



Centre for  
Tropical Livestock  
Genetics and Health

# Genomic data management

Getinet M. Tarekegn, PhD



ROSLIN

SRUC

ILRI

CGIAR

# Who am I?



Centre for  
Tropical Livestock  
Genetics and Health

## Academic background:

- BSc in Animal Sciences .....Haramaya University
- MSc in Animal Genetics and Breeding .....Haramaya University
- BSc in Computer Sciences .....Bahir Dar University
- Postgraduate studies in Higher Education.....Leeds Met. University, UK
- Visiting researcher in Small Ruminant Genomics ...Inner Mongolian Agri.University, China
- PhD in Applied Genetics (Livestock Genomics).... AAU, 2016
- Post-doc/Visiting researcher in African Goat Genomics ... BecA-ILRI, 2016/2017
- Post-doc in Dairy cattle genetics and genomics..... SLU, Sweden, 2017-2020
- Geneticist and Bioinformatician, Scotland's Rural College (SRUC)... 2022 and onwards

## Affiliations:

- ✓ Geneticist and Bioinformatician, Scotland's Rural College (SRUC),  
King's Building, University of Edinburgh, UK
- ✓ Adjunct Prof of Livestock Genomics and Bioinformatics, Addis Ababa University, Ethiopia

<https://scholar.google.com/citations?user=yOA-j4gAAAAJ&hl=en&oi=ao>



# Edinburgh Genetic Evaluation Service (EGENES)



Centre for  
Tropical Livestock  
Genetics and Health

- EGENES is a leading centre for the development and delivery of genetic improvement tools for the livestock industry.
- The team at EGENES produces national genetic and genomic evaluations for **all dairy cattle and sheep** and for the **UK's biggest beef breeds**.
- uses performance and pedigree data recorded by farmers, breeders and other industry players.
- These data are combined, quality controlled and analysed to produce routine genetic evaluations, which are then **fed back to industry**.





# Edinburgh Genetic Evaluation Services

A centre for the development and delivery of genetic improvement tools for our livestock industries.

- 5 Geneticist and 4 programmers

## Focus areas:

- Dairy breeds
- Major Beef breeds
- Five sheep breed
- Carcass evaluation
- National evaluations
- Interbull
- Interbeef
- M3GE project



## Outline

- SNP genotypes
- Quality control
- Converting the snps data to different formats
- Plink practicals

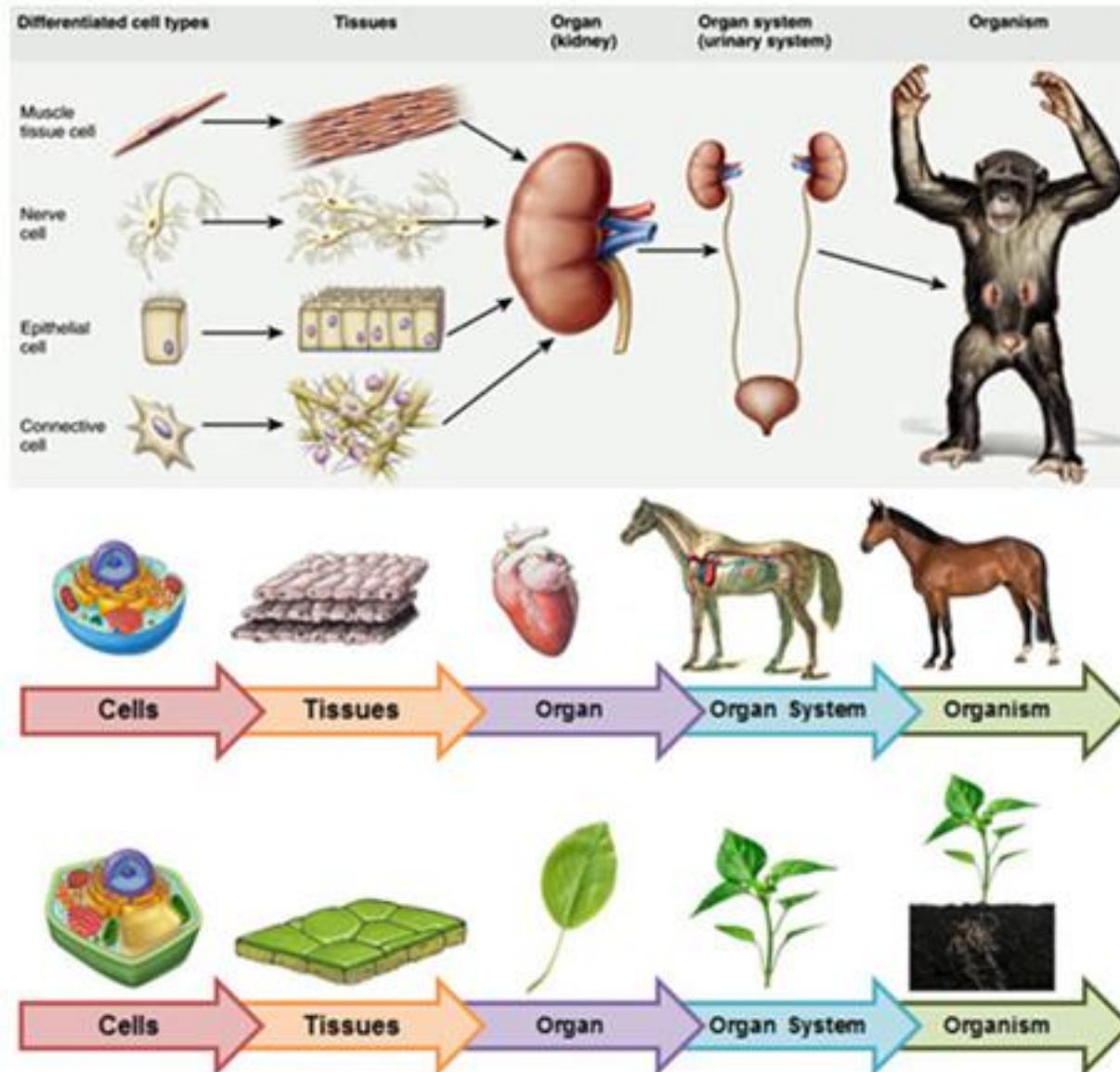
### Learning objectives:

