# Imputation using Sanger institute - draft

Jean-Tristan Brandenburg

July 8, 2022

## 1   Introduction : imputation genetics

Imputation in genetics refers to the statistical inference of unobserved geno-
types or missing data : using Linkage disequilibrium and/or reference panel,
objectives if to defined missing information and imputed position between your
positions.

   Common step before imputation, imputation and after imputation could be
:

- Quality Control of your genetics data, for instance using qc from h3agwas
  piplene

- Format and prepared your file for imputation

- Phasing your data : Phasing is the task or process of assigning alleles (the
  As, Cs, Ts and Gs) to the paternal and maternal chromosomes)

- Imputed data using reference panel.

- post qc your data and format your data to used Genome wide associaiton
  software

Imputation quality depend of software and panel imputation

### 1.1   resource

- Genotype Imputation [1]

## 2   Steps of imputation

### 2.1   QC

for more information on Quality Control see qc from h3agwas piplene and [2]

## 2.2   Format data

Main of imputation pipeline used VCF as input. When you prepared your data of array intialy in Plink in VCF must be done carefully. Pipeline can rejected your vcf file if allele reference is not good for instance. h3agwas pipeline can transform your plink data in vcf and check your VCF.

For instance :

```
wget -c http://ftp.ensembl.org/pub/grch37/current/fasta/homo_sapien/dna/Homo_sapiens.GRCh37
```

```
nextflow h3abionet/h3agwas/formatdata/plk_in_vcf_imp.nf \
 --input_dir dataqc \
 --input_pat exampledata2_qc \
 --reffasta Homo_sapiens.GRCh37.dna.primary_assembly.fa.gz \
 --output_pat exampledata2_qc \
 --output_dir ./vcfforimputation \
  -profile slurmSingularity \
 -resume
```

## 2.3   Phasing and imputation

Different way exist to phase and impute, using pipeline for instance

- chipimputation pipeline from h3abionet

- Sanger institute online server

- Michigan online server.

On Sanger and Michigan of online server proposed specific imputation panel as sanger with african panels.

# 3   Sanger imputation

## 3.1   Parameters

Sanger imputation, read about paramaters of choice :

- important to read about reference panels
- Important to read about pipeline

## 3.2   African panels from sanger

Some specificity of pipeline there are not insertion deletion.

## 3.3   Transfer of data

To transfer data, Sanger used globus.

### 3.3.1 globus installation and run

- create you need to created a account on globus

You can connect to globus or used globus information relatives to download

```
mkdir globus
 cd globus/
# download
wget https://downloads.globus.org/globus-connect-personal/linux/stable/globusconnectpersonal
# extract globus
tar xzf globusconnectpersonal-latest.tgz
# go on directory
 cd globusconnectpersonal-3.1.3/
 # connect your session follow instruction
 ./globusconnect -setup -nogui
 #open a screen
```

To open

```
 screen
 # open globus
 ./globusconnect -start  \# start globus
 # close screen ctrl-a ctrl-d
```

### 3.3.2 Created jobs and transfert data

AWIGEN have been imputed using african panels from sanger with PBWT and check what phasing.

- Sanger institute used globus to transfer file see globus section, you need to follow globus.

- created jobs for imputation on imputation sanger,see 1, 2

- You will received emel 3 of validation follow link selected your folder on the right, selected (drag and drop from your folder) your vcf file to transfert. 4

- wait confirmation of transfert by globus

- you will receive also a emel from sanger, when globus validate transfer by emel 5 click on link of sanger 6 to validate your sanger.7

### 3.3.3 Obtained your data imputed using globus

Waiting that you confirmation emel that well running from server and data ready from globus, click on link (be sure that you globus is on) see 7. transfers data on cluster using glob, in "activity" of website you can follow transfers otherwise you will receive a message of confirmation 8, 9. You can deleted your data after transfer have been done, 10.

Figure 1: parameter of server

```
## in running
[globusconnectpersonal-3.1.3]$ ./globusconnect -status
Globus Online:   connected
Transfer Status: active

## finish : status is idle
[globusconnectpersonal-3.1.3]$ ./globusconnect -status
Globus Online:   connected
Transfer Status: idle
```

