

MICHIGAN **IMPUTATIONSERVER**



Christian Fuchsberger



Sebastian Schönherr



Cassandra Spracklen



Lukas Forer



Eleftheria Zeggini



Albert Smith



Disclosure Slide

Financial Disclosure for:

Christian Fuchsberger
Sebastian Schönherr
Lukas Forer
Cassandra Spracklen
Eleftheria Zeggini
Albert Smith

We have nothing to disclose



Setup

- 6 Sessions:
 - (1) Genotype Imputation
 - (2) Use the server and the Imputation Bot,
 - (3) GWAS, chrX, HLA (4) GWAS pipeline and PGS Server,
 - (5) Helmholtz Munich Imputation Server, (6) TOPMed
 - Lectures
 - Demos
 - Interaction
 - PollEv.com/ashg

Question & Answer session at the end



Section 1 Imputation and the Server



Christian Fuchsberger Eurac Research cfuchsberger@eurac.edu





Learning objectives

Participants will

1. Understand the principles of genotype imputation and the Michigan Imputation Server



Genotype imputation

Key method used in GWAS to

- Increase the number of tested variants
- Fine-mapping becomes more complete
- Meta-analysis using different arrays



0. Imputation setting

GWAS Haplotypes

Reference Haplotypes (e.g. TOPMed)



1. Identify match among reference

GWAS Haplotypes

Reference Haplotypes (e.g. TOPMed)

```
C G A G A T C T C C T T C T T C T G T G C

C G A G A T C T C C G A C C T C A T G G

C C A A G C T C T T T T T C T C T G T G C

C G A A G C T C T T T T T C T C T G T G C

C G A A G C T C T C T T T T T C T T C T G T G C

C G A G A T C T C C C G A C C T T A T G C

T G G A G A T C T C C C G A C C T T A T G C

C G A G A T C T C C C G A C C T T A T G C

C G A G A C T C T T T T T C T T T G T G C

C G A G A C T C T C T T T T C T T T G T G C
```



2. Impute

GWAS Haplotypes

```
        c
        g
        a
        g
        A
        t
        c
        c
        c
        g
        A
        c
        c
        t
        c
        A
        t
        g
        g
```

Reference Haplotypes (e.g. TOPMed)

```
C G A G A T C T C C T T C T T C T G T G C

C G A G A T C T C C G A C C T C A T G G

C C A A G C T C T T T T T C T C T G T G C

C G A A G C T C T T T T T C T C T G T G C

C G A A G C T C T C T T T T T C T C T G T G C

C G A G A T C T C C G A C C T T A T G C

T G G G A T C T C C G A C C T C A T G G

C G A G A T C T C C C G A C C T T A T G G

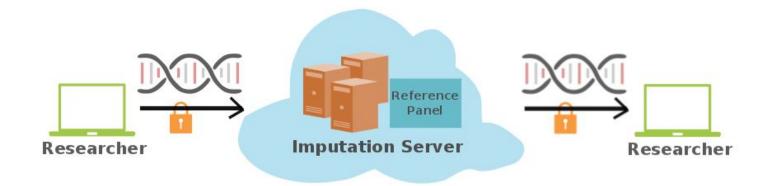
C G A G A C T C T C C C G A C C T T G T G C

C G A G A C T C T C T T T T C T T T G T G C

C G A G C T C T C T C C C G A C C T C G T G C
```



ASHG 2014: imputation web service



1.

Upload GWAS data

2.

Server performs

- Quality checks
- Pre-phasing
- Imputation
- Encryption

3.

Download results







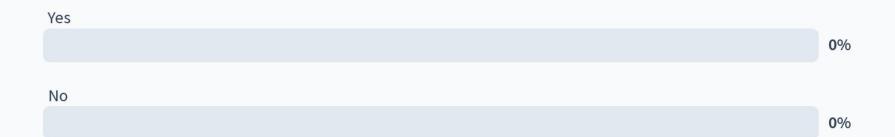
Interactive polls



https://pollev.com/ashg

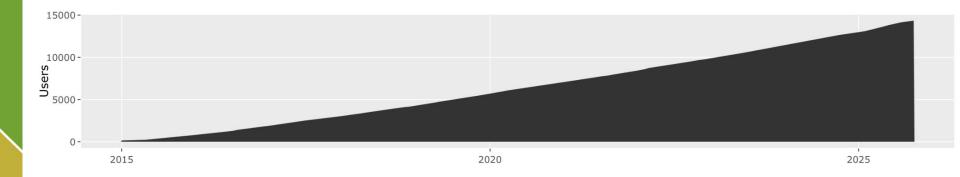


Have you ever used the Michigan Imputation Server



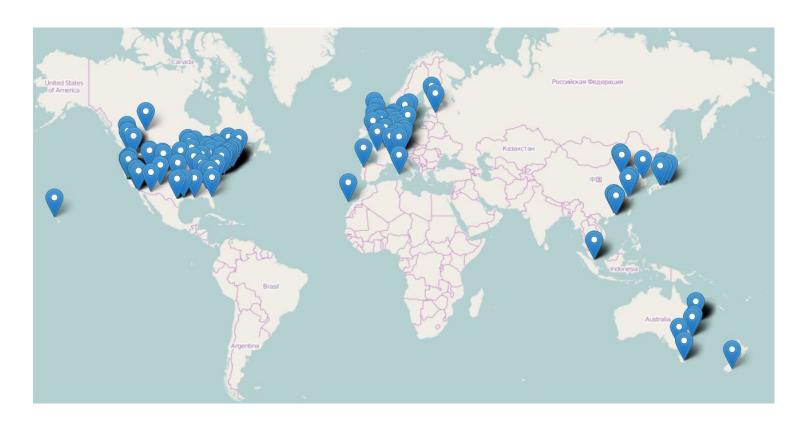


>14,000 users



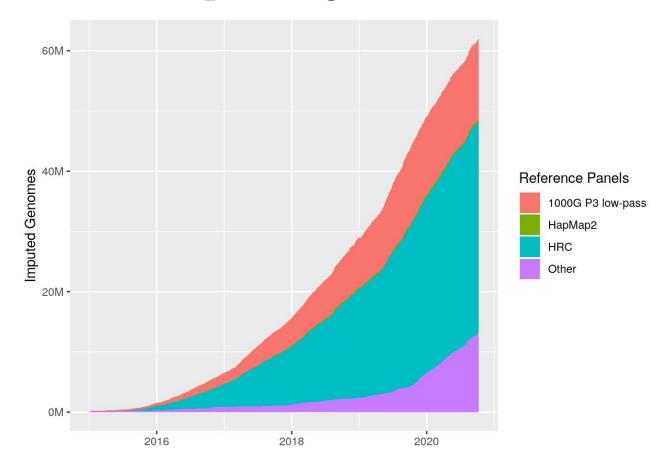


Used by researcher world-wide





>120M imputed genomes





Summary

Genotype imputation key method in GWAS

Michigan imputation server is easy to use and ensures high quality imputation

Cloud-services will accelerate genetic research so we can devote our time to more interesting tasks

More info and FAQ can be found here: https://genepi.github.io/michigan-imputationserver/



Section 2

Run a job, Data Preparation and Server Interaction MICHIGAN IMPUTATIONSERVER



Sebastian Schönherr Medical University of Innsbruck sebastian.schoenherr@i-med.ac.at



Learning objectives

Participants will learn

- 1. How to submit a genotype imputation job
- 2. How to prepare your input data
- 3. Different ways to interact with imputation servers



To run an imputation / PGS jobs, several imputation servers are available (1)

- Michigan Imputation Server (2015)
 https://imputationserver.sph.umich.edu
- TOPMed Imputation Server (2020)
 https://imputation.biodatacatalyst.nhlbi.nih.gov
- Helmholtz Munich Imputation Server (2024)
 https://imputationserver.helmholtz-munich.de
- Mexico City Prospective Study Imputation Server (2025)
 https://imputationserver-reg.sph.umich.edu/



To run an imputation / PGS jobs, several imputation servers are available (2)

- All servers are built on the same software stack
 - version 1: Hadoop based
 - version 2: Nextflow based (since 2024)
- Each server offers different reference panels and applications

 For this workshop, we provide a hands-on tutorial using the Michigan Imputation Server (MIS) instance



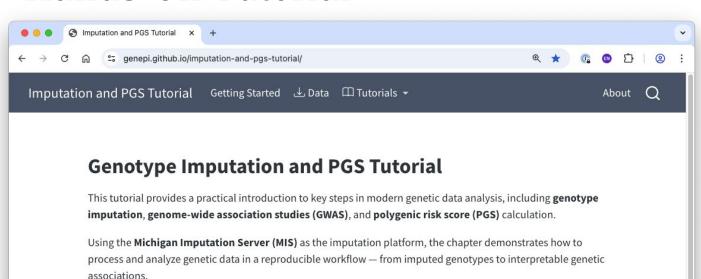
Which genotype imputation service have you utilized to date?

Michigan Imputation Server	
	0%
TOPMed Imputation Server	
	0%
Munich Imputation Server	
	0%
Other Imputation Service	
	0%



Hands-On Tutorial

You will learn how to:



https://tinyurl.com/imputationserver-2025

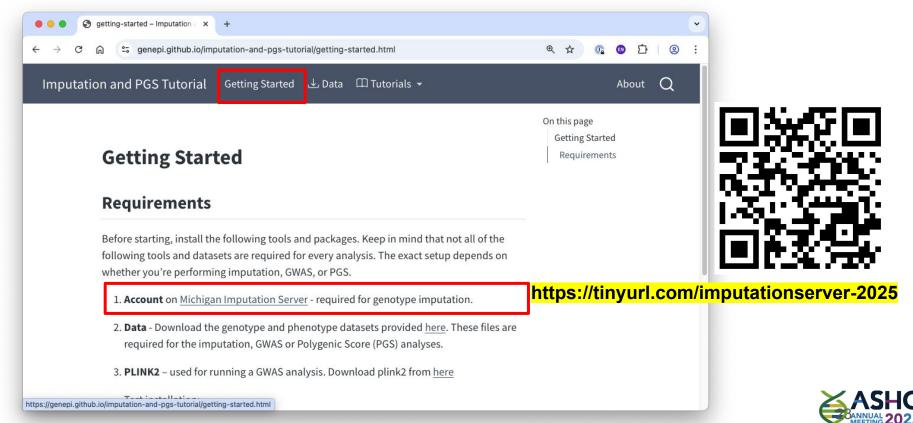
- · Perform and evaluate genotype imputation using the Michigan Imputation Server
- Conduct a GWAS to identify variants associated with phenotypes
- Calculate and visualize polygenic risk scores (PRS) for different trait types

The tutorial concludes with examples of four phenotypes, illustrating how PRS performance varies between binary and continuous traits — ranging from clear genetic signals to weak or absent associations.

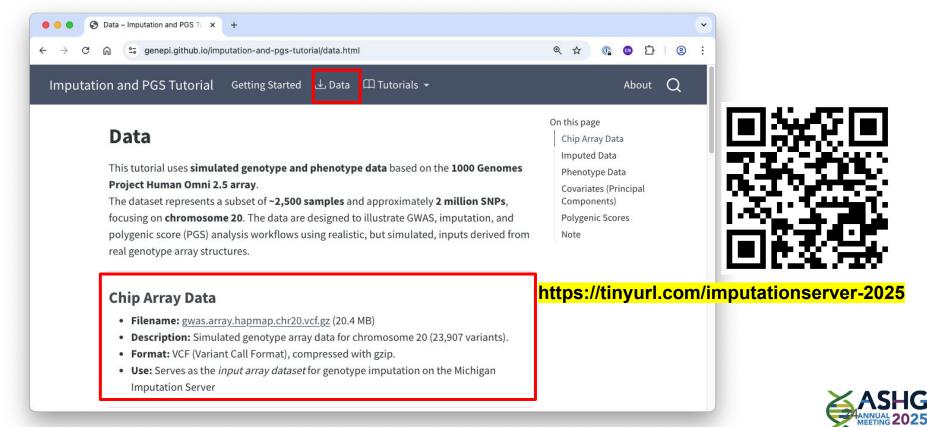




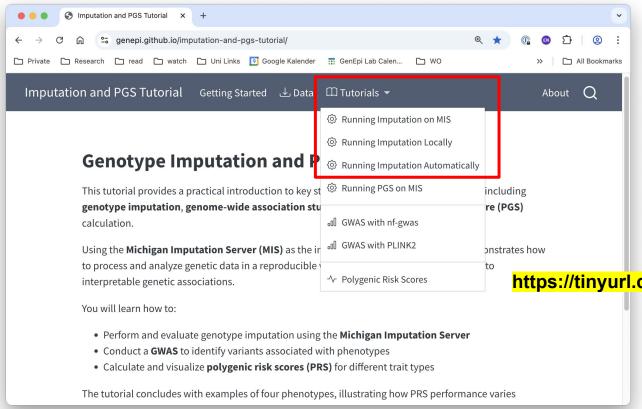
Running Imputation: Register an Account



Running Imputation: Download Array Data



Running Imputation on MIS





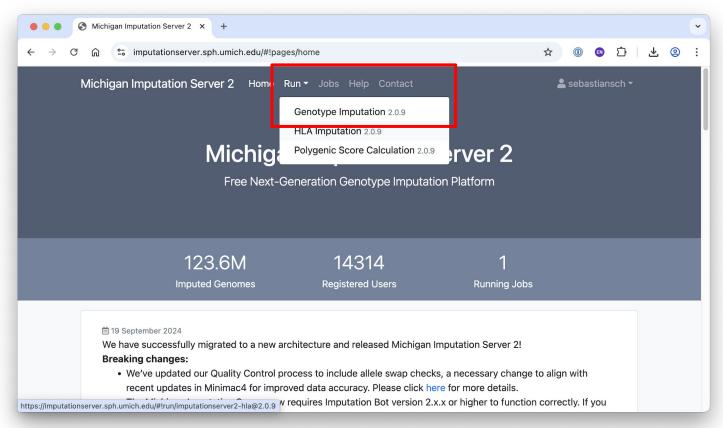
https://tinyurl.com/imputationserver-2025