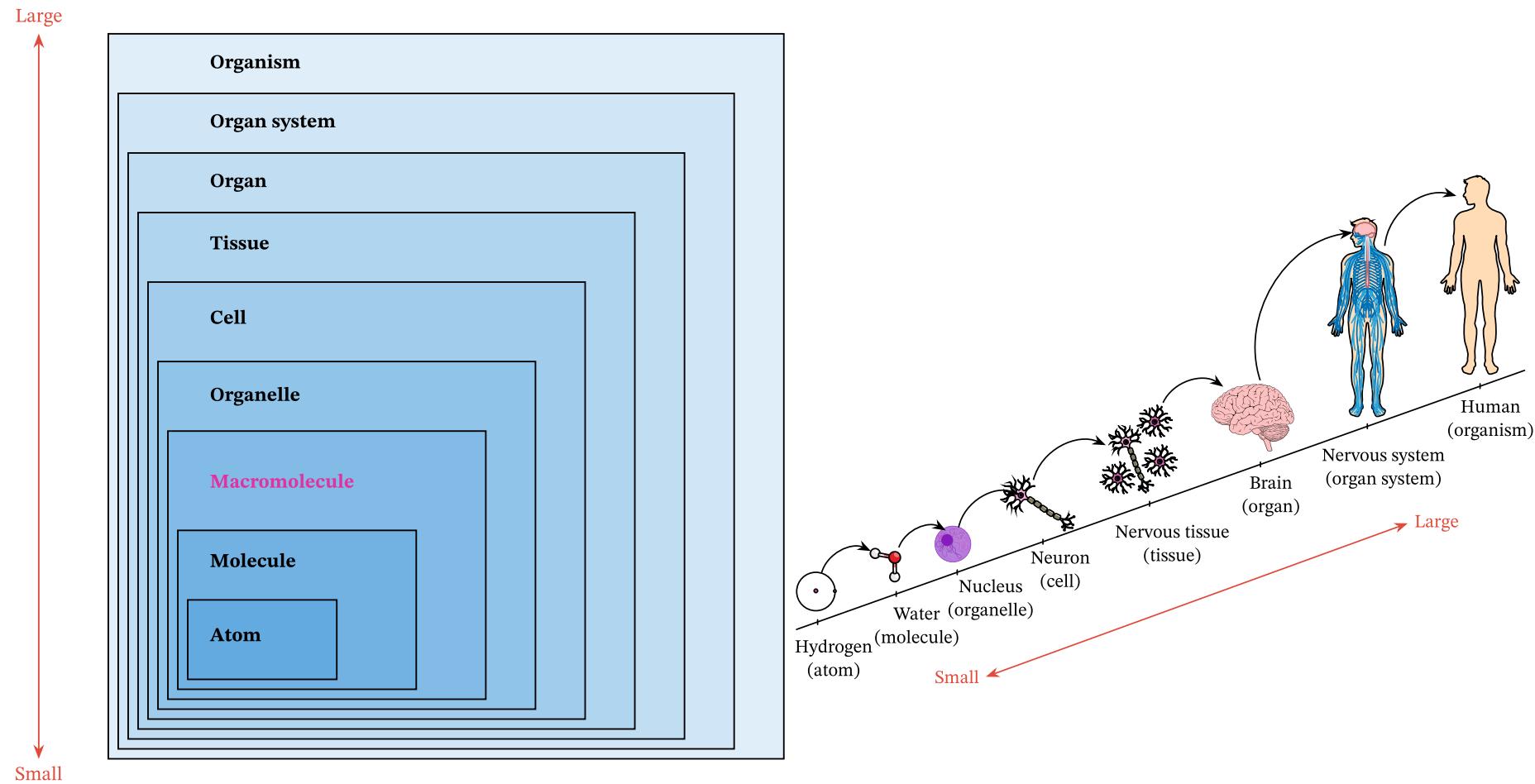
- Understand how to filter genotype data and converting to different formats

### Learning outcome:

- All trainees are capable of handling genotype data

# Biological levels of Organization of living things



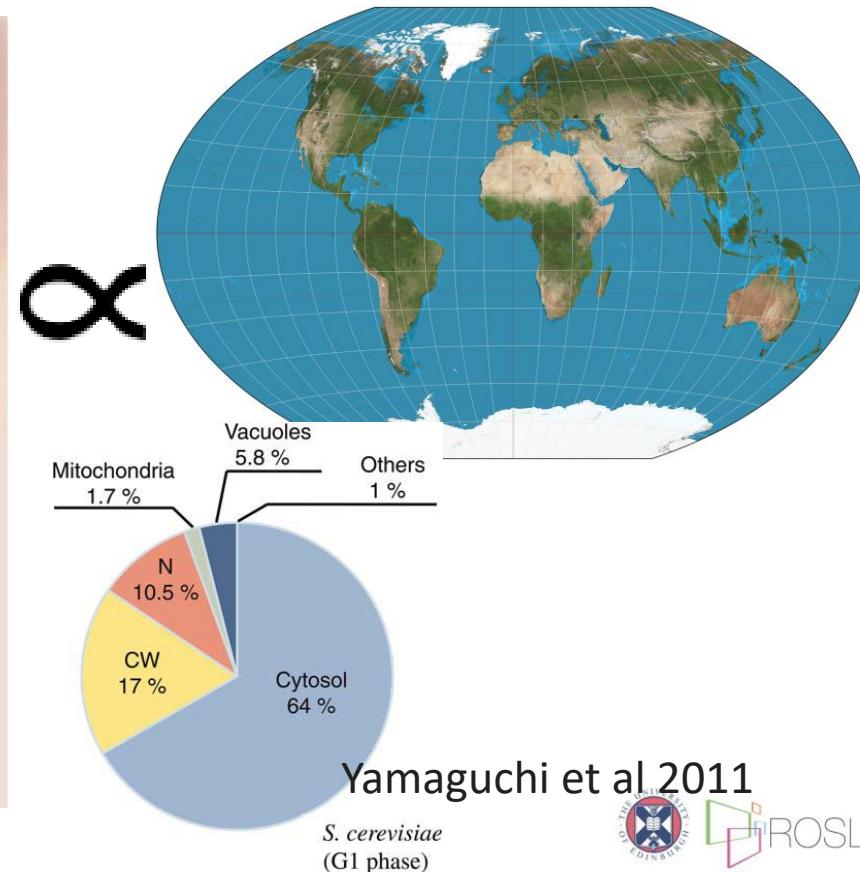
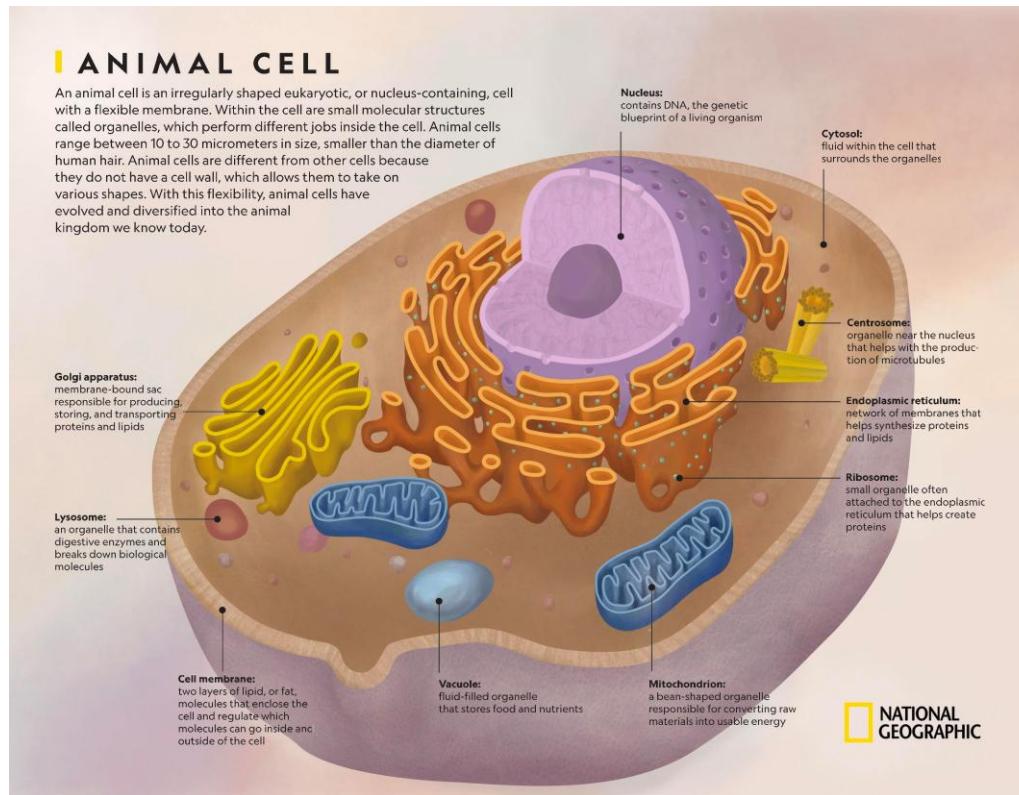