# 4    Transform your VCF imputed in plink

To transform your files vcf in plink, you can used script in h3abionet/h3agwas
pipeline

```
ls imp/vcf/*.vcf.gz > listvcf
wget -c http://ftp.ensembl.org/pub/grch37/current/fasta/homo_sapiens/dna/Homo_sapiens.GRCh37
```
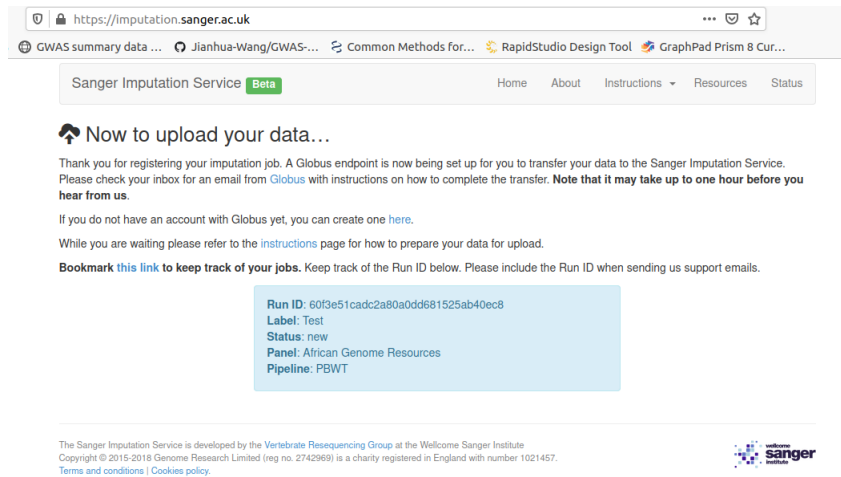
Figure 2: after validation



Figure 3: Emel receive from globus to select job, follow job

```
nextflow run h3abionet/h3agwas/formatdata/vcf_in_plink.nf \
  --file_listvcf listvcf \
  --output_pat exampledata2_imp \
  --output_dir fileimputation/ \
  --reffasta Homo_sapiens.GRCh37.dna.primary_assembly \
-profile slurmSingularity
```

# References

[1]  Yun Li et al. "Genotype Imputation". In: *Annual review of genomics and human genetics* 10 (2009), pp. 387–406. ISSN: 1527-8204. DOI: 10.1146/annurev.genom.9.081307.164242. URL: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2925172/ (visited on 07/06/2022).
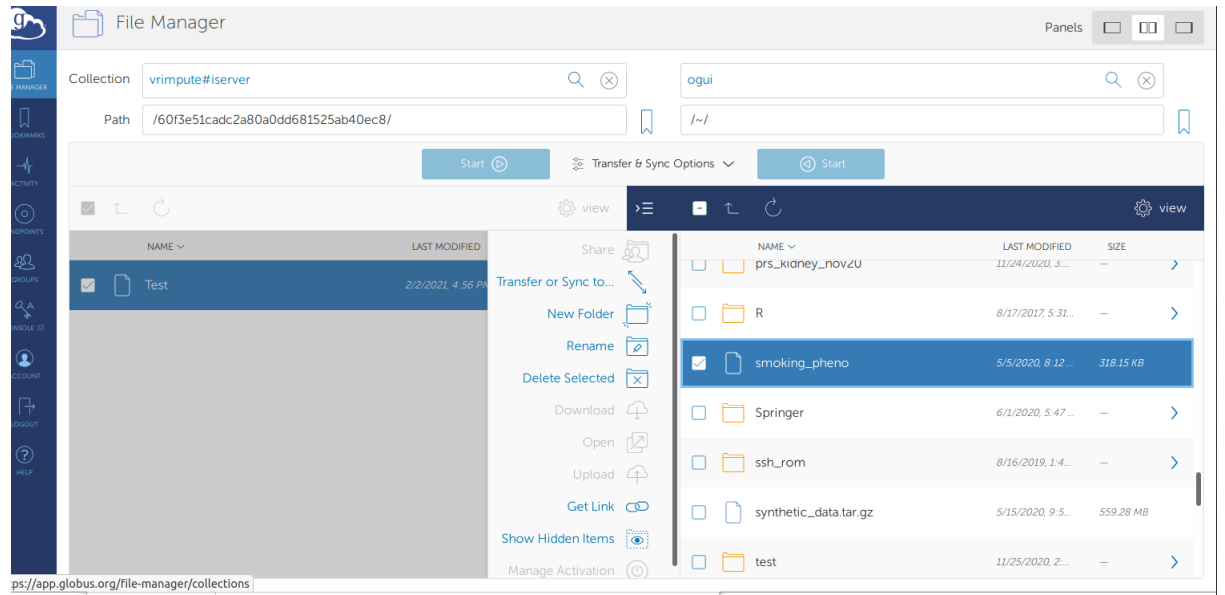
Figure 4: transfert data

[2] Stephen Turner et al. "Quality Control Procedures for Genome Wide Association Studies". In: *Current protocols in human genetics / editorial board, Jonathan L. Haines ... [et al.]* CHAPTER (Jan. 2011), Unit1.19.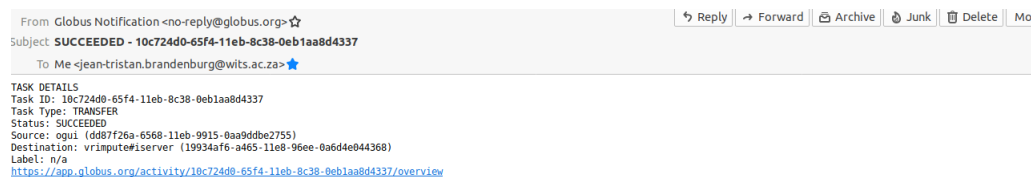 ISSN: 1934-8266. DOI: 10.1002/0471142905.hg0119s68. URL: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3066182/ (visited on 07/06/2022).