Running Imputation on MIS

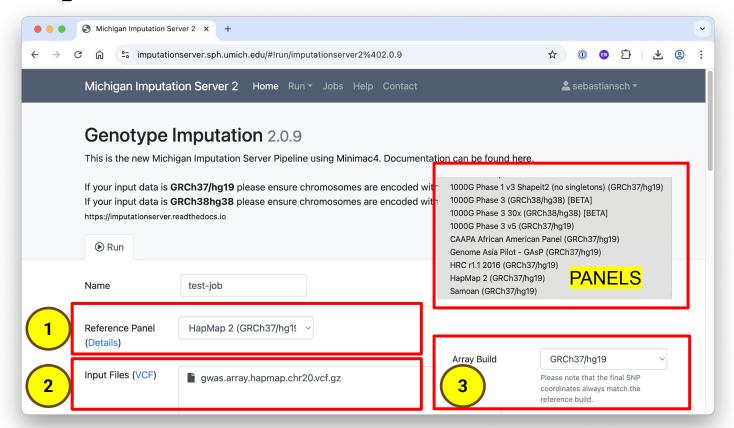


Running Imputation: Submit a first job



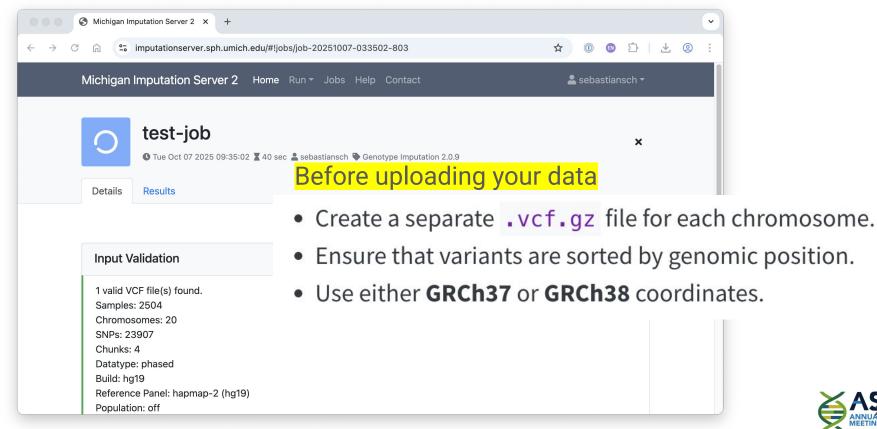


Upload Data and Select Reference Panel



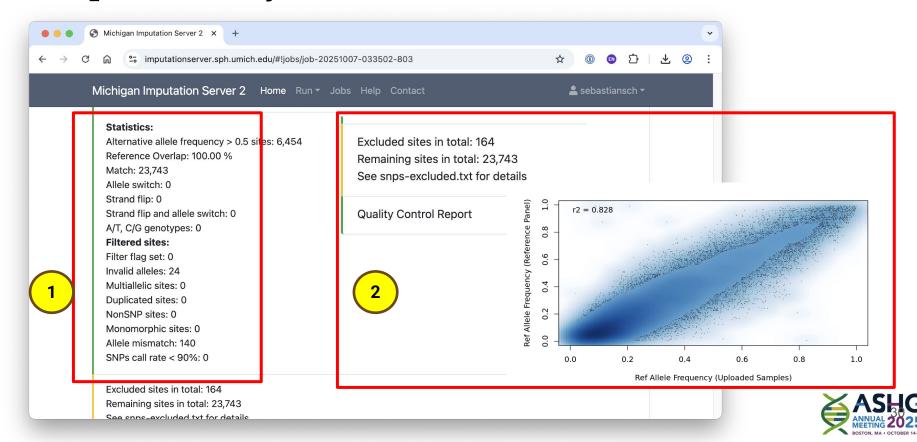


Step 1: Input Validation

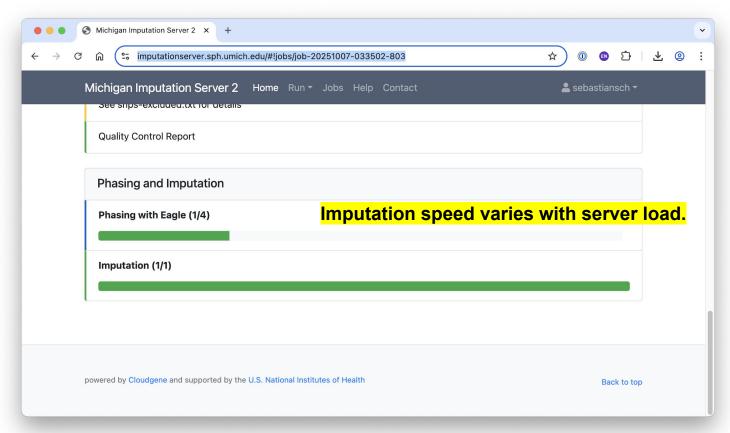




Step 2: Quality Control

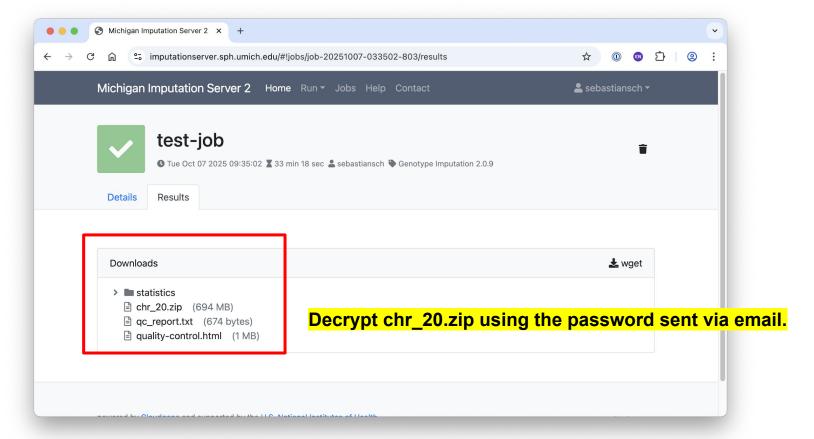


Step 3 + 4: Phasing, Imputation, Encryption





Download Results



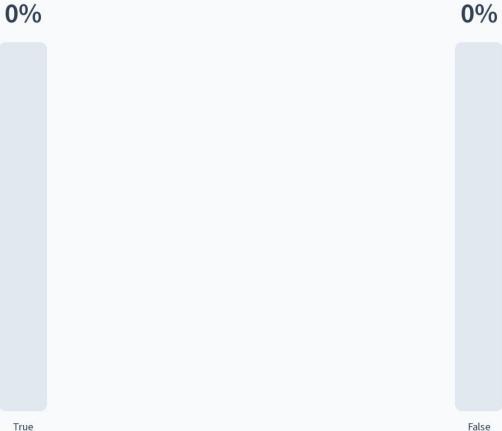


How many jobs are failing?

- 40% in 2015; 20% in 2019; 7% 2020-2024; >10% 2025
- Reason for job failures: Something wrong with your input data
 or phasing/imputation issue on our side



For users: Have you ever encountered a failed imputation job?





Input Validation

Input Validation

The provided VCF file is malformed. Error during index creation: [tabix] was bgzip used to compress this file? (see Help).

Input Validation

The provided VCF file contains more than one chromosome. Please split your input VCF file by chromosome (see Help).

Input Validation

Unable to parse header with error: Your input file has a malformed header: We never saw the required CHROM header line (starting with one #) for the input VCF file (see Help).

QC Step

Excluded sites in total: 695
Remaining sites in total: 185,791
See snps-excluded.txt for details

Typed only sites: 397

See typed-only.txt for details

Warning: 2 Chunk(s) excluded: reference o

Remaining chunk(s): 40

Training. 2 orialik(s) excluded. Tereferioe c

0.01--1/-) ---1-1-1--

Excluded sites in total: 11,088
Remaining sites in total: 1,968
See snps-excluded.txt for details
Typed only sites: 848,376

See typed-only tyt for details

Warning: 153 Chunk(s) excluded: reference overlap < 50.0% (see chunks-excluded.txt for details).

Remaining chunk(s). 0

Error: No chunks passed the QC step. Imputation cannot be started!

Error: More than 100 obvious strand flips have been detected. Please check strand. Imputation cannot be sta

Excluded sites in total: 165,335 Remaining sites in total: 161,526 See snps-excluded.txt for details

Warning: 2 Chunk(s) excluded: < 20 SNPs (see chunks-excluded.txt for details).

Remaining chunk(s): 152

Error: More than 100 allele switches have been detected. Imputation cannot be started!

How to fix input files?



Imputation Preparation Tool

TUTORIAL AVAILABLE ON MIS

- Developed by W. Rayner
- Works for all major reference panels (HRC, TOPMed, Asia, CAAPA, 1000G)
- Checks for consistency between input data and a reference panel
- Updates/removes SNPs, Updates strand, position and ref/alt assignment
- Input Data in PLINK Binary Format (bim, bed, fam)



Uploaded my data 5 mins ago and my job is still waiting...

There is a problem with MIS - email MIS team to let them know 0% There is a problem with my data 0% MIS is very busy - check again later 98% Don't know 2%



Error: No chunks passed the QC step. Imputation cannot be started!

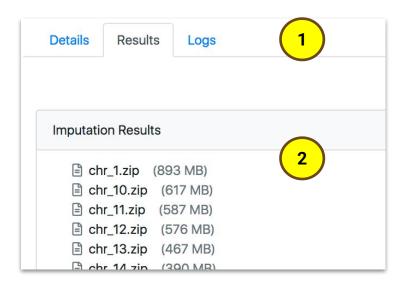
Data are on the wrong build Too few variants overlap with the reference panel 52% Must be a MIS problem, since imputation runs locally 2%	Email MIS team	
Too few variants overlap with the reference panel 52% Must be a MIS problem, since imputation runs locally		4%
Too few variants overlap with the reference panel 52% Must be a MIS problem, since imputation runs locally		
Too few variants overlap with the reference panel 52% Must be a MIS problem, since imputation runs locally	Data are on the wrong build	
Must be a MIS problem, since imputation runs locally		35%
Must be a MIS problem, since imputation runs locally		
Must be a MIS problem, since imputation runs locally	Too few variants overlap with the reference panel	
		52 %
2%	Must be a MIS problem, since imputation runs locally	
		2%
Don't know	Don't know	
8%		8%

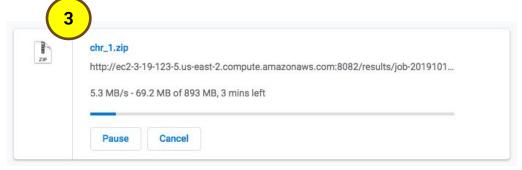


How to download the imputed genotypes?



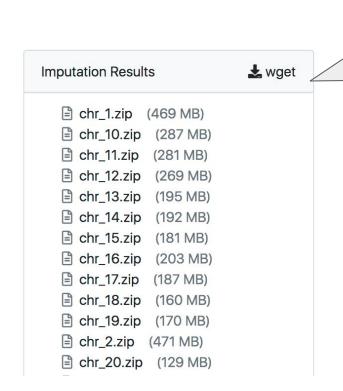
Option 1: Web-Interface

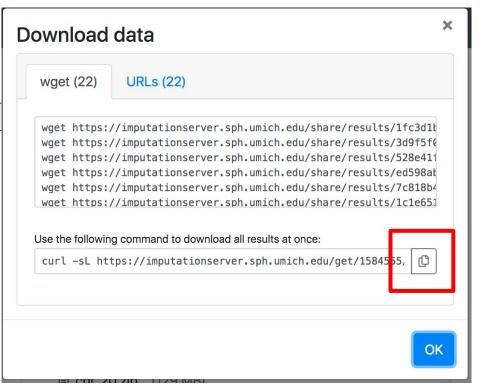






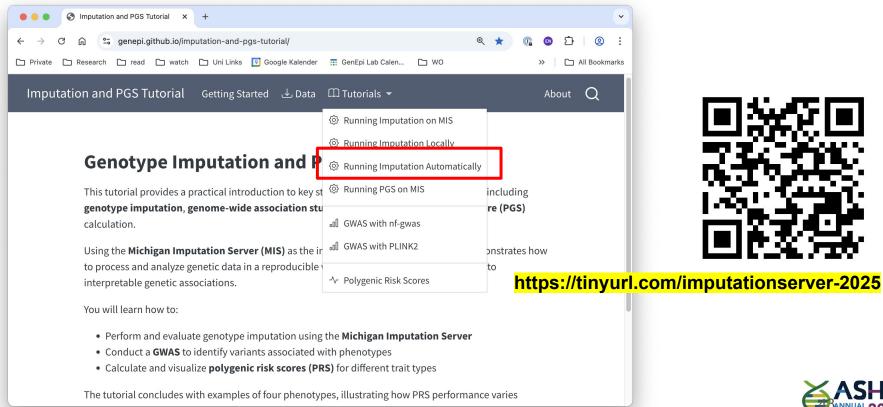
Option 2: Batch Download







Option 3: Use Imputationbot

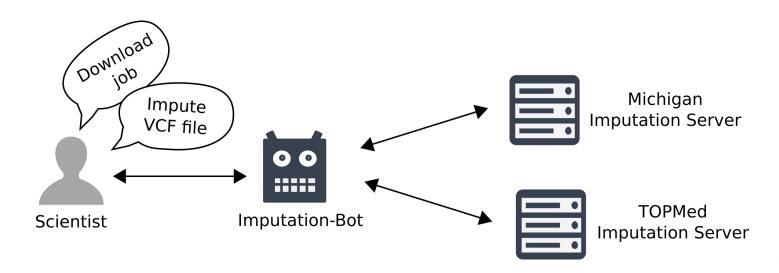




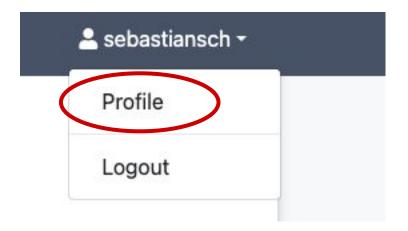


Imputation Bot

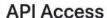
- Automate remote imputation
- Submit and monitor jobs from the command line
- Different commands can easily be combined











This service provides a rich RestAPI to submit, monitor and download jobs.

You need a access token to use the API. Learn more.



