# Genome-wide association study

# Data formats in PLINK

Segun Fatumo PhD

UNIVERSITY OF CAMBRIDGE

wellcome trust sanger institute

# Data formats in PLINK

# PLINK Flat files (MAP/PED)

- PLINK is a very widely used application for analyzing genotypic data. It can be considered the "de-facto" standard of the field, although newer formats are starting to be widespread as well.

# MAP files

The **MAP** file describes the SNPs.

The fields in a MAP file are:
- Chromosome
- Marker ID
- Genetic distance
- Physical position

# MAP files

- Chromosome number *[integer]*
- SNP ID *[string]*
- SNP genetic position (cM) *[float]*
- SNP physical position (bp) *[integer]*

- This file should have L lines and 4 columns, where L is the number of SNPs contained in the dataset.

- Each SNP must have a unique physical position. All the SNPs must be ordered by physical position.

# Quick Exercise

- Download c_2150SNPs.map & c_2150SNPs.ped

- View this data

- how many individuals are in the datasets? How
- many SNPs? What is genotyping rate?

# MAP file - Example

| Example of a MAP file of the standard PLINK format: | | | |
|---|---|---|---|
| 21 | rs11511647 | 0 | 26765 |
| X | rs3883674 | 0 | 32380 |
| X | rs12218882 | 0 | 48172 |
| 9 | rs10904045 | 0 | 48426 |
| 9 | rs10751931 | 0 | 49949 |
| 8 | rs11252127 | 0 | 52087 |
| 10 | rs12775203 | 0 | 52277 |
| 8 | rs12255619 | 0 | 52481 |

# Ped (Pedigree)

- The **PED** file describes the individuals and the genetic data. The PED file corresponding to the example dataset is:

# Ped (Pedigree)

- This file can be *SPACE* or *TAB* delimited. Each line corresponds to a single individual. The first 6 columns are:
- **Family ID** *[string]*
- **Individual ID** *[string]*
- **Father ID** *[string]*
- **Mother ID** *[string]*
- **Sex** *[integer]*
- **Phenotype** *[float]*
- Columns 7 & 8 code for the observed alleles at SNP1, columns 9 & 10 code for the observed alleles at SNP2, and so on. Missing data are coded as "0 0" as for example for SNP3 of IND1. This file should have N lines and 2L+6 columns, where N and L are the numbers of individuals and SNPs contained in the dataset respectively.

# Ped (Pedigree)

- Pedigree Name: A unique alphanumeric identifier for this individual's family. Unrelated individuals should not share a pedigree name.

- Individual ID: An alphanumeric identifier for this individual. Should be unique within his family (see above).

- Father's ID : Identifier corresponding to father's individual ID or "0" if unknown father. Note that if a father ID is specified, the father must also appear in the file.

- Mother's IDIdentifier corresponding to mother's individual ID or "0" if unknown mother Note that if a mother ID is specified, the mother must also appear in the file.

- SexIndividual's gender (1=MALE, 2=FEMALE).

- Affection status

# PED file - Example

| | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FAM1 | NA06985 | 0 | 0 | 1 | 1 | A | T | T | T | G | G | C | C | A | T | T | T | G | G | C | C |
| FAM1 | NA06991 | 0 | 0 | 1 | 1 | C | T | T | T | G | G | C | C | C | T | T | T | G | G | C | C |
| 0 | NA06993 | 0 | 0 | 1 | 1 | C | T | T | T | G | G | C | T | C | T | T | T | G | G | C | T |
| 0 | NA06994 | 0 | 0 | 1 | 1 | C | T | T | T | G | G | C | C | C | T | T | T | G | G | C | C |
| 0 | NA07000 | 0 | 0 | 2 | 1 | C | T | T | T | G | G | C | T | C | T | T | T | G | G | C | T |
| 0 | NA07019 | 0 | 0 | 1 | 1 | C | T | T | T | G | G | C | C | C | T | T | T | G | G | C | C |
| 0 | NA07022 | 0 | 0 | 2 | 1 | C | T | T | T | G | G | 0 | 0 | C | T | T | T | G | G | 0 | 0 |
| 0 | NA07029 | 0 | 0 | 1 | 1 | C | T | T | T | G | G | C | C | C | T | T | T | G | G | C | C |
| FAM2 | NA07056 | 0 | 0 | 0 | 2 | C | T | T | T | A | G | C | T | C | T | T | T | A | G | C | T |
| FAM2 | NA07345 | 0 | 0 | 1 | 1 | C | T | T | T | G | G | C | C | C | T | T | T | G | G | C | C |

# PLINK Binary files (BED/BIM/FAM)

- The binary PLINK format contains the same information as the flat file PLINK format but in a compressed and significantly more efficient form.

# BED File

- The BED files are encoded in binary format.
- Binary PED (BED) files.