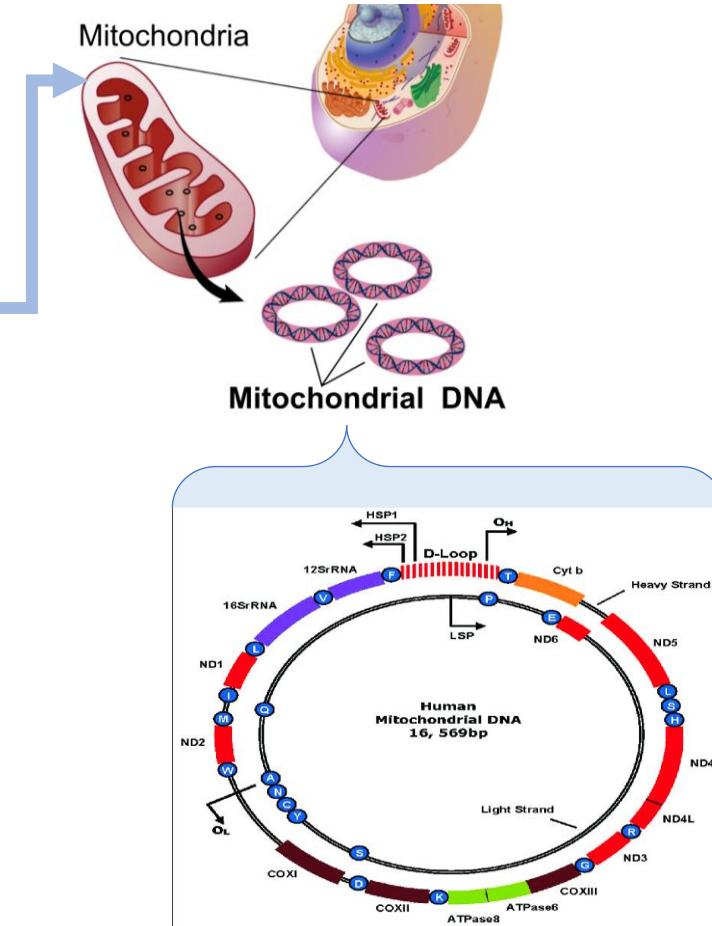
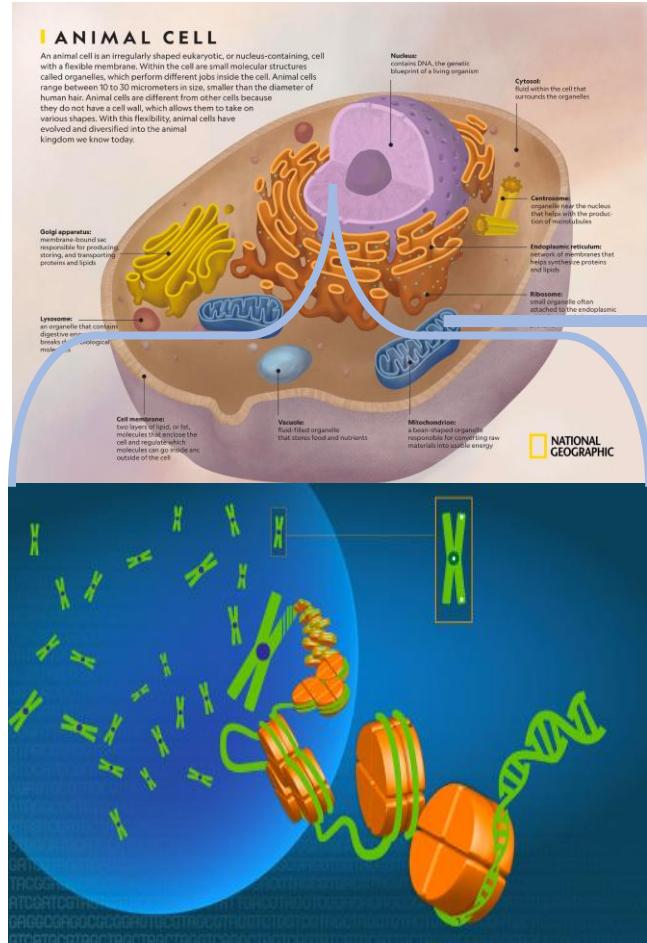
A diagram showing the biological levels of organization from an atom to a multicellular organism <https://www.nagwa.com/en/explainers/430187521519/>



Cell is the smallest unit of life.

- can divide, multiply, grow and respond to stimuli from the environment; but is it as smallest in size as we expect?



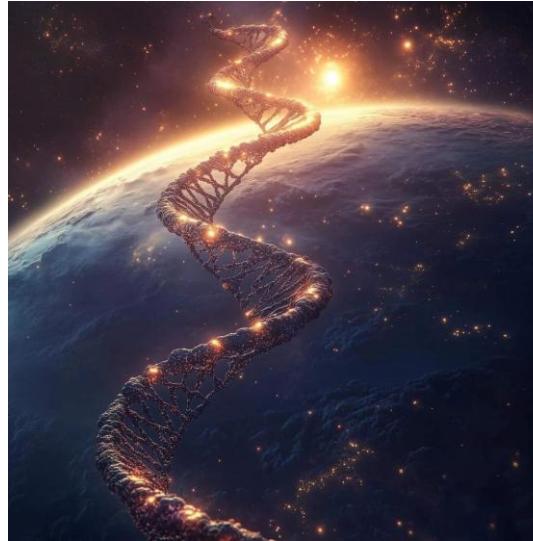


- 5 cm long (about 2 inches) each chr, and all 46 chrs be about 2 m

# Interesting Facts about DNA

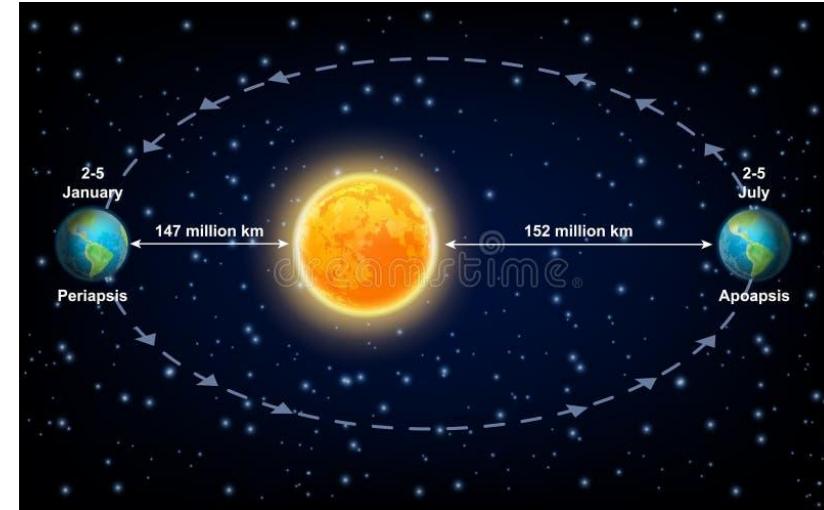
- Your DNA Could Go From Earth to the Sun 600 Times!  
Your DNA is incredibly long - if stretched out, the DNA in your body could reach from the Earth to the Sun and back over 600 times.
- Each human cell contains approximately 6 feet of DNA, which is compacted into a structure called chromatin to fit inside the nucleus.
- If all the DNA in your body was uncoiled, it would stretch 67 billion miles long - equivalent to about 150,000 round trips to the Moon.