API Token

Make sure to copy your personal access token now. You won't be able to see it again:

eyJhbGciOiJIUzl1NiJ9.eyJzdWliOiJzZWJhc3RpYW5zY2giLCJ uYmYiOjE3NTk4NDczMTgslm1haWwiOiJzZWluc2Nob2VuaG VyckBnbWFpbC5jb20iLCJhcGlfaGFzaCl6lm44SXFHQ2lwN1l QamlwS0pkSjJManU5eVRwazhxZSlslnJvbGVzljpbXSwiaXNzl



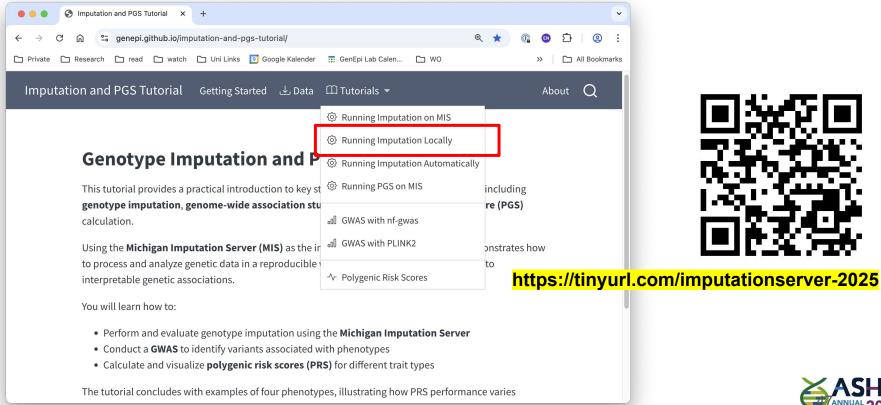
×



Workflow

- curl -sL imputationbot.now.sh/v2 | bash
- ./imputationbot add-instance
- ./imputationbot impute --files gwas.array.hapmap.chr20.vcf.gz --refpanel hapmap-2 --population eur
- ./imputationbot download job-20251007-104433-914 --password ABCD



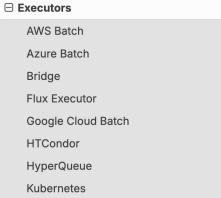








- Since 2024 MIS Workflow is based on Nextflow
- Nextflow is a workflow management system for scalable and reproducible bioinformatics pipelines
- Requirements
 - Install Nextflow
 - Use your own/publicly available reference panel
 - Run Imputation on the commandline







```
# Create temp directory
mkdir imputationserver
cd imputationserver
# Download Input Data
wget https://genepi.i-med.ac.at/downloads/imputation/gwas.array.hapmap.chr20.vcf.gz
# Download and unzip Reference Panel

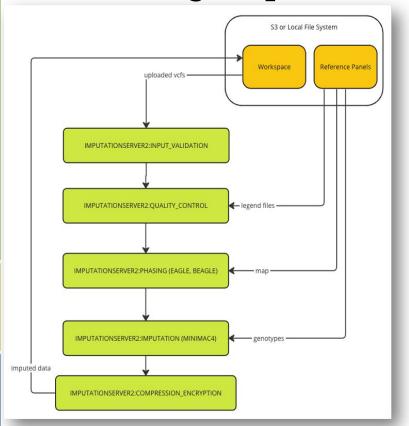
wget https://imputationserver.sph.umich.edu/resources/ref-panels/imputationserver2-lunzip imputationserver2-hapmap2.zip -d hapmap2
```

```
params {
                               = "my-test-project"
   project
    build
                                = "ha19"
                               = "gwas.array.hapmap.chr20.vcf.gz"
   files
   allele frequency population = "eur"
                               = "imputation"
    mode
    refpanel yaml
                               = "hapmap2/imputation-hapmap2.yaml"
                               = "results"
   output
   encryption.enabled
                                = false
```

3

nextflow run genepi/imputationserver2 -r v2.0.9 -c imputation.config -profile docker







```
NEXTFLOW ~ version 24.04.2

Launching `main.nf` [fabulous_mendel] DSL2 - revision: 16c60661b1

Loading reference panel from file /Users/seb/repositories/nf-imputationserver/tests/data/Welcome to Imputation Server 2 (v2.0.6)
executor > local (5)
[0a/d99e6b] INPUT_VALIDATION: INPUT_VALIDATION_VCF [100%] 1 of 1 ×
[91/a9fe51] QUALITY_CONTROL: QUALITY_CONTROL_VCF [100%] 1 of 1 ×
[cf/093d23] QUALITY_CONTROL: QUALITY_CONTROL_REPORT [100%] 1 of 1 ×
[6c/fd767c] PHASING: EAGLE (chunk_20_0040000001_00600000000.vcf.gz) [25%] 1 of 4
[- ] IMPUTATION: MINIMAC4 [0%] 0 of 1
[- ] ENCRYPTION: COMPRESSION_ENCRYPTION_VCF -
```



Summary

- MIS Web Interface provides a fast and reliable way to impute data
- We provide tutorials to run imputation on MIS, locally and via the imputationbot
- MIS applies a strict Quality Control with the goal to return high quality imputation data



Section 3 Performing GWAS using imputed data



Cassandra Spracklen
University of Massachusetts
cspracklen@umass.edu



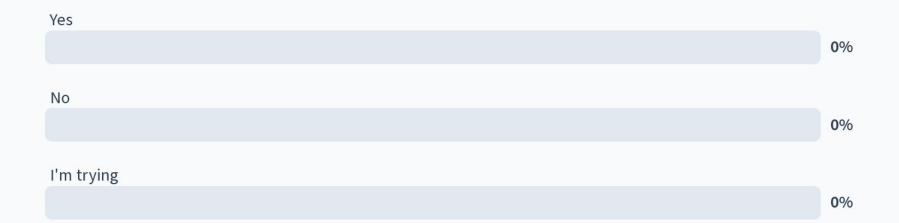
Learning objectives

Participants will learn to:

- Identify and understand the use of variant imputation quality information following imputation in the MIS
- Distinguish between some of the available options for GWAS
- Troubleshoot common GWAS errors



Have you ever performed a GWAS?





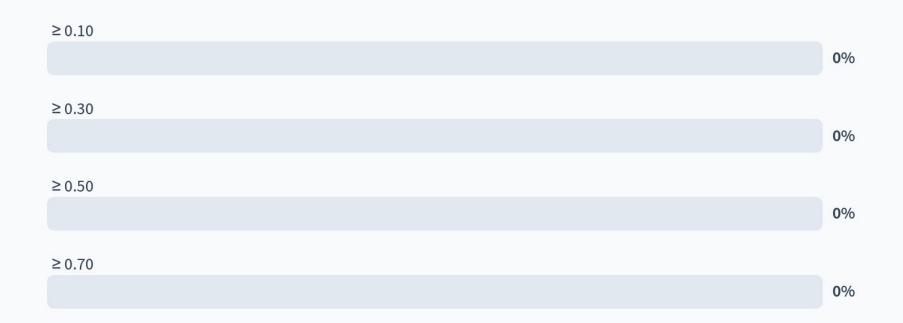
Imputation Quality

- For each variant, how confident can we be that the imputation dosages are sufficiently "accurate" for association analyses?
- Measure of confidence in imputed dosages: "Rsq" column [range 0-1]

```
SNP REF(0) ALT(1) ALT_Frq MAF AvgCall Rsq Genotyped ...
20:61795:G:T G T 0.26318 0.26318 0.88455 0.54658 Imputed ...
20:63231:T:G T G 0.03843 0.03843 0.98342 0.67736 Imputed ...
20:63244:A:C A C 0.16132 0.16132 0.91761 0.49907 Imputed ...
```

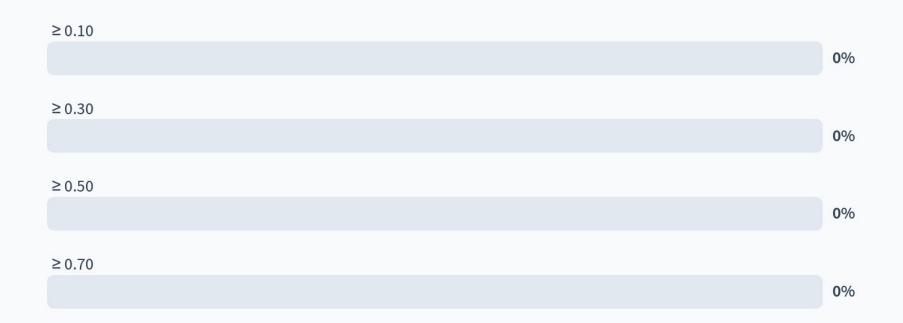


Minimally accepted Rsq value for common (MAF≥5%) variants?





Minimally accepted Rsq value for low frequency (MAF<5%) and rare variants?





Imputation Quality

- Minimal Rsq value for common variants
 - ≥ 0.30
- Minimal Rsq value for low frequency/rare variants
 - ≥0.50
- Before performing GWAS, remove variants that do not meet these thresholds
 - Suggested program: VCFtools
 - Saves computational time when performing GWAS



Which GWAS program(s) have you used?



Performing the GWAS

- Each program has its own input, output formats, and options
- Typical input files
 - Genotype file (.vcf; .bgen; .bed/.bim/.fam)
 - Phenotype/covariate file (.txt; .ped)
 - Some programs use separate phenotype and covariate files
 - Kinship/relationship matrix (EPACTS, SAIGE)



Available GWAS Programs

No File Reformatting (VCF from MIS)

- EPACTS
- Rvtests
- SNPTEST
- SAIGE

File Formatting Required

- BOLT-LMM
- BGENIE
- regenie
- PLINK



EPACTS/Rytests

- + Many model options single variant, gene-based
- + Chr X analyses
- + Phenotypic transformation (e.g inverse normal; Rvtests only)
- + Linear mixed model for sample relatedness (quantitative traits only)
- + Generate covariance matrices for downstream analyses (e.g conditional analyses; Rvtests only)
- Memory intensive
- Sample size ~≤20,000 (better ≤10,000)

EPACTS: https://genome.sph.umich.edu/wiki/EPACTS
Rvtests: https://genome.sph.umich.edu/wiki/Rvtests



SAIGE

- + Similar to Rvtests, but for very large sample sizes (e.g. biobanks)
- + Able to account for sample relatedness for binary traits
- + Designed to handle unbalanced number of cases and controls
- + Chr X analyses
- Should not be used to examine heritability (biased variance estimates)
- Computational time can vary widely between phenotypes and sample sizes
- Can be conservative for extremely unbalanced case and control ratio
- Odds ratios estimated to conserve computational time



SNPTEST

- + Frequentist and bayesian methods supported
- + Chr X analyses
- Limited to unrelated individuals
- Computationally intensive



BOLT-LMM/BGENIE/Regenie

- + Great for very large sample sizes (e.g. biobanks)
- + Chr X analyses
- + Computationally efficient (Regenie)
- Requires files to be in BGEN or PLINK format
- Nextflow pipeline for regenie using VCF: https://github.com/genepi/nf-gwas
- Not optimal for extremely unbalanced case control ratio (especially with rare variants)

BOLT-LMM: https://data.broadinstitute.org/alkesgroup/BOLT-LMM/#x1-5600011

BGENIE: https://jmarchini.org/bgenie/

Regenie: https://github.com/rgcgithub/regenie



PLINK

- + Quick
- + Multiple versions; often as intermediary tool to the other programs
- + Can run on the command line (unix not required)
- + Chr X analyses
- Requires files to be in PLINK format (.bed/.bim/.fam)
- Limited model options



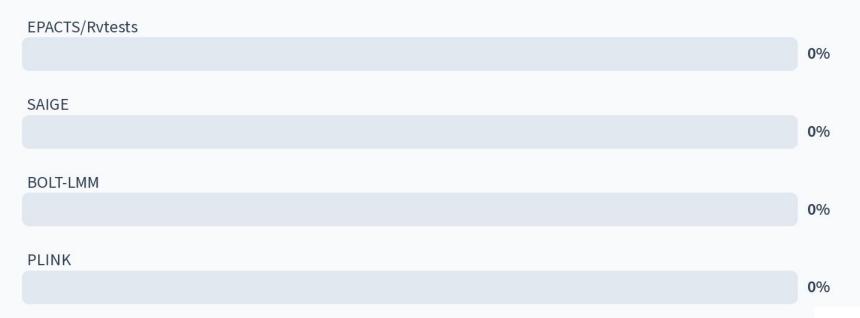
Summary of common GWAS analysis tools

	EPACTS	Rvtests	SNPTEST	SAIGE	BOLT-LMM	Bgenie	Regenie
Input VCF	Y	Y	Y				
Sample relatedness (Quantitative outcome)	Y	Y		Υ	Y	Υ	Υ
Sample relatedness (Binary outcome)				Υ		Υ	Υ
Case control imbalance				Υ			Υ
Large sample size (>20,000)				Υ	Y	Υ	Υ



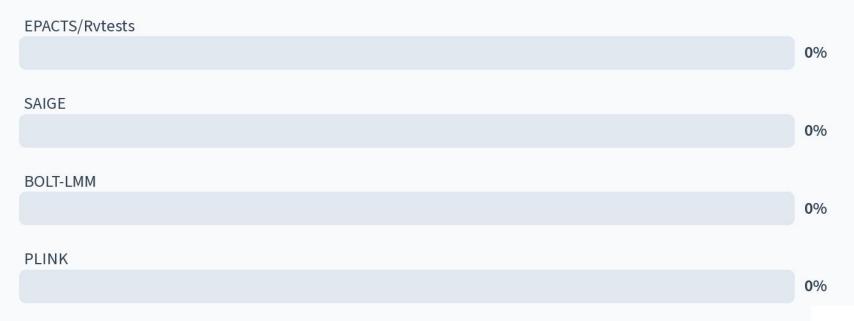
Which program(s) would be best?

A researcher new to genetic analyses and unfamiliar to the UNIX environment wants to perform a GWAS on total cholesterol using a cohort of 5,500 unrelated individuals.





Which program(s) would be best? Researchers want to perform a GWAS using a cohort of 10,000 individuals with household based recruitment (i.e. includes related individuals).