- We can inspect the BED file with the Unix xxd command, to view a binary file

- ***xxd -b test.bed***
- which generates:

0000000: 01101100 00011011 00000001 11011100 00001111 11100111 l..... 0000006: 00001111 01101011 00000001

# BIM files

The fields in a BIM file are:

- Chromosome
- Marker ID
- Genetic distance
- Physical position
- Allele 1
- Allele 2

***How has BIM file different from MAP file ?***

# BIM File

**Example of a BIM file of the binary PLINK format:**

- 21  rs115116470    26765        A    T
- X    rs3883674  0    32380        C    G
- X    rs122188820    48172        T    T
- 9    rs109040450    48426        A    T
- 9    rs107519310    49949        C    T
- 8    rs112521270    52087        A    C
- 10  rs127752030    52277        A    A
- 8    rs122556190    52481        G    T

# FAM files

- **The fields in a FAM file are**

- Family ID

- Sample ID

- Paternal ID

- Maternal ID

- Sex (1=male; 2=female; other=unknown)

- Affection (0=unknown; 1=unaffected; 2=affected)

- Note: The **FAM file** is just the first six columns of the **PED file**

# FAM file

| | | | | | |
|---|---|---|---|---|---|
| FAM1 | NA06985 | 0 | 0 | 1 | 1 |
| FAM1 | NA06991 | 0 | 0 | 1 | 1 |
| 0 | NA06993 | 0 | 0 | 1 | 1 |
| 0 | NA06994 | 0 | 0 | 1 | 1 |
| 0 | NA07000 | 0 | 0 | 2 | 1 |
| 0 | NA07019 | 0 | 0 | 1 | 1 |
| 0 | NA07022 | 0 | 0 | 2 | 1 |
| 0 | NA07029 | 0 | 0 | 1 | 1 |
| FAM2 | NA07056 | 0 | 0 | 0 | 2 |
| FAM2 | NA07345 | 0 | 0 | 1 | 1 |

# PED/MAP to BED/BIM/FAM conversion

- To convert *myPlinkTextData.ped* and *myPlinkTextData.map* in Plink binary format, use Plink as follows:

- plink --file myPlinkTextData --make-bed --out myPlinkBinaryData

gPLINK

PLINK

Haploview

GUI to initiate PLINK jobs

C/C++ analysis engine (can run standalone)

Initiate PLINK jobs locally or remotely

Track PLINK jobs and results

Plot PLINK WGAS results

Visualize LD patterns

Job tracking interface

Integrate with Haploview

Tabulate, filter PLINK WGAS results

Visualize PLINK results (population stratification)

# GUI for many **PLINK** analyses
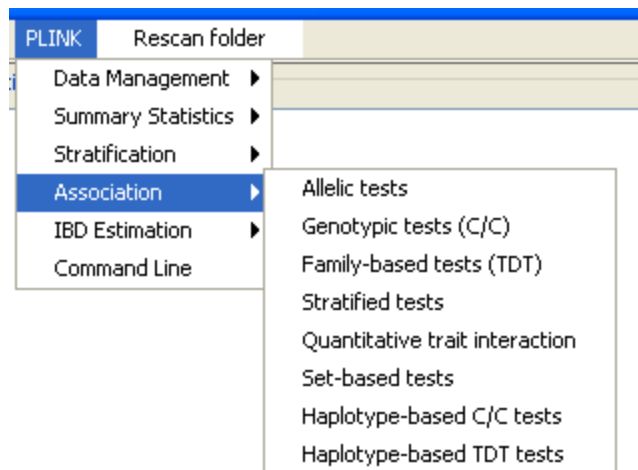
**Data management**
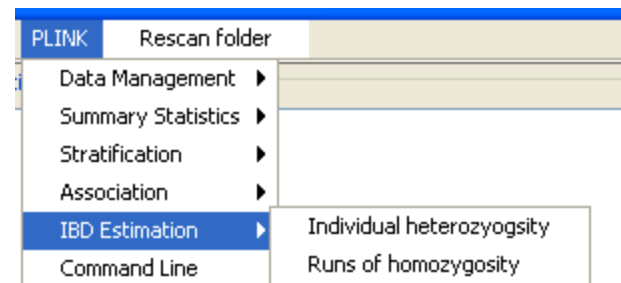