<https://www.ancestry.com/c/dna-learning-hub/dna-facts#:~:text=1.,stretch%20over%2067%20billion%20miles.>



A human cell contains 46 chromosomes:  
The human body contains an estimated 37.2 trillion cells,  
with the number varying by sex and age:

- **Males:** Approximately 36 trillion cells
- **Females:** Approximately 28 trillion cells
- **Children:** Approximately 17 trillion cells
- **Newborns:** Approximately 26 billion cells



## ... starts at the field



Sampling



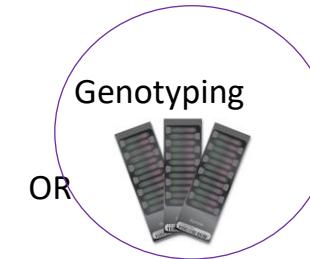
Extracting the DNA



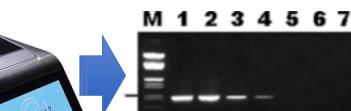
Quality check



WGS



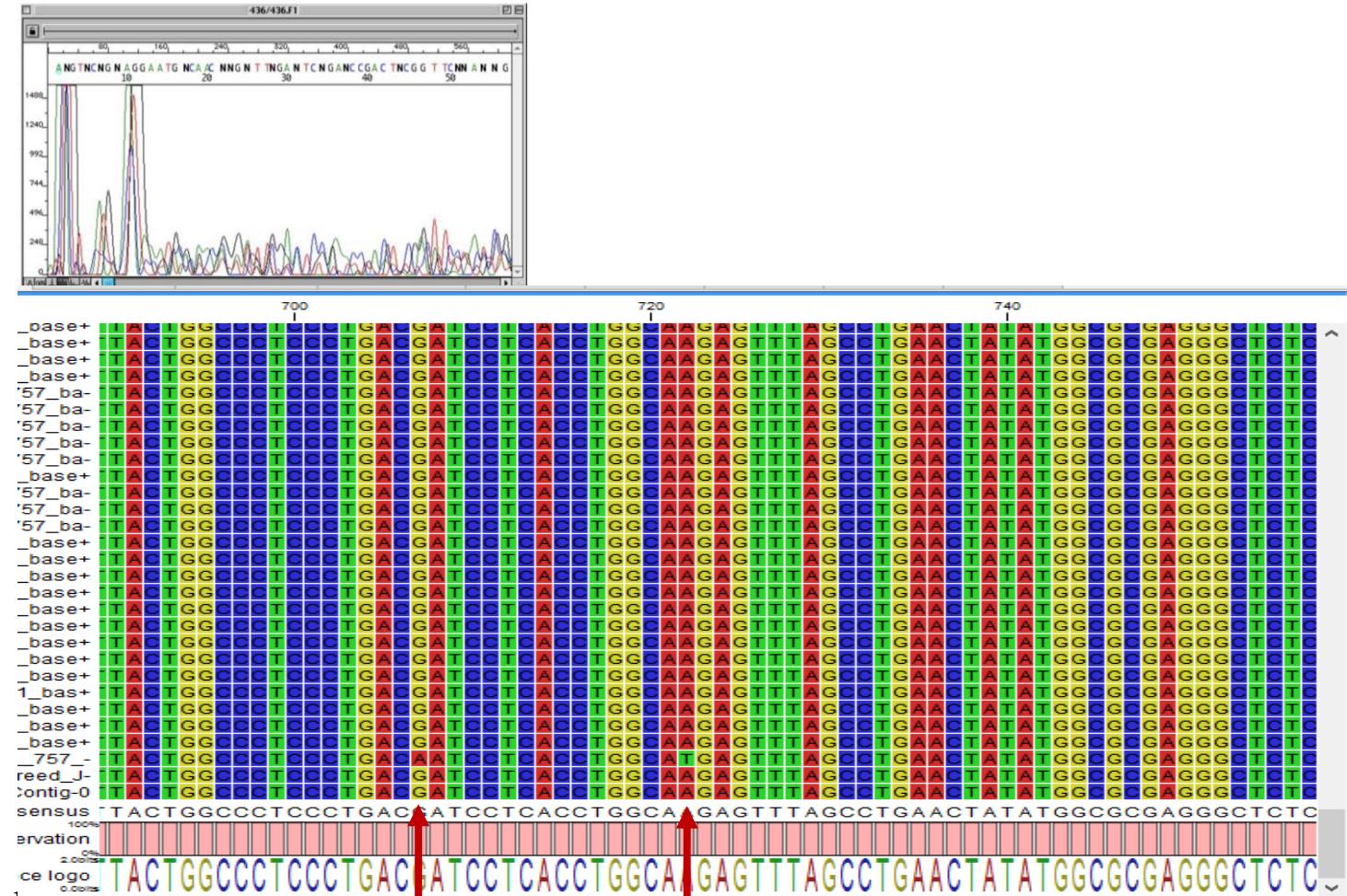
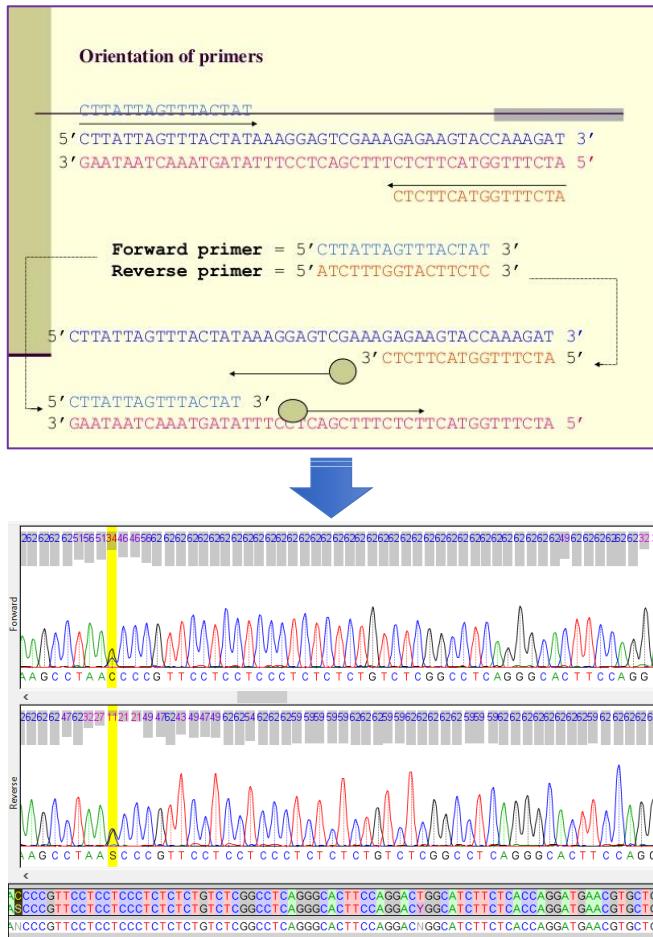
PCR



# Sanger sequences



Centre for  
Tropical Livestock  
Genetics and Health



- A chromatogram (sometimes also called electropherogram) is the visual representation of a DNA sample produced by a sequencing machine

2 variants



## WGS/NGS data

---

### fastq format

```
@A00291:9:H5N3MDMX:1:1101:1181:1094 1:N:0:ATGCCTAA
GNTGGCTTGGGGTTTGGAATCGTGATACCAGAGGATGCCTACGAAAGAGTTAAATAC
+
F#FFFF:FFFF:::FFF:::FFFFFF:FFFFFF:FFFFFF:FFFFFF:FFFFFF:FFFFFF
```

Each read is represented by four lines. These lines are:

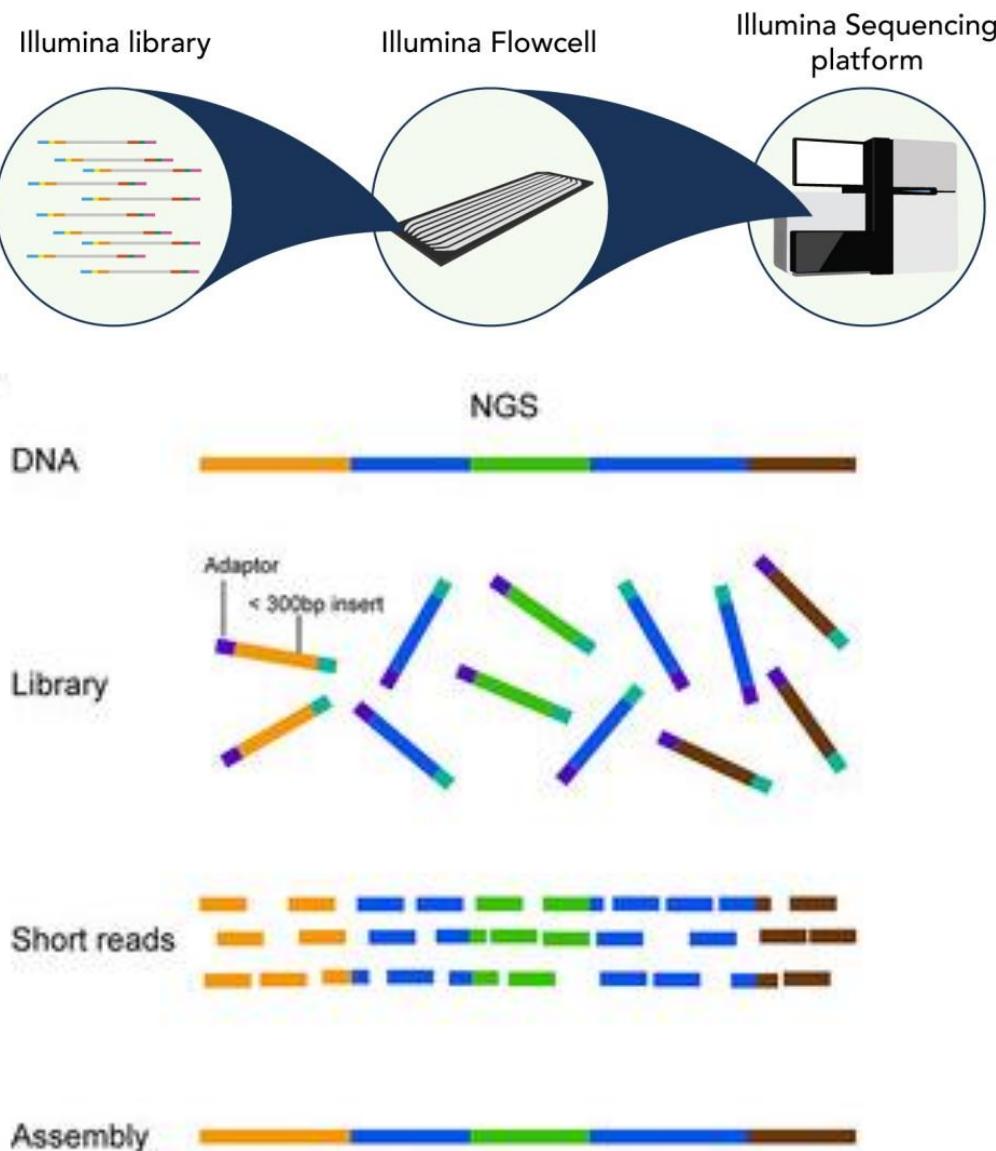
1. A sequence identifier with information about the sequencing run and the cluster.
2. The sequence (the base calls; A, C, T, G and N).
3. A separator, which is simply a plus (+) sign.
4. The base call quality scores.

[https://knowledge.illumina.com/software/general/software-general-reference\\_material-list/000002211](https://knowledge.illumina.com/software/general/software-general-reference_material-list/000002211)

# WGS data handling pipeline



Centre for  
Tropical Livestock  
Genetics and Health

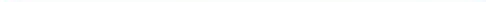
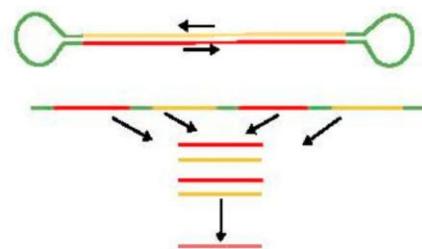


<https://www.nature.com/articles/s41597-024-03342-9>



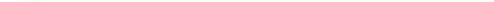
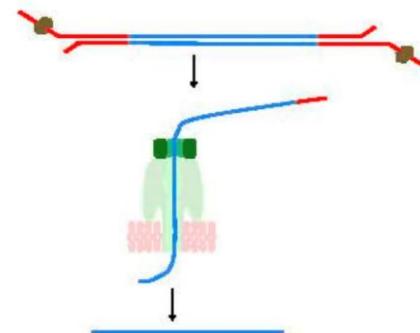


### PacBio



- Long-read lengths

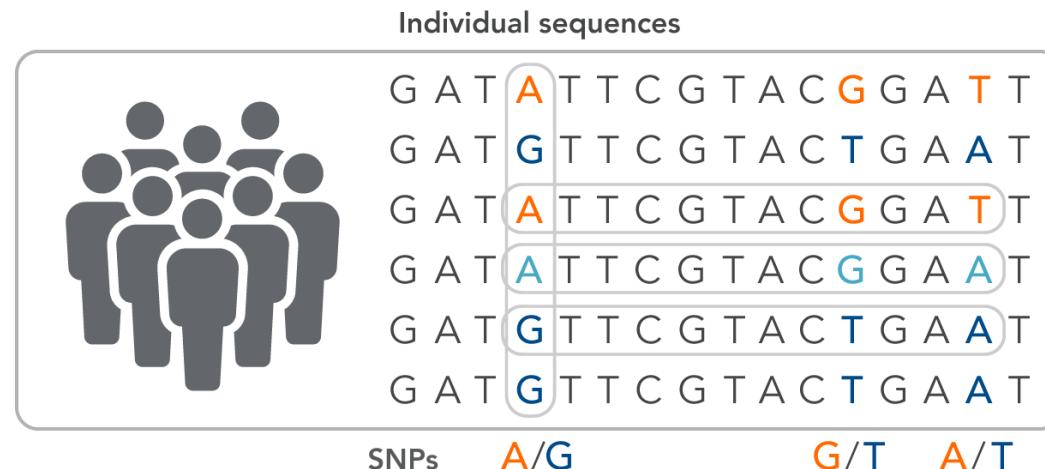
### Oxford Nanopore



- Ultra-long read lengths

## Single-nucleotide polymorphisms (SNP) genotypes

- A single nucleotide polymorphism (SNP) is a variation at a single position in a DNA sequence among individuals.
- SNPs:
  - the most abundant genetic markers
  - have been widely used in **genomic** research
    - disease gene mapping
    - Medical and clinical diagnostics
    - forensic tests
    - genetic diversity
    - GWAS
    - Prediction



<https://eu.idtdna.com/pages/education/decoded/article/genotyping-terms-to-know>



# Are all mutations SNPs?

**articles**

## Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium\*

\* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.