Which program(s) would be best? Researchers want to perform a GWAS using data from BioBank Japan (>200,000 individuals)

(A) EPACTS/Rvtests	
	0%
(B) SAIGE	
	0%
(C) BOLT-LMM	
	0%
(D) PLINK	
	0%
(E) Regenie/Bgenie	
	0%



Common Errors When Running a GWAS

- Wording of error messages vary by program, but the same issues will cause errors throughout all of the program
- [Unix] Errors independent of GWAS program
 - File permissions
 - Correct by changing file permissions
 - Directory/file not found
 - Correct by making sure all of the file locations and names are accurate
 - Not enough memory/time
 - Correct by restarting job with adequate memory/time allocation



Common Errors When Running a GWAS

- Common errors
 - o IDs don't match
 - Correct by ensuring that the ID in the phenotype, genotype, covariance, kinship matrix are consistent format in all files



- Common errors
 - o IDs don't match
 - Correct by ensuring that the ID in the phenotype, genotype, covariance, kinship matrix are consistent format in all files
 - File format(s) incorrect
 - Correct by making sure the format of all files are as the program is expecting (e.g. columns, delimiters, headers, file extension



Common errors

- IDs don't match
 - Correct by ensuring that the ID in the phenotype, genotype, covariance, kinship matrix are consistent format in all files
- File format(s) incorrect
 - Correct by making sure the format of all files are as the program is expecting (e.g. columns, delimiters, headers, file extension
- Improperly specified options/command
 - Correct by checking all needed options are specified, correct order, no typos



Common errors

- o IDs don't match
 - Correct by ensuring that the ID in the phenotype, genotype, covariance, kinship matrix are consistent format in all files
- File format(s) incorrect
 - Correct by making sure the format of all files are as the program is expecting (e.g. columns, delimiters, headers, file extension
- Improperly specified options/command
 - Correct by checking all needed options are specified, correct order, no typos
- Peripheral programs not available (e.g. R with EPACTS, SAIGE)
 - Correct by installing other peripheral programs



Common errors

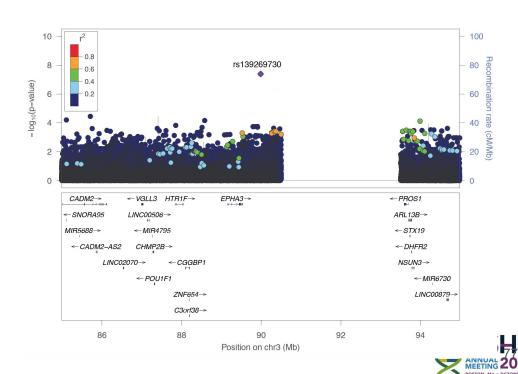
- o IDs don't match
 - Correct by ensuring that the ID in the phenotype, genotype, covariance, kinship matrix are consistent format in all files
- File format(s) incorrect
 - Correct by making sure the format of all files are as the program is expecting (e.g. columns, delimiters, headers, file extension
- Improperly specified options/command
 - Correct by checking all needed options are specified, correct order, no typos
- Peripheral programs not available (e.g. R with EPACTS, SAIGE)
 - Correct by installing other peripheral programs
- Invalid estimate (e.g. heritability in BOLT-LMM)
 - Sample too related and/or sample size too small
 - Correct by using a different program



Interpreting GWAS Results

- GWAS results must be carefully reviewed for:
 - Imputation quality!
 - Genomic inflation
 - False positives

- Replication datasets
- PheWas



Sex chromosomes

- Humans have two sex chromosomes, X and Y, that in combination determine the sex of an individual.
 - To ensure balanced expression, one of the female's X-chromosome is randomly selected to undergo inactivation(either paternal or maternal determined at early embryonic development).
 - Pseudoautosomal regions (PAR1 and PAR2) are short regions on the X- and Y-chromosomes that are inherited in an autosomal rather than a sex-linked manner



X-chromosome - genotype calls and QC

- Genotype calls on X-chromosomes
 - One copy in males, two copies in females
- Clarity of genotype calls on X-Chromosome
 - 0/1/2 allele coding in females; 0/1 or 0/2 allele coding in males
 - o 0/0.5/1 allele coding in females; 0/0.5 allele coding in males
 - Comparing standard error with autosomal data
- Analysis plan to be specific on allelic coding
- Develop QC checks before association analyses, e.g., similar allele frequencies between males and females (differential calls bias)

X-chromosome - imputation

Chromosome X Pipeline

Additionally to the standard QC, the following per-sample checks are executed for chrX:

- Ploidy Check: Verifies if all variants in the nonPAR region are either haploid or diploid.
- Mixed Genotypes Check: Verifies if the amount of mixed genotypes (e.g. 1/.) is < 10 %.

For phasing and imputation, chrX is split into three independent chunks (PAR1, nonPAR, PAR2). These splits are then automatically merged by Michigan Imputation Server and are returned as one complete chromosome X file. Only Eagle is supported.



Interpretation of association

- Biological considerations:
 - Assumption of total X-inactivation but there is evidence of escape and compensation mechanisms
 - Different inheritance/activation in different tissues?
 - Sex-biased gene expression? Unclear if it is real differences in gene expression or influence of hormones.
- Methodology development
 - Statistical methods to estimate the inactivation when estimating the effects; adjusting for sex



Summary

- Variants must be filtered post-imputation to remove those with poor imputation quality
- There are many GWAS programs available, each with their own strengths and limitations - so be sure to pick one that fits your analyses needs
- As these GWAS programs are widely used or adopted by consortia, there are tutorials and help-pages available
- X-chr imputation and analyses are also possible!

More info and FAQ can be found here: https://imputationserver.readthedocs.io



HLA Imputation

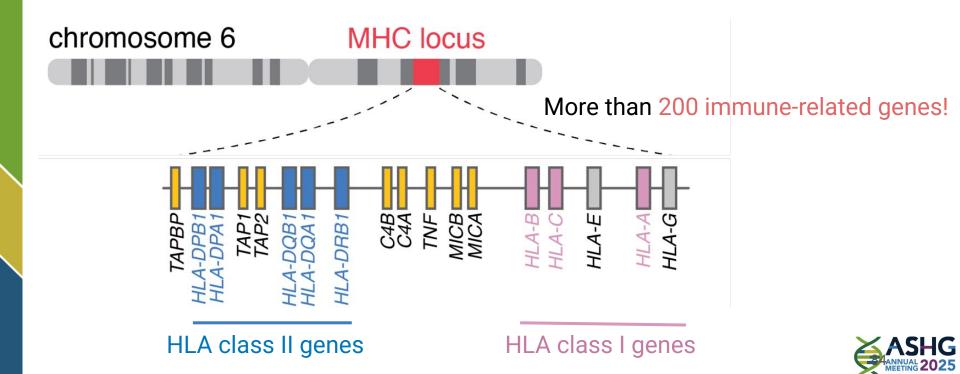
Slides Created By:

Saori Sakaue Broad Institute ssakaue@broadinstitute.org

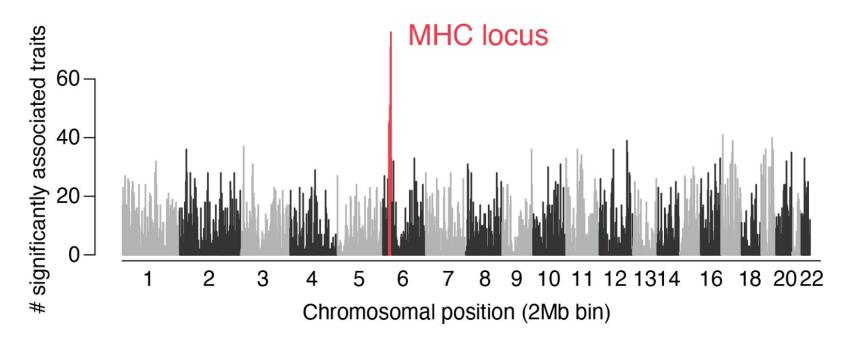




HLA genes are within MHC (major histocompatibility complex) locus on chromosome 6

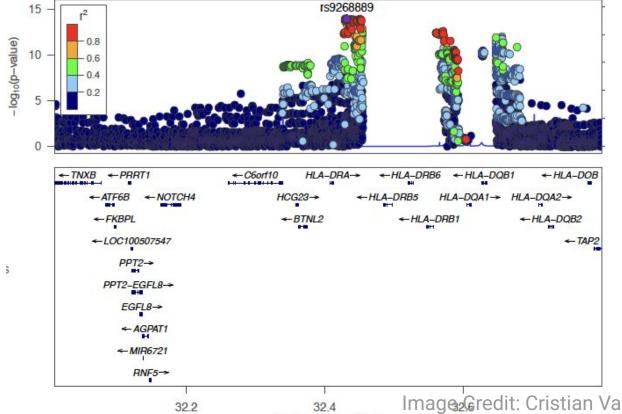


MHC locus confers the largest number of associations of any locus genome-wide





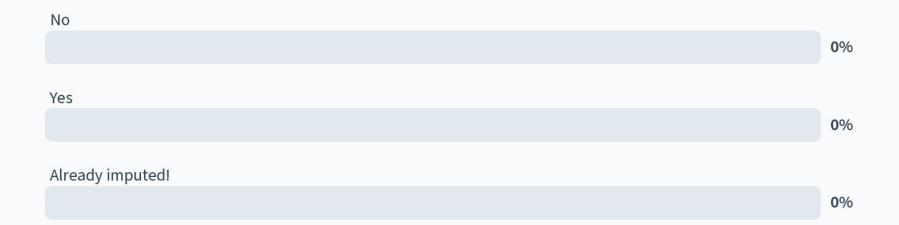
Genotype Imputation of HLA Region is great, but not perfect



Position on chr6 (Mb)



Are you planning to impute the HLA region?





Genotype imputation in a nutshell

Genotype imputation

Known Genotyped SNPs

Given Reference haplotype with

whole-genome SNPs

Unknown Untyped SNPs (Impute!)



Unknown Untyped SNPs

(Impute!)

	Genotype imputation	HLA imputation
Known	Genotyped SNPs	Genotyped SNPs
Given	Reference haplotype with whole-genome SNPs	Reference haplotype with HLA alleles and amino acid sequences

Untyped HLA alleles and

amino acid sequences

Known SNP genotype of the target cohort										
Genotype in the MHC region	HLA alleles	HLA amino acids	HLA intragenic SNPs							
¶ CGA.ATCTGTCTTCTGT.CTAA	[?]	?	?							
TCAA.ATCTGTCCT.TGT.CTAA	?	?	?							
<pre> ¶CAA.ATTTTGCTTCAGT.CTAA</pre>	?	?	?							

QCed Plink *.{bed,bim,fam}



```
Genotype in the MHC region

CGA.ATCT..GTCTTCTGT.CTAA

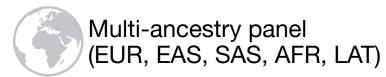
CAA.ATCT..GTCCT.TGT.CTAA

CAA.ATCT..TGCTTCAGT.CTAA

CAA.ATTT..TGCTTCAGT.CTAA
```

Given HLA imputation reference panel

```
Scaffold genotype in the MHC HLA alleles HLA amino acids HLA intragenic SNPs CGAGATCTCAGTCTTCTAA ••• DRB1*04:01••• GGSCMAALTVTLMVL••• GGAGGACCTGTGAACCA CAAGATTTCCTTCATCTGTTCTAA ••• DRB1*01:01••• GGSCMTALTVTLMVL••• GGAAGACCTGCGAACCA CGAGATCTCCTGCTTCAGTTCTAA ••• DRB1*01:02••• GGSCMTALTVTLMVL••• GGAAGACCTGCGAACCA CAAGATCTCCGTCCTCTGTTCTAA ••• DRB1*15:01••• GGSCMTALTVTLMVL••• GGAAGACCTGTGAACCG
```



Luo et al. Nat Genet 2021, Sakaue et al. Nat Protocols 2023



```
Known SNP genotype of the target cohort
 Genotype in the MHC region
  CGA.ATCT..GTCTTCTGT.CTAA
  CAA.ATCT..GTCCT.TGT.CTAA
  CAA.ATTT..TGCTTCAGT.CTAA
Given | HLA imputation reference panel
                                         HLA amino acids
 Scaffold genotype in the MHC
                             HLA alleles
                                                             HLA intragenic SNPs
 CGAGATCTCAGTCTTCTAA ← DRB1*04:01 ← GGSCMAALTVTLMVL ← GGAGGACCTGTGAACCA
                     CTAA •—• DRB1*01:01•—•GGSCMTALTVTLMVL•—• GGAAGACCTGCGAACCA
 CGAGATCTCCTGCTTCAGTTCTAA ◆ → DRB1*01:02 ◆ → GGSCMTALTVTLMVL ◆ → GGAGGACCTGCGAACCA
 CAAGATCTCCGTCCTCTGTTCTAA ◆ → DRB1*15:01 ◆ → GGSCMTALTVTLMVL ◆ → GGAAGACCTGTGAACCG

    ★ Haplotype phasing + Imputation
      Estimation of HLA imputation for the target cohort
¶CGAGATCTCAGTCTTCTGTTCTAA • → DRB1*04:01 • → GGSCMAALTVTLMVL • → GGAGGACCTGTGAACCA
             CCTCTGTTCTAA → DRB1*15:01 → GGSCMTALTVTLMVL → GGAAGACCTGTGAACCG
                TCAGTTCTAA •—● DRB1*01:01•—●GGSCMTALTVTLMVL •—● GGAAGACCTGCGAACCA
```



```
Known SNP genotype of the target cohort
 Genotype in the MHC region
  CGA.ATCT..GTCTTCTGT.CTAA
  CAA.ATCT..GTCCT.TGT.CTAA
  CAA.ATTT..TGCTTCAGT.CTAA
Given | HLA imputation reference panel
 Scaffold genotype in the MHC
                             HLA alleles
                                           HLA amino acids
                                                              HLA intragenic SNPs
 CGAGATCTCAGTCTTCTGTTCTAA ● → DRB1*04:01 ● → GGSCMAALTVTLMVL ● → GGAGGACCTGTGAACCA
                     CTAA •—• DRB1*01:01•—•GGSCMTALTVTLMVL •—• GGAAGACCTGCGAACCA
 CGAGATCTCCTGCTTCAGTTCTAA ◆ → DRB1*01:02 ◆ → GGSCMTAI TVTI MVI ◆ → GGAGGACCTGCGAACCA
 CAAGATCTCCGTCCTCTGTTCTAA ◆ → DRB1*15:01 ◆ → GGSCMTALTVTLMVL ◆ → GGAAGACCTGTGAACCG

    ★ Haplotype phasing + Imputation
       Estimation of HLA imputation for the tar
                                              Beagle
                                                                        Beagle
                                               SHAPEIT
                                                                        Minimac
                                               Eagle
```