**Summary statistics**



**Population stratification**



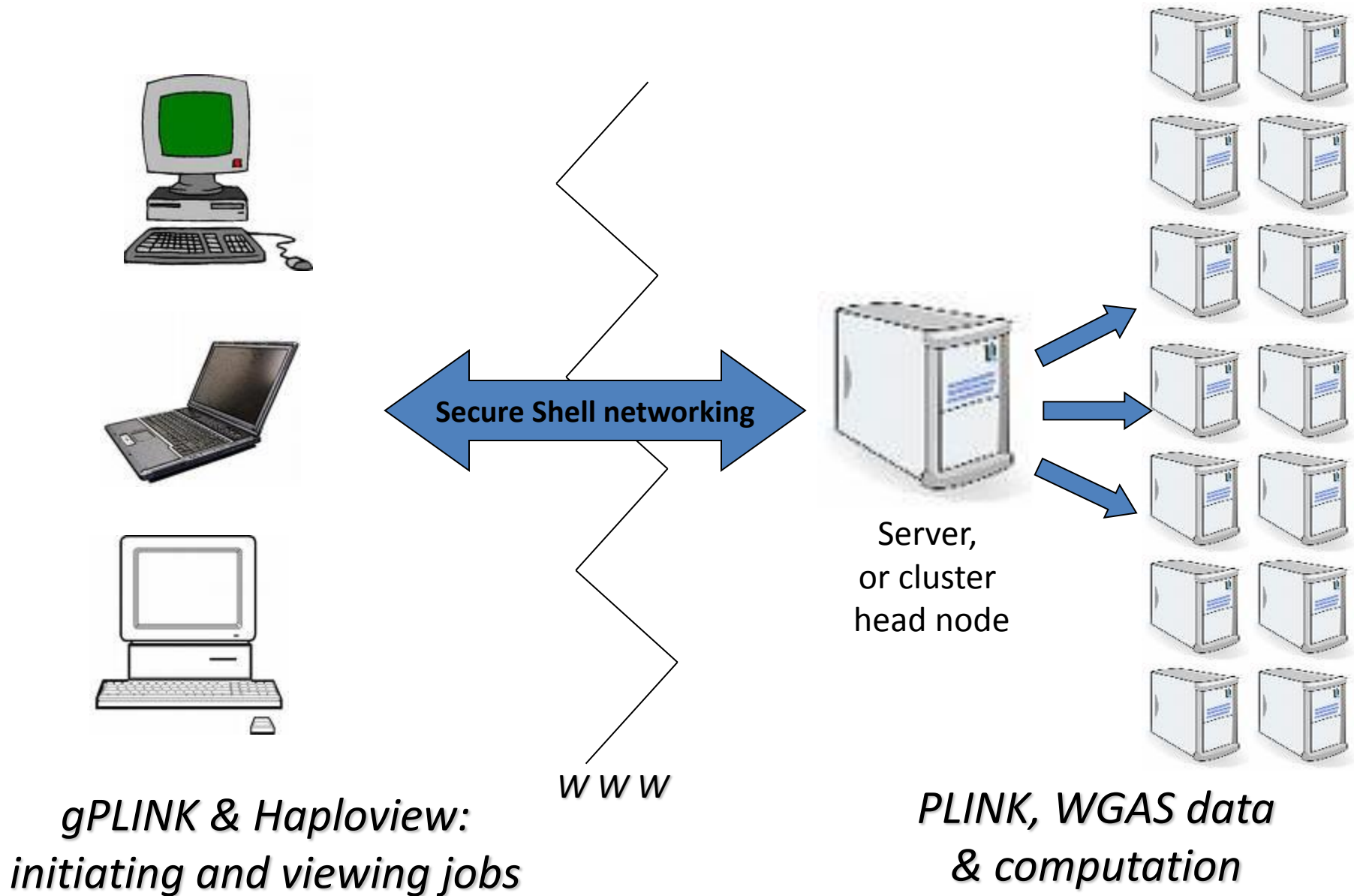**Association analysis**



**IBD-based analysis**

# Computational efficiency

*350 individuals genotyped on 100,000 SNPs*

| Load, filter and analyze | ~12 seconds |
|---|---|
| 1 permutation (all SNPs) | ~1.6 seconds |

# gPLINK / PLINK in "remote mode"



**Secure Shell networking**

*w w w*

Server,
or cluster
head node

*gPLINK & Haploview:
initiating and viewing jobs*

*PLINK, WGAS data
& computation*

A simulated WGAS dataset

Summary statistics and quality control

Whole genome SNP-based association

Whole genome haplotype-based association

Assessment of population stratification

Further exploration of 'hits'

Visualization and follow-up using Haploview

# Exercise

- Create a simulated class.ped and class.med files
- Convert your ped & map file to binary file

- **Assumption:**

-You may assume that there are 20 students in this class.

-You may assume also that two students are brothers

# Manipulating the data files

- Get only the genotypes for a single chromosome or a region around a
- SNP
- --chr 13
- Exercise: Get data from chromosome 13 and write to a new BED file. If
- you are having trouble running the full dataset, you can use this fileset

instead of fulldataset.
- plink --bfile bipolar --chr 13 --make-bed --out bipolar_chr13