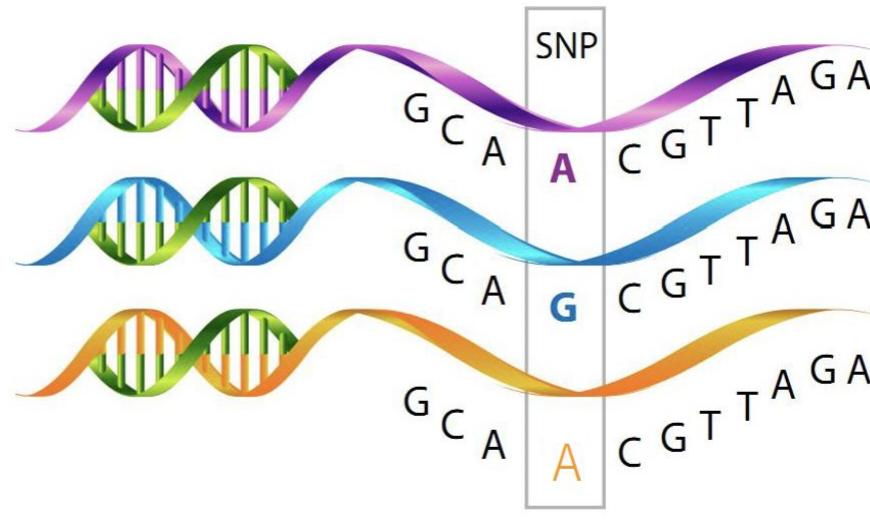
The rediscovery of Mendel's laws of heredity in the opening weeks of the 20th century<sup>1,2</sup> sparked a scientific quest to understand the nature and content of genetic information that has propelled biology for the last hundred years. The scientific progress made falls naturally into four main phases, corresponding roughly to the first, second, third and fourth centuries of the history of heredity: the chromosomes; the DNA double helix; the third unlocked the informational basis of heredity, with the discovery of the biological mechanism by which genes are turned on and off; and finally the genome. The last quarter of a century has been marked by a relentless drive to decipher first genes and then entire genomes, spawning the field of genomics. The fruits of this work already include the 3 sequences of over viruses and viroids, 205 naturally or plantains, 185 bacterial, 31 eukaryotic, seven archaeal genomes, two animals and one plant.

Here we report the results of a collaboration involving 20 countries, the United States, the United Kingdom, Japan, Germany and China to produce a draft sequence of the genome. The draft sequence is now available online<sup>3</sup> and covers 91% of the genome, with a map covering more than 99% of the euchromatic part of the genome and, together with additional sequence in public databases, it covers about 94% of the human genome. The sequence was produced over a relatively short period, with coverage rising from about 10% to more than 90% over roughly fifteen months since the start of the project. The sequence is being updated daily throughout the project.<sup>4</sup> The task ahead is to finish the sequence by closing all gaps and resolving all ambiguities. Already about one billion bases are in final form and the bringing the vast majority of the sequence to this standard is straightforward and should proceed rapidly.

The analysis of the genome is of interest in its own right. It is the largest genome to be extensively sequenced, being 25 times as large as any previously sequenced genome, eight times as large as the sum of all such genomes. It is the first vertebrate genome to be extensively sequenced. And, uniquely, the genome is now complete.

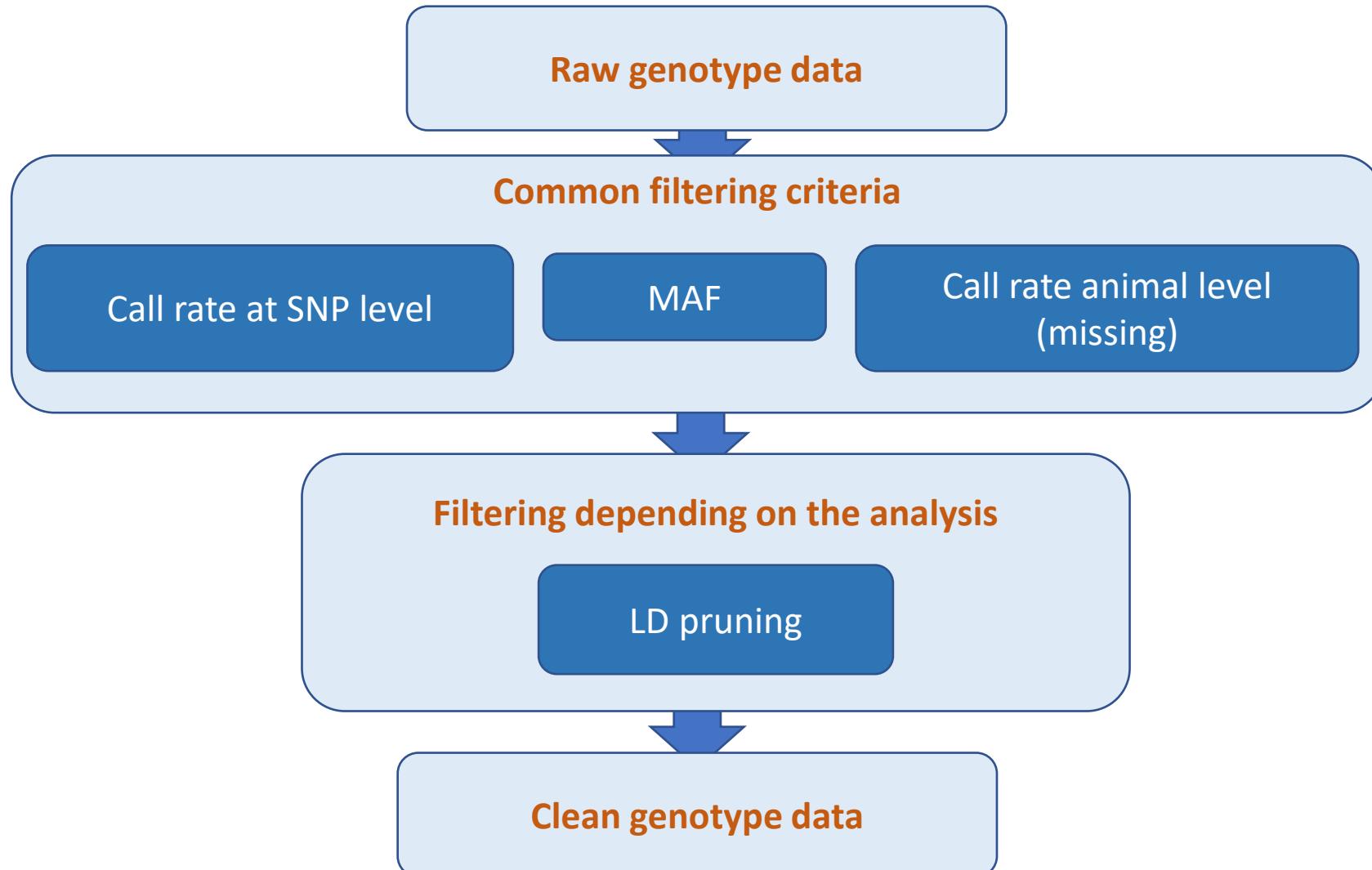
Much work remains to be done to produce a complete genome sequence, but the vast trove of information that has become available through this collaborative effort allows a global perspective on the human genome. Although the details will change

nature  
the human genome



- To say a SNP, the available polymorphism needs to be observed at appreciable frequency (traditionally, **at least 1%**) in the human population for instance.

# Quality control parameters



**1. Call rate at SNP level:** the call rate for a given SNP is **the proportion of individuals** in the study for which the corresponding SNP information is not missing.

- E.g. using a call rate of 95%, meaning we retain SNPs for which there is less than 5% missing data.