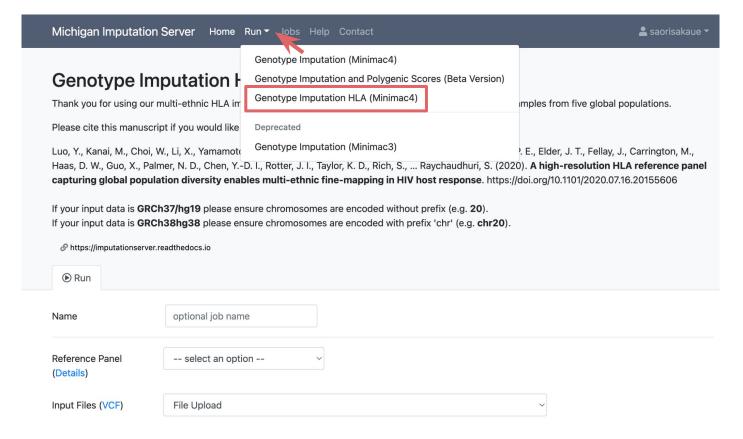
```
Known SNP genotype of the target cohort
 Genotype in the MHC region
  CGA.ATCT..GTCTTCTGT.CTAA
  CAA.ATCT..GTCCT.TGT.CTAA
  CAA.ATTT..TGCTTCAGT.CTAA
Given HLA imputation reference panel
                                           HLA amino acids
 Scaffold genotype in the MHC
                             HLA alleles
                                                              HLA intragenic SNPs
 CGAGATCTCAGTCTTCTAA ← DRB1*04:01 ← GGSCMAALTVTLMVL ← GGAGGACCTGTGAACCA
                     CTAA •—• DRB1*01:01•—• GGSCMTALTVTLMVL •—• GGAAGACCTGCGAACCA
 CGAGATCTCCTGCTTCAGTTCTAA ◆ → DRB1*01:02 ◆ → GGSCMTALTVTLMVL ◆ → GGAGGACCTGCGAACCA
 CAAGATCTCCGTCCTCTGTTCTAA ◆ → DRB1*15:01 ◆ → GGSCMTALTVTLMVL ◆ → GGAAGACCTGTGAACCG

    ★ Haplotype phasing + Imputation
       Estimation of HLA imputation for the tar
                                               Beagle
                                                                        Beagle
```

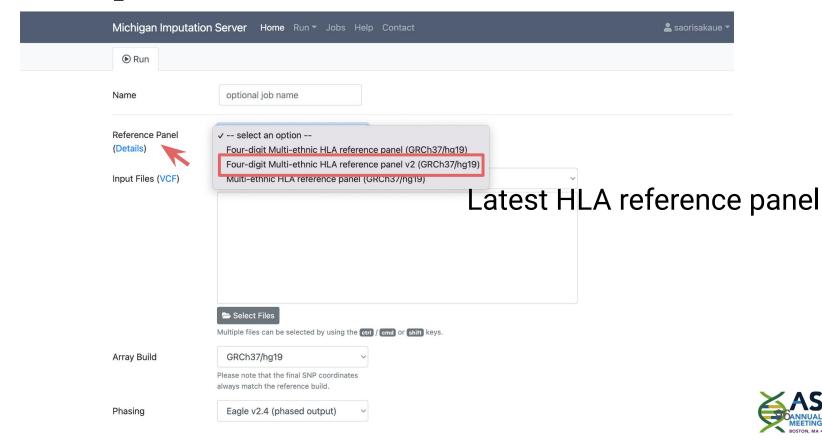
SHAPEIT

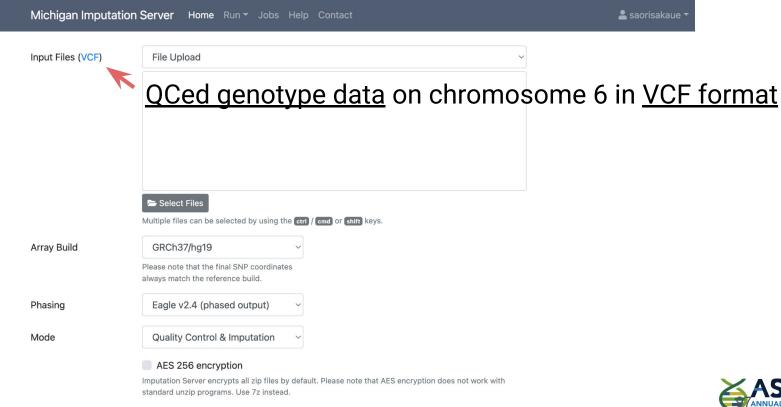
Eagle

GGAAGACCTGCGAACCA

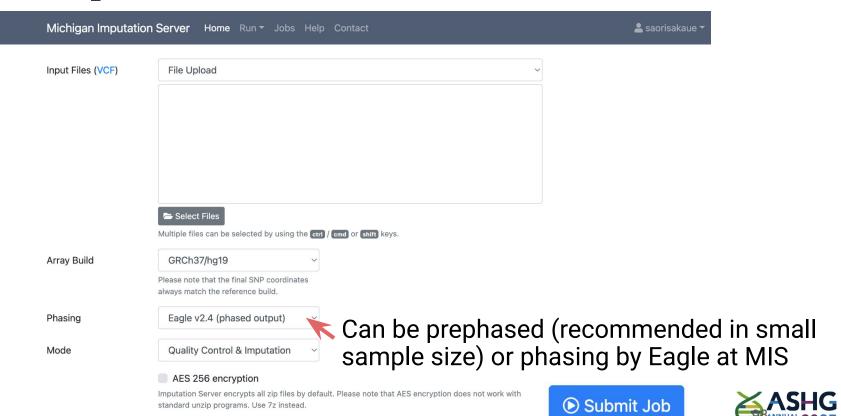












Let's look at the output from MIS!

```
ssakaue@wmbed-37d:~/Downloads/chr_6 (2)$ zcat chr6.dose.vcf.gz | less -S
```



Let's look at the output from MIS!

```
##fileformat=VCFv4.1
##filedate=2022.11.14
##contig=<ID=6>
##INFO=<ID=AF, Number=1, Type=Float, Description="Estimated Alternate Allele Frequency">
##INFO=<ID=MAF, Number=1, Type=Float, Description="Estimated Minor Allele Frequency">
##INFO=<ID=R2, Number=1, Type=Float, Description="Estimated Imputation Accuracy (R-square)">
##INFO=<ID=ER2, Number=1, Type=Float, Description="Empirical (Leave-One-Out) R-square (available only for genotyped variants)">
##INFO=<ID=IMPUTED,Number=0,Type=Flag,Description="Marker was imputed but NOT genotyped">
##INFO=<ID=TYPED, Number=0, Type=Flag, Description="Marker was genotyped AND imputed">
##INFO=<ID=TYPED ONLY, Number=0, Type=Flag, Description="Marker was genotyped but NOT imputed">
##FORMAT=<ID=GT, Number=1, Type=String, Description="Genotype">
##FORMAT=<ID=DS, Number=1, Type=Float, Description="Estimated Alternate Allele Dosage: [P(0/1)+2*P(1/1)]">
##FORMAT=<ID=HDS, Number=2, Type=Float, Description="Estimated Haploid Alternate Allele Dosage">
##FORMAT=<ID=GP, Number=3, Type=Float, Description="Estimated Posterior Probabilities for Genotypes 0/0, 0/1 and 1/1 ">
##pipeline=michigan-imputationserver-1.5.8
##imputation=minimac4-1.0.2
                                                          Imputed SNPs within MHC
##phasing=eagle-2.4
##panel=apps@multiethnic-hla-panel-4digit-v2@1.0.0
##r2Filter=0.0
#CHROM POS
                ID
                         REF
                                                  FILTER INFO
                                                                                                           010061321010_R02C01_10020793
                                                                                                                                            010061321010
                                                                  FORMAT
                                                                          010061321010 R01C01 10007854
        27970031
                         rs149946
                                      G
                                                                  PASS
                                                                          AF=0.22479; MAF=0.22479; R2=0.99214; IMPUTED
                                                                                                                            GT:DS:HDS:GP
                                                                                                                                            0 | 0:0:0,0:1,0
        27976200
                         rs9380032
                                         G
                                                                  PASS
                                                                          AF=0.02975; MAF=0.02975; R2=0.97885; IMPUTED
                                                                                                                            GT:DS:HDS:GP
                                                                                                                                            0|0:0:0,0:1,0
                                                                                                                                            0|0:0.003:0.0
        27979188
                         rs4141691
                                                                  PASS
                                                                          AF=0.11754; MAF=0.11754; R2=0.94642; IMPUTED
                                                                                                                            GT:DS:HDS:GP
        27979625
                         rs10484402
                                                                  PASS
                                                                          AF=0.04041; MAF=0.04041; R2=0.92706; IMPUTED
                                                                                                                            GT:DS:HDS:GP
                                                                                                                                            0 | 0:0.003:0.0
        27981673
                         rs9368540
                                                                  PASS
                                                                          AF=0.03706; MAF=0.03706; R2=0.98634; IMPUTED
                                                                                                                            GT:DS:HDS:GP
                                                                                                                                            0|0:0:0,0:1,0
                                                                  PASS
        27984726
                         rs74505854
                                                                          AF=0.00719; MAF=0.00719; R2=0.94391; IMPUTED
                                                                                                                            GT:DS:HDS:GP
                                                                                                                                            0|0:0:0,0:1,0
        27984907
                                                  C
                                                                  PASS
                         rs17765055
                                                                          AF=0.04632; MAF=0.04632; R2=0.99897; ER2=0.97293; TYPED
                                                                                                                                    GT:DS:HDS:GP
                                                                                                                                                     0 | 0:0
        27986199
                         rs72848791
                                                                  PASS
                                                                          AF=0.04361; MAF=0.04361; R2=0.99732; ER2=0.97465; TYPED
                                                                                                                                    GT:DS:HDS:GP
                                                                                                                                                    0|0:0
        27986529
                                                                  PASS
                                                                          AF=0.04631; MAF=0.04631; R2=0.99885; IMPUTED
                                                                                                                                            0|0:0:0,0:1,0
                         rs9368544
                                                                                                                            GT:DS:HDS:GP
                                         G
                                                                          AF=0.11782; MAF=0.11782; R2=0.99974; ER2=0.99765; TYPED
                                                                                                                                    GT:DS:HDS:GP
        27998258
                         rs149990
                                                                  PASS
                                                                                                                                                     010:0
        27999044
                                         A
                                                                  PASS
                                                                                                                                            0 | 0:0:0,0:1,0
                         rs9368545
                                                                          AF=0.04627; MAF=0.04627; R2=0.99825; IMPUTED
                                                                                                                            GT:DS:HDS:GP
        27999421
                         rs16893573
                                                                  PASS
                                                                          AF=0.02852; MAF=0.02852; R2=0.98601; IMPUTED
                                                                                                                            GT:DS:HDS:GP
                                                                                                                                            0 | 0:0:0,0:1,0
        28001003
                         rs17708949
                                                                  PASS
                                                                          AF=0.03124; MAF=0.03124; R2=0.98465; IMPUTED
                                                                                                                            GT:DS:HDS:GP
                                                                                                                                            0|0:0:0,0:1,0
                                                                  PASS
                                                                          AF=0.26963; MAF=0.26963; R2=0.99868; ER2=0.98731; TYPED
        28001610
                         rs149942
                                                                                                                                    GT:DS:HDS:GP
                                                                                                                                                     0 | 0:0
        28002388
                         rs149943
                                         G
                                                                  PASS
                                                                          AF=0.11781; MAF=0.11781; R2=0.99984; ER2=0.99889; TYPED
                                                                                                                                    GT:DS:HDS:GP
                                                                                                                                                     010:0
                                                                                                                                                    0|0:0
        28003271
                         rs183926
                                                                  PASS
                                                                          AF=0.00996; MAF=0.00996; R2=0.99551; ER2=0.95540; TYPED
                                                                                                                                    GT:DS:HDS:GP
```

Let's look at the output from MIS!

```
ssakaue@wmbed-37d:~/Downloads/chr_6 (2)$ zcat chr6.dose.vcf.gz
ssakaue@wmbed-37d:~/Downloads/chr_6 (2)$ zcat chr6.dose.vcf.gz | grep HLA | less -S
```



Imputed HLA alleles

	370000000000000000000000000000000000000							
6	29910247	HLA_A*01		T	PASS	AF=0.15456;MAF=0.15456;R2=0.99719;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910248	HLA_A*01:01	A	Т	PASS	AF=0.15234;MAF=0.15234;R2=0.99517;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910249	HLA_A*01:02	Α	T	PASS	AF=0.00110;MAF=0.00110;R2=0.85198;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910250	HLA_A*01:136	Α	T	PASS	AF=0.00002;MAF=0.00002;R2=0.04644;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910251	HLA_A*02	Α	Т	PASS	AF=0.26026;MAF=0.26026;R2=0.99741;IMPUTED	GT:DS:HDS:GP	0 1:0.998:0.0
6	29910252	HLA_A*02:01	Α	Т	PASS	AF=0.22914;MAF=0.22914;R2=0.98881;IMPUTED	GT:DS:HDS:GP	0 1:0.996:0.0
6	29910253	HLA_A*02:02	Α	T	PASS	AF=0.00471;MAF=0.00471;R2=0.99686;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910254	HLA_A*02:03	Α	T	PASS	AF=0.00094;MAF=0.00094;R2=0.88473;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910255	HLA_A*02:04	Α	Т	PASS	AF=0.00017;MAF=0.00017;R2=0.89691;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910256	HLA_A*02:05	Α	T	PASS	AF=0.01556;MAF=0.01556;R2=0.99839;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910257	HLA_A*02:06	Α	T	PASS	AF=0.00364; MAF=0.00364; R2=0.98705; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910258	HLA_A*02:07	Α	T	PASS	AF=0.00133;MAF=0.00133;R2=0.94265;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910259	HLA_A*02:10	Α	T	PASS	AF=0.00001;MAF=0.00001;R2=0.06774;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910260	HLA_A*02:11	Α	T	PASS	AF=0.00072;MAF=0.00072;R2=0.79302;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910261	HLA_A*02:135	Α	T	PASS	AF=0.00005; MAF=0.00005; R2=0.22275; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910262	HLA_A*02:17	Α	T	PASS	AF=0.00108;MAF=0.00108;R2=0.91078;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910263	HLA_A*02:195	Α	T	PASS	AF=0.00003;MAF=0.00003;R2=0.00997;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910264	HLA_A*02:20	Α	T	PASS	AF=0.00019; MAF=0.00019; R2=0.41177; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910265	HLA_A*02:22	Α	T	PASS	AF=0.00026;MAF=0.00026;R2=0.84742;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910266	HLA_A*02:279	Α	T	PASS	AF=0.00008;MAF=0.00008;R2=0.24552;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910267	HLA_A*02:55	Α	T	PASS	AF=0.00001; MAF=0.00001; R2=0.02045; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910268	HLA_A*02:56	Α	T	PASS	AF=0.00007;MAF=0.00007;R2=0.08360;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910269	HLA_A*02:60	Α	T	PASS	AF=0.00004; MAF=0.00004; R2=0.49123; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910270	HLA_A*02:76	Α	T	PASS	AF=0.00030; MAF=0.00030; R2=0.40308; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910271	HLA_A*02:87	Α	T	PASS	AF=0.00001;MAF=0.00001;R2=0.00984;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910272	HLA_A*03	Α	T	PASS	AF=0.12657; MAF=0.12657; R2=0.99727; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910273	HLA_A*03:01	Α	T	PASS	AF=0.12061; MAF=0.12061; R2=0.99073; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910274	HLA_A*03:02	Α	T	PASS	AF=0.00441; MAF=0.00441; R2=0.97833; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910275	HLA_A*03:36N	Α	Т	PASS	AF=0.00063; MAF=0.00063; R2=0.40664; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910276	HLA_A*03:89	Α	T	PASS	AF=0.00006; MAF=0.00006; R2=0.05589; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910277	HLA_A*11	Α	T	PASS	AF=0.06332; MAF=0.06332; R2=0.99759; IMPUTED	GT:DS:HDS:GP	1 0:1.000:1.0
6	29910278	HLA_A*11:01	Α	T	PASS	AF=0.05966; MAF=0.05966; R2=0.96821; IMPUTED	GT:DS:HDS:GP	1 0:0.958:0.9
6	29910279	HLA_A*11:02	Α	T	PASS	AF=0.00059; MAF=0.00059; R2=0.84655; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910280	HLA_A*11:32	Α	T	PASS	AF=0.00021; MAF=0.00021; R2=0.04209; IMPUTED	GT:DS:HDS:GP	0 0:0.024:0.0
6	29910281	HLA_A*11:50Q	Α	T	PASS	AF=0.00250; MAF=0.00250; R2=0.41268; IMPUTED	GT:DS:HDS:GP	0 0:0.018:0.0
6	29910282	HLA_A*23	Α	T	PASS	AF=0.02822;MAF=0.02822;R2=0.99569;IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910283	HLA_A*23:01	Α	T	PASS	AF=0.02797; MAF=0.02797; R2=0.99150; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0
6	29910284	HLA_A*23:15	Α	Т	PASS	AF=0.00000; MAF=0.00000; R2=0.00005; IMPUTED	GT:DS:HDS:GP	0 0:0:0,0:1,0

