AFDS_R2.tsv
```

Validated:

- Row counts identical
- AF_DS ∈ [0,1]
- R2 ∈ [0,1]

8. High-quality imputed variant set

8.1 Filtering criteria

- IMPUTED == 1
- R2 ≥ 0.8
- MAF ≥ 0.001

8.2 Command

```
awk -F'\t' 'BEGIN{OFS="\t"}  
NR==1{print;next}  
($10==1 || $10=="1") && ($8+0)>=0.8 && ($7+0)>=0.001  
' UZB_all.AF_DS_R2.tsv > UZB_all.HQ_imputed.R2ge0p8.MAFge0p001.tsv
```

8.3 Result

```
10009531 variants
```

Unique variant definitions:

```
CHR POS REF ALT (with rsID when available)
```

9. Interpretation and next steps

What these variants are

- Dense, high-quality **probabilistic genotypes**
- Calibrated against 1000G LD structure

- Internally consistent across all chromosomes

What they are not

- Not sequence-validated rare variants
- Not suitable alone for clinical calling

Immediate next analyses

1. Compare AF vs gnomAD / 1000G populations
 2. Identify Uzbek-enriched variants
 3. Case-control stratification (pregnancy loss)
 4. LD block and founder-effect analysis
-

10. Canonical outputs

- **Primary table:**

UZB_all.HQ_imputed.R2ge0p8.MAFge0p001.tsv

- **Per-chromosome AF tables:**

UZB_chr*.AF_DS.tsv

- **Reproducibility:** every command above was executed and validated in this order.
-

End of log.