	SNP	ANIMAL
025	110101111051110111110010001221151205122115021511110150122010201021000221121015000122010	00011112011012110
036	211011010220121222201210122201012022111112021222111112102020101101020111112011012110	0001112020001112020
050	12101002111202111120002121222100021122122122110000020220000211022122212122020001112020:	
054	1200012002201212111001210022211021221102021121222120022002121212111120212022022022100:	
066	20000202022102121211200220012221110122020211020222020220001222121011201021022010011010:	
097	10110212102211221110210011111001022112120221111102022100120122012111021021021012000:	
101	1210021202200112211100011112220100101120112121211121201221002102002021211222022010022110:	
151	111001020221220210201010122020012122111122122112111222002201112012111022000022012:	
172	211012020211112101211021102201010012212122211022202012210202121120102111220221120110:	
224	220001110221102211010211025202011121022212221201101210110220050210121022010022125:	
277	2101022001212212112120210122200201221021211020112102122100211011021112021012010:	
314	122011120122202102100100212100112012020200121002021210012010202121211102201010100:	
419	221112210121120222210221021102010211212111220000001112200022111220222112120012121110:	
439	20020120102121210101021022110112220202211001011121001201021022011021010011020:	
456	12000102221112200101021002022121222200101022111021021211221111102112010:	
501	11100021221121201212100221101202221101022112221102200112021100202101022010021020:	
571	11000012022022121020021012000111221101022221102200220120210010102121022102010110:	
579	11210021021001010111102002212211112020222211022210102012011122211111201011020:	
581	21100202021001221202011002200201125121215025222225022101120112105120222112111012110:	
657	1100111202201121111010200122210001122212102121112120022001122022200221221211200112011:	
660	2100021202211202211210210122210110122212211121202012210122011211121112022000012101:	
730	2100020202200202222001200222000122022220021102252200122001202111151001012022001012025:	
732	212102121011002201200012101121201215110215122521121150220011102111050202221122011022010:	
764	111102121520012212110200012100220201225222115021522211502201102021200502020220111112110:	
780	121101021122202101010220022120120112122101211110111221020202001010112212121002021021:	
800	22100012022122210202021102221101012112022120222220022100221112102102102201102201011010:	
816	11000122022012120201102201112110010121122121220021122200222111021111212022011022010:	
832	12101001120011211110021112220111121222121020111102022100121211121012111110:	
900	2101001102201221211021102121012212121101111102210012021211102110102110022020:	
901	12100102022112121210010002120201111221112122001111102210022010220122121021000012020:	

## 2. Sample-level filtering (call rate):

- Individuals who have missing genotype data across more than a pre-defined percentage of the typed SNPs need to be excluded.
- This proportion of missingness across SNPs is referred to as the **sample call rate**, and we apply a threshold of 95%.
- That is, individuals who are missing genotype data for more than 5% of the typed SNPs are removed.

	SNP	ANIMAL
025	1101011101S1110111110010001221151205122115225111101012201020102100022112101500122010	
036	211011101022012122220121012220101202221111120212221111121020201011010201111112011012110:	
050	12101002111202111120002121222100021122122122110000020220000211022122212122020001112020:	
054	12000120022012121110012100222211021122110201121222120022002121212111202112022002022100:	
066	2000020220210212211200220012222111012202011020222020220001222121011201021022010011010:	
097	101102120220121122111021001111100102211212022111110202210012012201211102102101200:	
101	121100212022001122110001111220100101120112121211121201221002102002021211222022010022110:	
151	111001020221220210201101220201212211112212111122002211120112111020220002201211102022012:	
172	21101202201111210121102110220101021121202211221102022012112010211122022112011010:	
224	22000111022101221010102110252020111212022212221101121011102220050210121022010022125:	
277	210102200121221211212021012220020122102121102011210212210022110110202111201021010:	
314	122011120122202102010021210011201202002102002121001120102202121211022010101100:	
419	22111221012112022221022102110201021121112200000011122002211122020222112120012121110:	
439	2002021010212121010102112211010112202020221100101112100112010220122201021011020:	
456	12000102022111220010102002211000202211212222001010221110210210112122111110211010:	
501	11100002121121201212121002221101202222101022112221102200112021100201102022100021020:	
571	110000120220022121202200121002001112211011022221200220021200101021212022101010:	
579	1121002102100101011102200221200022112020222211022210120212112221111102101010:	
581	2110020210101221202101002200201125121502252222502210112011210512022211211012110:	
657	110011120220111211101020012210001122121202112112002200222002212211212001112011:	
660	2100021202211202211210102211010222122212121112021221012011211112111202200012101:	
730	21000202200202222001200222001220222200211022522001220012011151001012022001012025:	
732	212102121510022012000121011212012151102151225211150220011102111050202211220102010:	
764	111102120100122122110200012202012252221150220152222115022011020212005020202202211112110:	
780	121101021122202101010220022120120112122101211110111221020202001011221212002021021:	
800	221000120221222210202021102221101012112022102222200221022211121021202011022010111010:	
816	1100012202201212201100220111211000110211212120020112220022211102111121022011111010:	
832	1210100111200112111100211122011112122212102011102022100211221001212111210211111010:	
900	2101001102201221212110211021210221212110111110221001202121110211021100220:	
901	121001020221121212210010021202011112211121220011111102210022010220122121201000120:	

**3. MAF (SNP-level filtering):** A large degree of homogeneity at a given SNP across study participants generally results in inadequate power to infer a statistically significant relationship between the SNP and the trait under study.

- remove SNPs for which the MAF is less than 1%.
  - In some instances, particularly small sample settings, a cut off point 5% is applied.
    - MAF is the lowest of the two allele frequencies
    - $p = freq(A)$
    - $q = 1 - p = freq(B)$
    - $MAF = \min(p, q)$
    - A fixed marker ( $p = 0$  or  $p = 1$ ) gives no information
    - An almost-fixed marker ( $p = 0.0001$  or  $p = 0.9999$ ) gives almost no info



#### 4. LD pruning:

- Linkage disequilibrium (LD) is **the nonrandom association of alleles of different loci (Slatkin, 2008)**.
- **LD pruning** is the removing loci based on high levels of pairwise LD.
- Better results of population structure and Principal Component Analysis (PCA) are assumed to be obtained if the markers used are not in linkage disequilibrium with each other.
- If any pair of markers within the window are in LD greater than the specified threshold, the first marker in the pair will be pruned.



Centre for  
Tropical Livestock  
Genetics and Health

---

# SNP genotype file format, and how can we manage the data?



# Formats of the input file

0	snp6964-scaffold12561-13292	0	0
0	snp7949-scaffold12878-8765	0	0
0	snp7961-scaffold12886-28972	0	0
0	snp7962-scaffold1289-6273	0	0
0	snp8461-scaffold1305-10877	0	0
0	snp8463-scaffold1305-74425	0	0
0	snp8464-scaffold1305-146599	0	0
0	snp8466-scaffold1306-22563	0	0
0	snp8716-scaffold1310-23370	0	0
0	snp9138-scaffold1331-33874	0	0
0	snp9158-scaffold1333-23917	0	0
0	snp9659-scaffold13458-6636	0	0
0	snp9687-scaffold1348-7500	0	0
1	snp19065-scaffold1917-222828	0	0
1	snp4303-scaffold1134-142270	0	0
1	snp14078-scaffold1560-21647	0	21647
	Map file		
	snp14079-scaffold1560-51801	0	51801
1	snp14080-scaffold1560-100946	0	100946
1	snp14082-scaffold1560-185575	0	185575
1	snp14083-scaffold1560-228319	0	228319
1	snp14084-scaffold1560-297224	0	297224
1	snp14085-scaffold1560-330454	0	330454
1	snp14086-scaffold1560-378997	0	378997
1	snp14087-scaffold1560-412202	0	412202
1	snp14088-scaffold1560-452762	0	452762
1	snp14089-scaffold1560-488156	0	488156
1	snp14090-scaffold1560-516844	0	516844
1	snp14091-scaffold1560-581464	0	581464
1	snp14092-scaffold1560-614162	0	614162
1	snp14093-scaffold1560-666974	0	666974
1	snp14094-scaffold1560-721048	0	721048
1	snp14095-scaffold1560-755958	0	755958
1	snp14096-scaffold1560-786871	0	786871
1	snp14097-scaffold1560-830197	0	830197
1	snp14098-scaffold1560-872326	0	872326
1	snp14099-scaffold1560-920888	0	920888
1	snp14100-scaffold1560-986550	0	986550
1	snp14101-scaffold1560-1032913	0	1032913
1	snp2819-scaffold1082-727669	0	1087568
1	snp2817-scaffold1082-658683	0	1156554