"A": Presence of the allele "A": Absence of the allele



29910247 HLA_A*01 PASS AF=0.15456; MAF=0.15456; R2=0.99719; IMPUTED GT:DS:HDS:GP 0 | 0:0:0,0:1,6 29910248 HLA A*01:01 PASS AF=0.15234; MAF=0.15234; R2=0.99517; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 HLA_A*01:02 AF=0.00110; MAF=0.00110; R2=0.85198; IMPUTED 0|0:0:0,0:1,0 6 29910249 PASS GT:DS:HDS:GP HLA A*01:136 PASS 0|0:0:0,0:1,0 29910250 AF=0.00002; MAF=0.00002; R2=0.04644; IMPUTED GT:DS:HDS:GP 29910251 HLA A*02 AF=0.26026; MAF=0.26026; R2=0.99741; IMPUTED 0|1:0.998:0.0 PASS GT:DS:HDS:GP 29910252 HLA_A*02:01 AF=0.22914; MAF=0.22914; R2=0.98881; IMPUTED GT:DS:HDS:GP 0|1:0.996:0.0 PASS HLA A*02:02 0|0:0:0,0:1,0 29910253 PASS AF=0.00471; MAF=0.00471; R2=0.99686; IMPUTED GT:DS:HDS:GP 29910254 HLA A*02:03 PASS AF=0.00094; MAF=0.00094; R2=0.88473; IMPUTED GT:DS:HDS:GP 0|0:0:0.0:1.0 HLA_A*02:04 0|0:0:0,0:1,6 29910255 PASS AF=0.00017; MAF=0.00017; R2=0.89691; IMPUTED GT:DS:HDS:GP 29910256 HLA A*02:05 PASS AF=0.01556; MAF=0.01556; R2=0.99839; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 0|0:0:0,0:1,0 29910257 HLA A*02:06 PASS AF=0.00364; MAF=0.00364; R2=0.98705; IMPUTED GT:DS:HDS:GP 29910258 HLA_A*02:07 PASS AF=0.00133; MAF=0.00133; R2=0.94265; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 29910259 PASS AF=0.00001; MAF=0.00001; R2=0.06774; IMPUTED 0|0:0:0,0:1,6 HLA A*02:10 GT:DS:HDS:GP 29910260 HLA A*02:11 PASS AF=0.00072; MAF=0.00072; R2=0.79302; IMPUTED GT:DS:HDS:GP 0|0:0:0.0:1.0 AF=0.00005; MAF=0.00005; R2=0.22275; IMPUTED 29910261 HLA_A*02:135 PASS GT:DS:HDS:GP 0 | 0:0:0,0:1,6 29910262 HLA A*02:17 **PASS** GT:DS:HDS:GP 0|0:0:0,0:1,6 AF=0.00108; MAF=0.00108; R2=0.91078; IMPUTED PASS 29910263 HLA_A*02:195 AF=0.00003; MAF=0.00003; R2=0.00997; IMPUTED GT:DS:HDS:GP 0 | 0:0:0,0:1,0 29910264 HLA_A*02:20 **PASS** AF=0.00019; MAF=0.00019; R2=0.41177; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 29910265 HLA A*02:22 PASS AF=0.00026; MAF=0.00026; R2=0.84742; IMPUTED GT:DS:HDS:GP 0 0:0:0,0:1,6 HLA_A*02:279 29910266 PASS AF=0.00008; MAF=0.00008; R2=0.24552; IMPUTED GT:DS:HDS:GP 0 | 0:0:0,0:1,0 29910267 HLA A*02:55 PASS AF=0.00001; MAF=0.00001; R2=0.02045; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 29910268 HLA A*02:56 PASS AF=0.00007; MAF=0.00007; R2=0.08360; IMPUTED GT:DS:HDS:GP 0|0:0:0.0:1.0 29910269 HLA_A*02:60 **PASS** AF=0.00004; MAF=0.00004; R2=0.49123; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 29910270 HLA A*02:76 **PASS** 0|0:0:0,0:1,0 AF=0.00030; MAF=0.00030; R2=0.40308; IMPUTED GT:DS:HDS:GP 29910271 HLA_A*02:87 PASS 0|0:0:0,0:1,0 AF=0.00001; MAF=0.00001; R2=0.00984; IMPUTED GT:DS:HDS:GP HLA_A*03 29910272 PASS AF=0.12657; MAF=0.12657; R2=0.99727; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 29910273 HLA A*03:01 PASS 0|0:0:0,0:1,0 AF=0.12061; MAF=0.12061; R2=0.99073; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 29910274 HLA A*03:02 PASS AF=0.00441; MAF=0.00441; R2=0.97833; IMPUTED GT:DS:HDS:GP HLA_A*03:36N 0|0:0:0,0:1,0 29910275 PASS AF=0.00063; MAF=0.00063; R2=0.40664; IMPUTED GT:DS:HDS:GP 29910276 HLA A*03:89 PASS AF=0.00006; MAF=0.00006; R2=0.05589; IMPUTED GT:DS:HDS:GP 0|0:0:0,0:1,0 AF=0.06332; MAF=0.06332; R2=0.99759; IMPUTED 29910277 HLA A*11 PASS GT:DS:HDS:GP 1|0:1.000:1.0 29910278 HLA A*11:01 **PASS** 1|0:0.958:0.9 AF=0.05966; MAF=0.05966; R2=0.96821; IMPUTED GT:DS:HDS:GP 29910279 HLA A*11:02 PASS 0|0:0:0,0:1,0 AF=0.00059; MAF=0.00059; R2=0.84655; IMPUTED GT:DS:HDS:GP 29910280 HLA_A*11:32 PASS AF=0.00021; MAF=0.00021; R2=0.04209; IMPUTED GT:DS:HDS:GP 010:0.024:0.6 PASS 29910281 HLA A*11:50Q AF=0.00250; MAF=0.00250; R2=0.41268; IMPUTED GT:DS:HDS:GP 0|0:0.018:0.6

"A": Presence of the allele

AF=0.02822; MAF=0.02822; R2=0.99569; IMPUTED

AF=0.02797; MAF=0.02797; R2=0.99150; IMPUTED

AF=0.00000; MAF=0.00000; R2=0.00005; IMPUTED

PASS

PASS

PASS

29910282

29910283

29910284

HLA A*23

HLA_A*23:01

HLA A*23:15



0|0:0:0,0:1,0

0 0:0:0,0:1,0

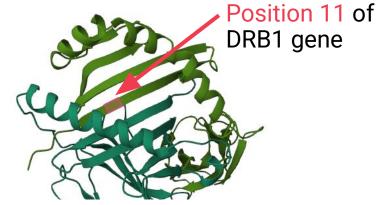
0|0:0:0,0:1,0

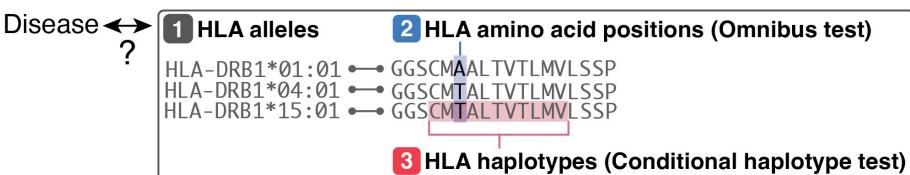
GT:DS:HDS:GP

GT:DS:HDS:GP

GT:DS:HDS:GP

How are imputed HLA variants associated with disease?



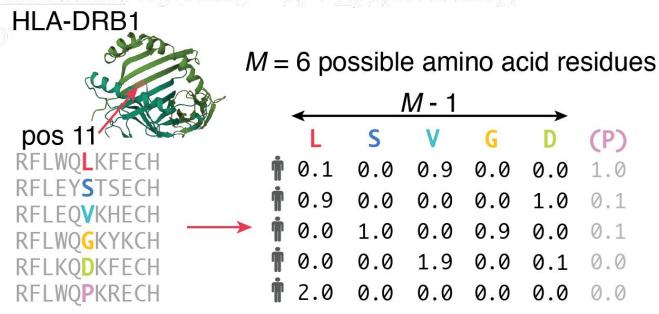




Omnibus test to ask which HLA amino acid position affects the disease

Full model:
$$log(odds_i) = \beta_0 + \sum_k \beta_k covariate_{k,i} + \sum_{m=1}^{M-1} \beta_m AM_{m,i}$$

Reduced model: $log(odds_i) = \beta_0 + \sum_k \beta_k covariate_{k,i}$





Omnibus test to ask which HLA amino acid position affects the disease

Full model: $log(odds_i) = \beta_0 + \sum_k \beta_k covariate_{k,i} + \sum_{m=1}^{M-1} \beta_m AM_{m,i}$

Reduced model: $log(odds_i) = \beta_0 + \sum_k \beta_k covariate_{k,i}$

ANOVA(Full model, Reduced model)

How much does this amino acid position's polymorphism increase the explained variance for the trait?

→ Determine the <u>single</u> most significant amino acid position



Summary

- 1. HLA amino acid sequences and alleles characterize antigen presentation and disease risk within HLA.
- 2. HLA amino acid sequences and alleles can be accurately imputed from genotyped SNPs by MIS.
- 3. Imputed HLA alleles can be used to fine-map causal disease mechanisms.



Summary

- 1. HLA amino acid sequences and alleles characterize antigen presentation and disease risk within HLA.
- 2. HLA amino acid sequences and alleles can be accurately imputed from genotyped SNPs by MIS.
- 3. Imputed HLA alleles can be used to fine-map causal disease mechanisms.

nature protocols

Review article

https://doi.org/10.1038/s41596-023-00853-4

Tutorial: a statistical genetics guide to identifying HLA alleles driving complex disease Sakaue et al. Nature Protocols 2023





Section 4 nf-gwas and PGS Server



Lukas Forer

Medical University of Innsbruck

lukas.forer@i-med.ac.at

@lukfor





Learning objectives

Participants will

- 1. Learn how to run a GWAS pipeline
- 2. Learn how to use the PGS extension
- 3. Learn how to analyze PGS calculated by the server



Summary of common GWAS analysis tools

	EPACTS	Rvtests	SNPTEST	SAIGE	BLOT-LMM	Bgenie	regenie
Input VCF	Υ	Υ	Y				
Sample relatedness (Quantitative outcome)	Y	Y		Υ	Y	Υ	Υ
Sample relatedness (Binary outcome)				Υ		Υ	Υ
Case control imbalance				Υ			Υ
Large sample size (>20,000)				Υ	Υ	Υ	Y



nf-gwas

- A GWAS pipeline based on REGENIE and works with VCF input
- Includes several pre- and post-processing steps
- Based on Nextflow
 - Allows to build a portable, reproducible, scalable pipeline
 - Runs on clusters (e.g. SLURM) or in the cloud (e.g. AWS Batch)

JOURNAL ARTICLE

Performing highly parallelized and reproducible GWAS analysis on biobank-scale data 3

Sebastian Schönherr, Johanna F Schachtl-Riess, Silvia Di Maio, Michele Filosi,
Marvin Mark, Claudia Lamina, Christian Fuchsberger, Florian Kronenberg, Lukas Forer

Author Notes

NAR Genomics and Bioinformatics, Volume 6, Issue 1, March 2024, Iqae015, https://doi.org/10.1093/nargab/Iqae015





Workflow **Download** (rsIDs, genes) Imputed MICHIGAN IMPUTATIONSERVER Genotypes variants Impute to plink2 Liftover Annotated. tbi-indexed summary statistics Create **Validate Phenotypes** Step 1 txt Step 2 report Interactive html report Annotate Merge Covariates by phenotype Log files Genotyped plink variants QC filter Prune Condition txt Single-variant testing list Interaction testing Interaction variable **GxE** (environment) Gene-based testing OΓ GxG (SNP) Annotations, Set list, txt Masks



Download Phenotypes

Phenotype Data

- Filename: phenotypes.txt (131.6 KB)
- **Description:** Simulated phenotype file containing four traits.
- Format: Plain text (tab-delimited).
- Columns:
 - sample_id individual sample identifiers
 - pheno_1, pheno_2, pheno_3, pheno_4 four simulated phenotypes
- Use: Provides phenotype data for GWAS and PGS evaluation.