Figure 5: Emel of globus validate download



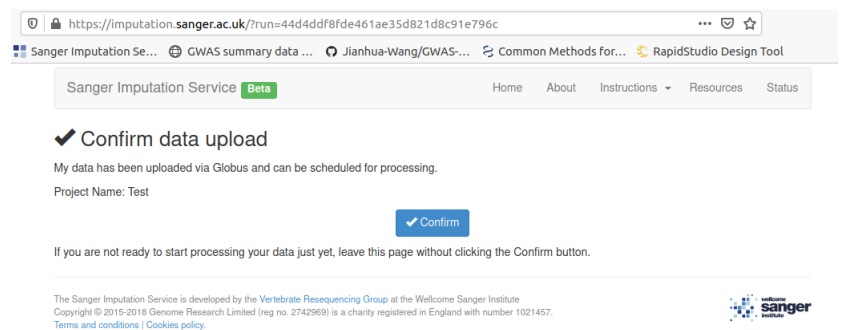Figure 6: Emel to validate well transfert to server



Figure 7: Validate your data after link emel

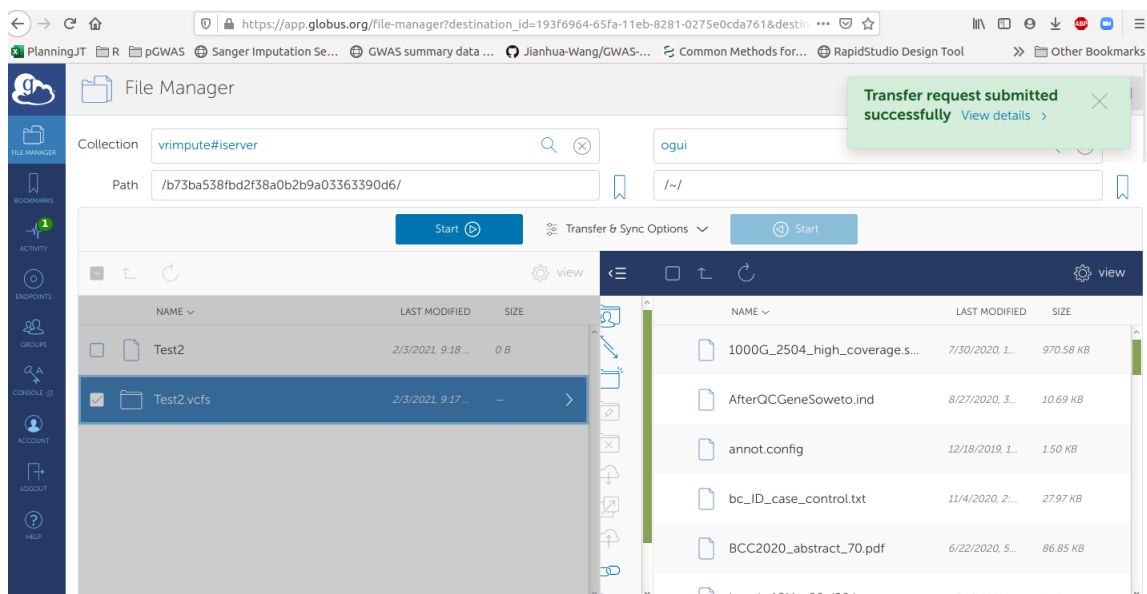Figure 8: Mel validate that you data is ready on globus
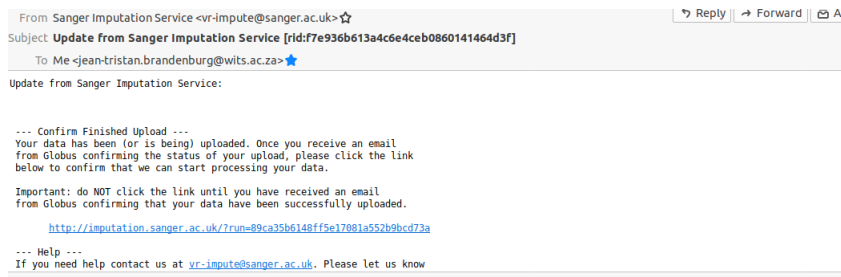


Figure 9: Transfer data from server to cluster



Figure 10: Validate your data after link emel