Covariates (Principal Components)

- Filename: covariates.txt (292.3 KB)
- Description: Covariate file containing the first 10 genetic principal components (PCs)
- Format: Tab-delimited, compressed with gzip.
- Columns:
 - sample individual identifiers matching the genotype data
 - PC1 PC10 first ten principal components representing population structure
- Use: Used as covariates in association analyses to correct for population stratification

You can also download the pre-imputed data

Simulated Phenotypes

Phonotypo	Typo	Highly Genetic?
Phenotype	Туре	riginty defletic:
pheno_1	Binary (e.g., case/control)	▼ Yes
pheno_2	Continuous (e.g., height)	▼ Yes
pheno_3	Binary or categorical	X No
pheno_4	Continuous	X No



https://tinyurl.com/imputationserver-2025



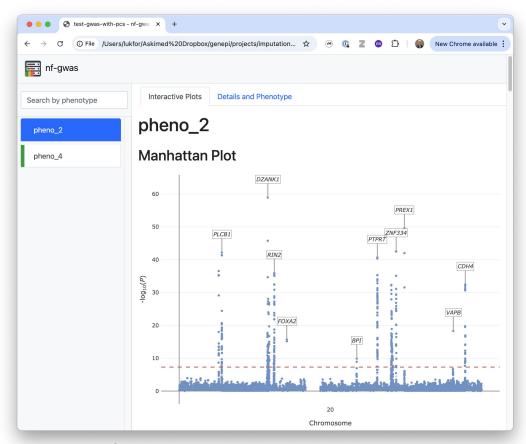
Easy to configure

my-gwas.config

```
params {
 project
                                = 'test-gwas-with-pcs'
                                = 'gwas.imputed.chr20.dose.vcf.gz'
  genotypes association
                                                                                        Imputed genotypes
  genotypes_association_format
                                = 'vcf'
                                                                                        from Imputation Server
  association build
                                = 'hq19'
  phenotypes filename
                                = 'phenotypes.txt'
                                                                                        phenotypes
  phenotypes columns
                                = 'pheno 2, pheno 4'
  covariates_filename
                                = 'covariates.txt'
                                                                                        covariates
  covariates columns
                                = 'PC1, PC2, PC3, PC4, PC5, PC6, PC7, PC8, PC9, PC10'
  regenie test
                                = 'additive'
  regenie_min_imputation_score
                                = 0.3
  regenie_skip_predictions
                                = true
                                                                                        Annotation and
  rsids filename
                                = "rsids-v154-hg19-chr20.index.gz"
                                                                                        visualization
 binning size
                                = 50000
```



Results Phenotype 2



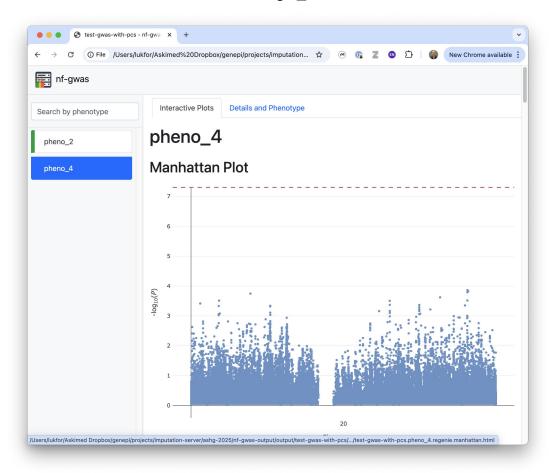
Top Loci

Variant	ID *	Ref Allele	Effect Allele	Nearest Gene	- log ₁₀ (<i>P</i>) [♥]	Beta
20:18445266	rs6111973	С	Т	DZANK1	58.9235	-0.522676
20:46858654	rs7269193	Α	G	PREX1	49.648	0.626932
20:45142619	rs847055	Α	Т	ZNF334	42.5918	0.493058
20:8887891	rs6056246	Т	С	PLCB1	42.2005	0.437334
20:41227725	rs6102924	С	Т	PTPRT	40.7084	-0.627067
20:19800969	rs1884767	Α	G	RIN2	35.9524	-0.385179
20:59507557	rs6071441	С	Т	CDH4	32.4673	0.480876
20:57020044	rs6070466	С	G	VAPB	18.2868	0.925181
20:22392054	rs13038903	Α	G	FOXA2	15.7283	-3.3011
20:36950990	rs6024863	Α	Т	BPI	9.8814	0.225461
20:1728212	rs605739	Т	G	SIRPG	3.95055	0.124383
20:4453463	rs297660	С	Т	PRNP	3.89483	-0.14028
20:2513877	rs6083567	Α	G	TMC2	3.79769	-0.217792
20:21387531	rs13040764	Α	G	NKX2-4	3.7689	-0.763363
20:31440180	rs6088022	Т	G	MAPRE1	3.69196	-0.442159
20:25957535	rs2252934	G	С	ZNF337	3.57123	-0.184881
20:29859763	rs6089068	С	Т	DEFB115	3.51671	0.189548
20:53303956	rs6023507	С	T	DOK5	3.47376	0.357591
20:25221547	rs6050477	T	С	PYGB	3.47251	-0.120575

Loci are defined ±200 kb around lead SNPs



Results Phenotype 4



No significant loci.



Tutorial: Setup, Parameters, and Example Analyses

- On the workshop website, you'll find a step-by-step tutorial on how to:
 - Setup and install the pipeline
 - Explore all available parameters
 - Run association tests for binary traits (e.g. pheno_1 and pheno_2)
 - Generate and interpret QQ plots
- Also available for PLINK2





JOURNAL ARTICLE

Imputation Server PGS: an automated approach to calculate polygenic risk scores on imputation servers 3

Nucleic Acids Research, Volume 52, Issue W1, 5 July 2024, Pages W70–W77, https://doi.org/10.1093/nar/gkae331



Imputation Server PGS

- Makes it easier to use polygenic scores (PGS) with imputed genotypes by calculating them directly on the server
- PGS: aggregates the effects of many genetic variants into a single number which predicts genetic predisposition for a phenotype
- Available since April 2022
- More than 5,000 submitted jobs
- Supports more than 3,000 scores

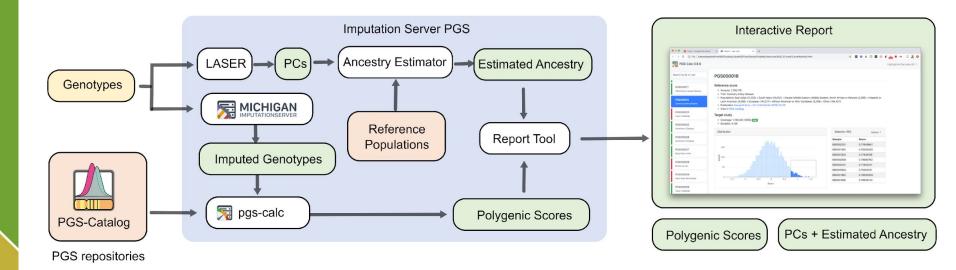


Have you previously used the PGS tool on the Michigan Imputation Server?

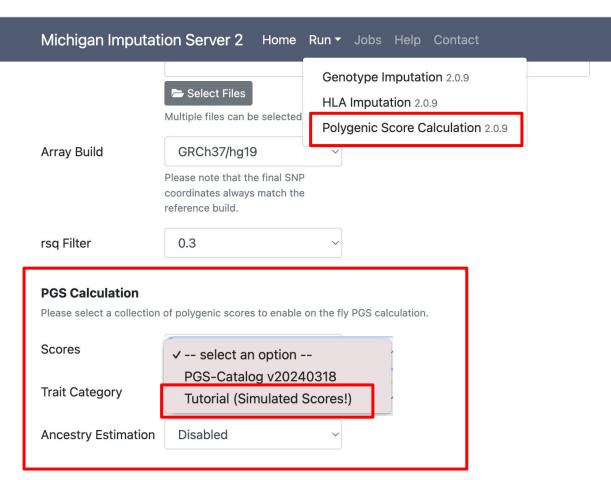




Workflow









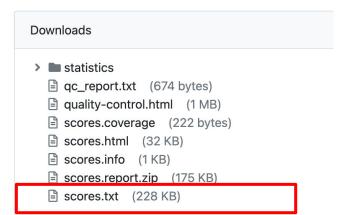
≜ lukfor ▼

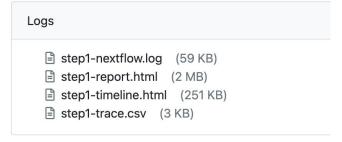
Results

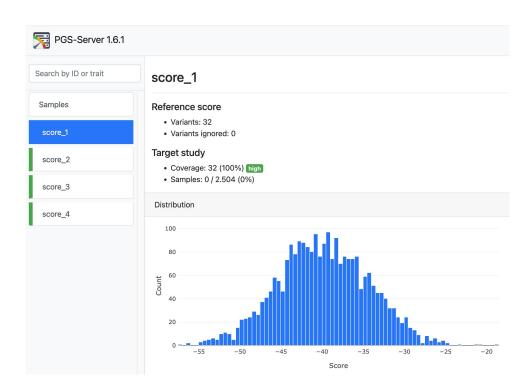
Details

Results

Logs









What can we do with these scores now?

Example: compare how well each score explains the simulated phenotypes

Merge scores with phenotypes

```
# Read PRS scores and phenotypes
scores <- read_csv("data/scores.txt")
phenos <- read_table("data/phenotypes.txt")

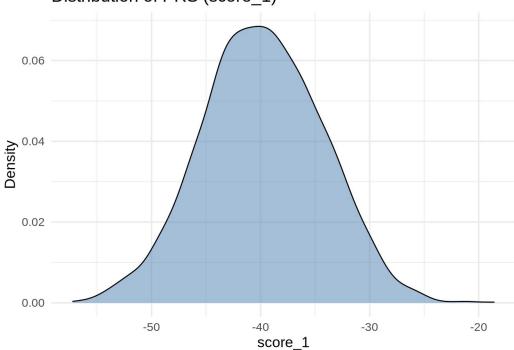
# Merge datasets by sample ID
merged <- inner_join(scores, phenos, by = c("sample" = "IID"))</pre>
```

You can also download the calculated polygenic scores from the website



```
ggplot(merged, aes(x = score_1)) +
  geom_density(alpha = 0.5, fill = "steelblue") +
  labs(
    title = "Distribution of PRS (score_1)",
    x = "score_1",
    y = "Density"
)
```

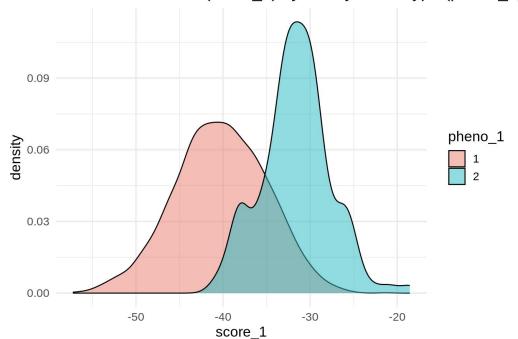


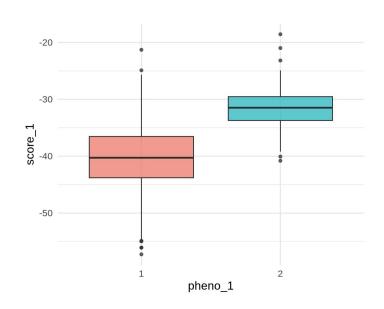




```
ggplot(merged, aes(x = score_1, fill = as.factor(pheno_1))) +
  geom_density(alpha = 0.5) +
  labs(
    title = "Distribution of PRS (score_1) by Binary Phenotype (pheno_1)",
    x = "score_1",
    fill = "pheno_1"
)
```

Distribution of PRS (score_1) by Binary Phenotype (pheno_1)

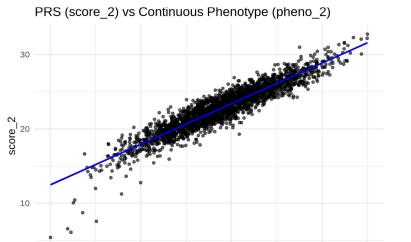






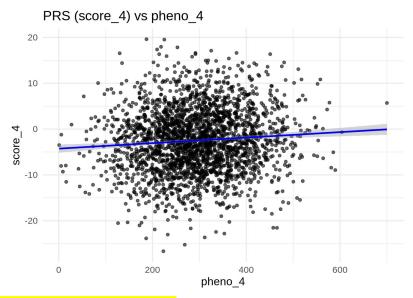
```
ggplot(merged, aes(x = pheno_2, y = score_2)) +
  geom_point(alpha = 0.6) +
  geom_smooth(method = "lm", se = TRUE, color = "blue") +
  labs(
    title = "PRS (score_2) vs Continuous Phenotype (pheno_2)",
    x = "pheno_2",
    y = "score_2"
)
```

```
ggplot(merged, aes(x = pheno_4, y = score_4)) +
geom_point(alpha = 0.6) +
geom_smooth(method = "lm", se = TRUE, color = "blue") +
labs(
   title = "PRS (score_4) vs pheno_4",
   x = "pheno_4",
   y = "score_4"
)
```



pheno_2

200





600



Summary

- Imputation Servers are easy to use and ensure high-quality imputation
- GWAS and PGS play a key role in modern genetic research
- Several pipelines and extensions of the Imputation Server

Aim: Taking the burden out of genetic analysis



Section 5 Helmholtz Munich Imputation Server



Eleftheria Zeggini Director, Institute of Translational Genomics





Why do we need another imputation server?

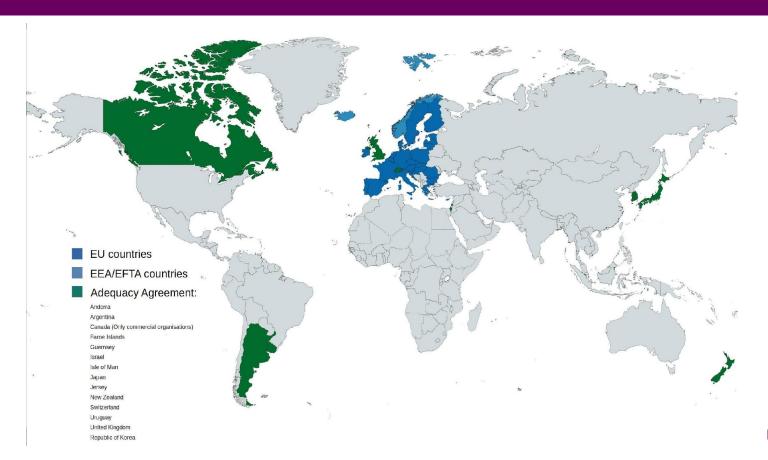


- The transfer of personal data to countries outside the European Economic Area (EEA) is generally prohibited by default
- Therefore, all personal data on EU citizens must reside and be processed on servers physically located within the EU
- Human genetic data are considered a special category of personal data subject to stricter requirements
- There are reasons when personal data can be moved outside the EEA:
 - Consent is given by the individual for the transfer
 - A Standard Contractual Clause (SCC) has been agreed between the two parties
 - There is adequacy agreement with country



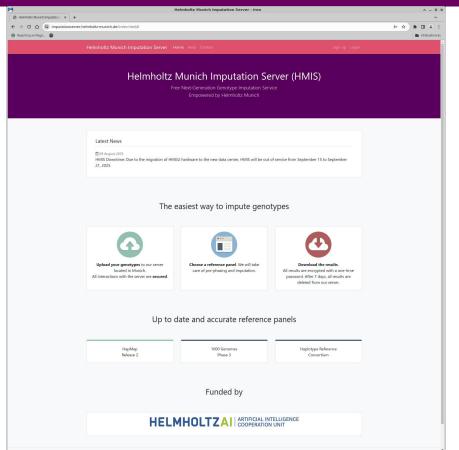
EU Adequacy Agreements





Helmholtz Munich Imputation Server





The server is designed to assist users in the EU to comply with the 2016 General Data Protection Regulation (GDPR) law

nature genetics

Correspondence Published: 15 November 2024

Toward GDPR compliance with the Helmholtz Munich genotype imputation server

N. William Rayner ☑, Young-Chan Park, Christian Fuchsberger, Andrei Barysenka & Eleftheria Zeggini ☑

Nature Genetics 56, 2580–2581 (2024) Cite this article



HELMHOLTZ MUNICI

Helmholtz Munich Imputation Server



- Based on the freely available Michigan software
- Currently version 2 based on Cloudgene and Hadoop
- Moving to version 3, currently integrating the changes we have made from the standard Michigan software

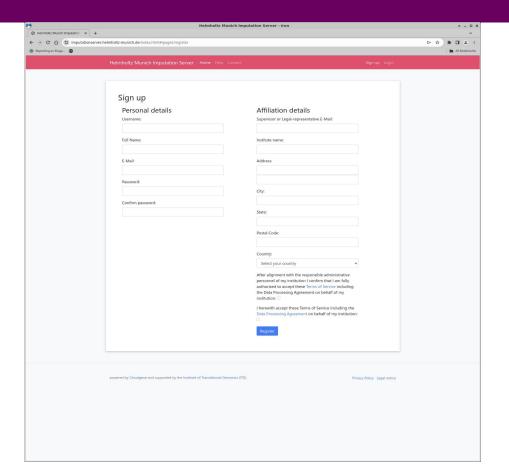
Notable Differences vs Michigan



More information is required on sign up to comply with GDPR

HMIS sign up





Institute name and address



email of legal representative and country are required so as to meet our obligations under GDPR

Notable Differences vs Michigan



- More information is required on sign up to comply with GDPR
- Short and long (>25,000 samples) processing queues introduced to ensure fair access to the resource
- Longer more complex passwords are now required, >15 characters

Current Reference Panels

- 1000G GRCh37
- 1000G GRCh38
- HRC
- HapMap2



Upcoming Reference Panels

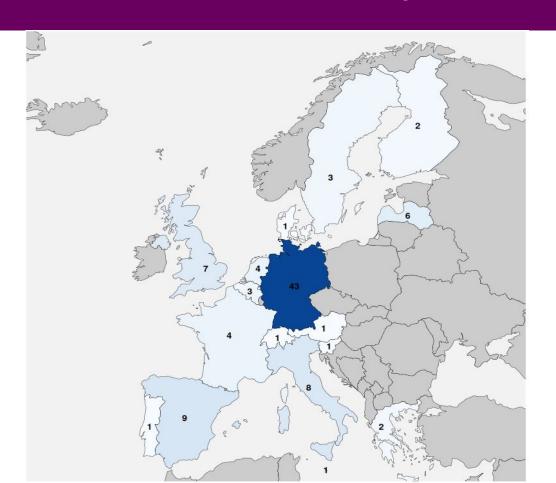
- German National Cohort (NAKO)
 ~15,500
 Will be expanded to 30,000
- Genome of Europe





Registered Users





We have 101 registered users, largely from across the EU

Number of genomes

 Over 900,000 genomes imputed to date



Direct future developments

0 surveys completed

0 surveys underway



Section 6 The TOPMed Imputation Server



Albert Smith
University of Michigan
albertvs@umich.edu
albertvs@umich.edu
avsmith

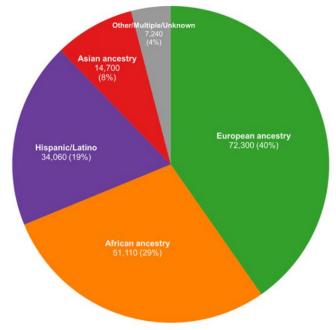


TOPMed Program

- Trans-Omics for Precision Medicine (TOPMed) Program
- A Precision Medicine Initiative sponsored by National Heart, Lung and Blood Institute
- Integrating whole-genome sequencing and other omics data
- >180k participants from >90 studies

Ancestry & Ethnicity

Phases 1-7 (~180K Participants)





TOPMed Imputation

- Current reference panel based on TOPMed Freeze 8 Calls
- Michigan Imputation Server ported to Amazon Web Services
- R2 panel released April 2020
- Updated R3 panel released December 2023
- https://imputation.biodatacatalyst.nhlbi.nih.gov
- Registration as before, open access to TOPMed panel



TOPMed Panel Compared

	TOPMed_r2	HRC	1000G Genomes
N samples	97K	39K	2,500
Ancestry	Multiethnic	European	Multiethnic
N variants	308M	39M	88M
Avg. depth	38X	8X	4X
Genome build Position	b38	b37	b37

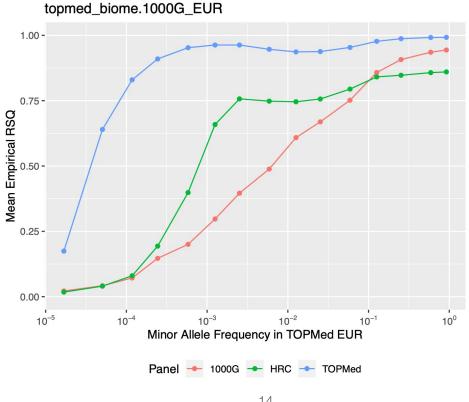


TOPMed Panel Freq Distribution

Variation type	Non-reference allele frequency bins				
- 31	(0, 0.005]	(0.005, 0.01]	(0.01, 0.05]	(0.05, 1)	Totals
SNVs	270,352,495	3,365,284	5,330,340	7,020,861	286,068,980
Insertions	5,462,262	74,150	130,506	148,595	5,815,513
Deletions	15,406,052	185,606	297,186	333,748	16,222,592
Totals	291,220,809	3,625,040	5,758,032	7,503,204	308,107,085

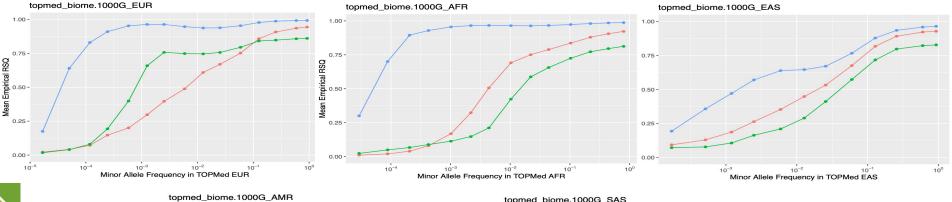


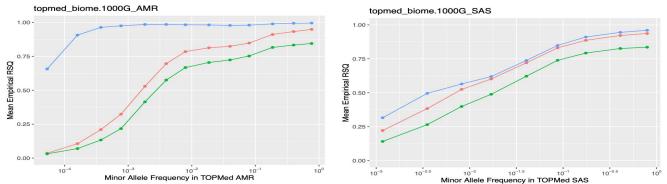
Imputation Panel Quality





Imputation Panel Quality

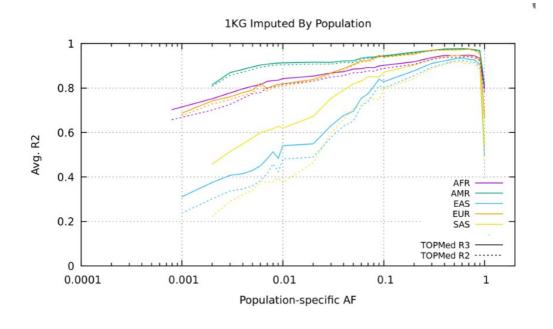






Updated TOPMed R3 Panel

- Existing panel has better coverage for EUR, AFR and AMR samples
- Targeted improvement for SAS and EAS samples
- Expanded panel (97k => 143k)
- Released Dec 2023 (R2 deprecated)





TOPMed Panel Compared

7	ΓΟΡMed_r3	TOPMed_r2	HRC	1000G Genomes
N samples	143K	97K	39K	2,500
Ancestry	Multiethnic	Multiethnic	European	Multiethnic
N variants	414M	308M	39M	88M
Avg. depth	38X	38X	8X	4X
Genome build Position	b38	b38	b37	b37

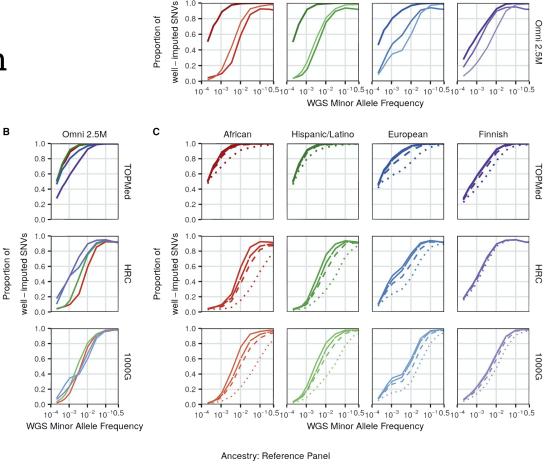
TOPMed_r2: Deprecated Dec 2025 TOPMed_r3: Contains HRC and HGDP

HRC: Haplotype Reference Consortium Panel 1000G: 1000 Genomes Project Reference Panel



TOPMed Imputation Compared to WGS

- Proportion well imputed (r2 > 0.8) down to MAF:
 - 0.14% in African
 - o 0.11% in Hispanic/Latino
 - o 0.35% in European
 - o 0.85% in Finnish
- Similar performance for arrays with >700k variants
- Source: Hanks et al. https://doi.org/10.1016/j.ajhg.202 2.07.012



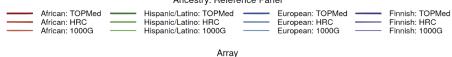
Α

African

Hispanic/Latino

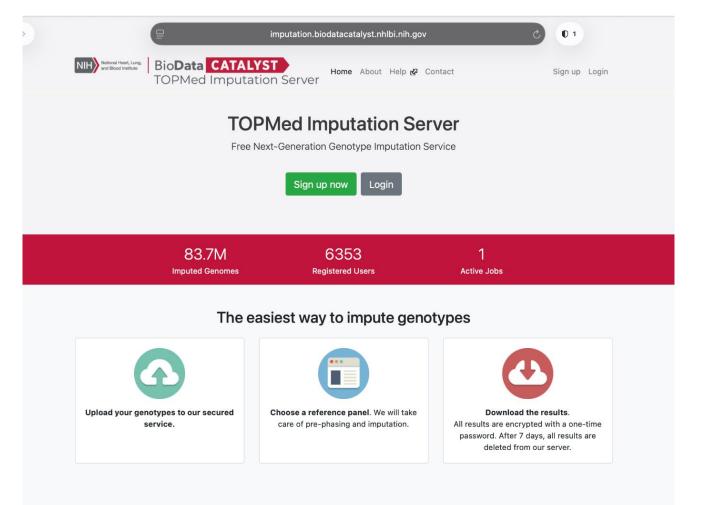
European

Finnish



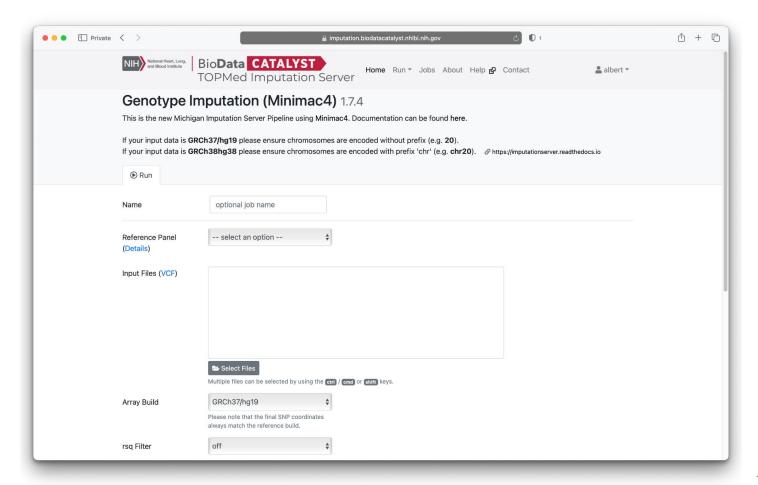
Omni 2.5M (2.4M) ———— MEGA (1.7M) ———— Omni Express (0.7M) · · · · · · Core (0.3M)







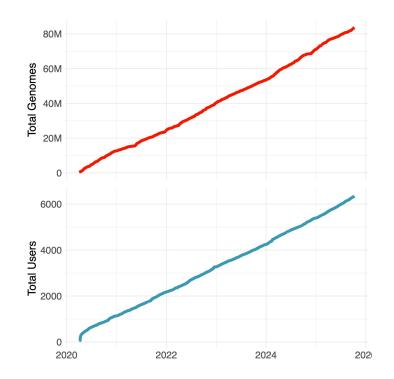
The TOPMed Imputation Server is powered by software invented and developed by the University of Michigan & and driven by data provided by the investigators of the TOPMed Program &.





TOPMed Imputation

- 84M genomes imputed
- Supplanted 1000g & HRC imputation for many studies
- Particularly benefits ethnically diverse cohorts
- Satisfying GDPR-related concerns of European users remains a challenge (Better served by Munich Imputation Server)





TOPMed Imputation

- TOPMed Server migrating to Nextflow based architecture in near future
- Beta testers needed
- If you have imputation ready samples, contact me <albertvs@umich.edu>
 or our support address <imputationserver@umich.edu>
- NOTE: This will be the same R3 panel. Testers should have new samples.



Imputation Resources

- Michigan Imputation Server
 https://imputation.sph.umich.edu/
- Helmholtz Munich Imputation Server
 https://imputationserver.helmholtz-munich.de/
- TOPMed Imputation Server
 https://imputation.biodatacatalyst.nhlbi.nih.gov/
- MCPS Imputation Server

https://imputationserver-reg.sph.umich.edu/



Your questions:

Nobody has responded yet.

Hang tight! Responses are coming in.



Thank you!



Your questions

If we run out of time, please send your questions to

Christian Fuchsberger - cfuchsb@umich.edu