AA	1011001	0	0
AB	1011002	0	0
AC	1011003	0	0
AD	1011004	0	0
AA	1011005	0	0
AB	1011006	0	0
AC	1011007	0	0
AD	1011008	0	0
AA	1011009	0	0
AB	1011011	0	0
AC	1011015	0	0
AD	1011016	0	0
AA	1011017	0	0
AB	1011018	0	0
AC	1011021	0	0
AD	1011022	0	0
AA	1011023	0	0
AB	1011026	0	0
AC	1011028	0	0
AD	1011029	0	0
AA	1011030	0	0
AB	1011032	0	0
AC	1011033	0	0
AD	1011035	0	0
AA	1011036	0	0
AB	1011037	0	0
AC	1011038	0	0
AD	1011039	0	0





## File formats for GBLUP analysis

UGA42014 210021212111121110010011121102001210220012120101101211110211211210101101110102  
UGA42019 201011202112112010111120121000110121112100212020211021211211112220210000110102  
UGA42029 10111020112111010000112112110010111102200121111121001112120121211120110111021001  
UGA42039 100020101221022010101112001101020111102111011202111112121121122101211110111121111  
UGA42047 20002020222202200000000200220020022002220202220202220200000020002  
UGA42051 200020202222022000000002002200200220022202022002220202220200000020002  
UGA42052 200020202222022000000002002200200220022202022002220202220200000020002  
UGA42056 1021011100212100210001111221200111011011001210211200011111020210121021111012110  
UGA42057 200020202222022000000002002200200220022202022002220202220200000020002  
UGA42061 100020101221022010101112001101020111102111011202111112121122101211110111121111  
UGA42085 10101120112102210101000211220111112112210121201121111211221022120210201001120012  
UGA42088 00111010012011102010122211101101020020211110111210111011221021110111020222122110  
UGA42094 01121110111110110000221222121010111212110210011210111002211011121020010211121101  
UGA42095 20002020222202200000000200220020022002200220022002200220022002200200000020002  
UGA42098 100020101221022010101112001101020111102111011202111112121122101211110111121111  
UGA42101 01210120101111012001021212102001121121110111100011220011210020120111011110  
UGA42108 200020202222022000000002002200200220022202022202220200000020002  
UGA42109 101110201121111010000112112100101111022001211111210011121211120110111021001  
UGA42127 200020202222022000000002002200200220020022002200220022002200200000020002  
UGA42136 100020101221022010101112001101020111102111011202111112121122101211110111121111  
UGA42137 10210111002121002100011112212001110110110012102112000111110202101210211110121111  
UGA42138 01210120101111012001021212102001121121110111100011220011210020120111011110  
UGA42139 1000201012210220101011120011010201111021101202111112121121122101211110111121111

SNP data file

Index	Name	Chromosome	Position	GenTrain	Score	SNP	ILMN	Strand	Customer	Strand	NormID
1	ARS-BFGL-BAC-10172	14	6371334	0.9176	[A/G]	TOP	TOP	2			
2	ARS-BFGL-BAC-1020	14	7928189	0.9413	[T/C]	BOT	TOP	2			
3	ARS-BFGL-BAC-10245	14	31819743	0.7646	[T/C]	BOT	BOT	2			
4	ARS-BFGL-BAC-10345	14	6133529	0.8906	[A/C]	TOP	TOP	2			
5	ARS-BFGL-BAC-10365	14	27005721	0.9206	[A/C]	TOP	TOP	BOT		1	
6	ARS-BFGL-BAC-10375	14	6616434	0.9258	[A/G]	TOP	TOP	2			
7	ARS-BFGL-BAC-10591	14	17544926	0.7439	[A/G]	TOP	TOP	1			
8	ARS-BFGL-BAC-10867	14	34639444	0.9085	[G/C]	BOT	BOT	101			
9	ARS-BFGL-BAC-10919	14	31267746	0.8255	[A/G]	TOP	TOP	2			
10	ARS-BFGL-BAC-10951	10	17911906	0.9056	[T/G]	BOT	BOT	2			
11	ARS-BFGL-BAC-10952	10	18882288	0.9184	[A/G]	TOP	TOP	2			
12	ARS-BFGL-BAC-10960	10	20609250	0.5678	[A/G]	TOP	TOP	2			
13	ARS-BFGL-BAC-10972	10	20792754	0.8432	[G/C]	BOT	BOT	102			
14	ARS-BFGL-BAC-10975	10	21225382	0.7991	[A/G]	TOP	TOP	2			
15	ARS-BFGL-BAC-10986	10	26527257	0.8941	[A/C]	TOP	BOT	2			
16	ARS-BFGL-BAC-10993	10	78512500	0.8649	[A/G]	TOP	BOT	2			
17	ARS-BFGL-BAC-11000	10	79252023	0.9433	[T/G]	BOT	BOT	2			
18	ARS-BFGL-BAC-11003	10	80410977	0.8842	[T/C]	BOT	BOT	2			
19	ARS-BFGL-BAC-11007	10	80783719	0.9110	[T/C]	BOT	BOT	2			
20	ARS-BFGL-BAC-11025	10	84516867	0.8711	[T/G]	BOT	BOT	2			

Map file

## Quality control QC using plink in R

---

The following script helps to evaluate the QC in plink.

- Download plink from: <https://www.cog-genomics.org/plink/>
- put plink (only the executable file) at your working directory

### Quality Control exercise

```
####convert the data from .ped to .bed format
system("./plink --file Sheep05_724 --make-bed --chr-set 26 --out Sheep05_724")
#missing at SNP level
system("./plink --bfile Sheep05_724 --chr-set 26 --make-bed --geno 0.01 --out Sheep05_afterQC")
#minor allele frequency
system("./plink --bfile Sheep05_724 --chr-set 26 --make-bed --geno 0.10 --maf 0.10 --out Sheep05_afterQC")
#missing at animal level
system("./plink --bfile Sheep05_724 --chr-set 26 --make-bed --geno 0.10 --maf 0.10 --mind 0.05 --out
Sheep05_afterQC")
```



# Id-pruning

---

LD-pruning:

```
system("./plink --bfile Sheep05_afterQC --chr-set 26 --indep-pairwise 50 10 0.5 ")
system("./plink --bfile Sheep05_afterQC --extract plink.prune.in --make-bed --chr-set 26 --out
Sheep05_afterQC_pruned")
```

#To generate the PCA:

```
system("./plink --bfile Sheep05_afterQC_pruned --chr-set 26 --pca --outSheep05_afterQC_pruned_pca ")
```

#to generate structure format

```
system("./plink --bfile Sheep05_afterQC_pruned --recode-structure --chr-set 26 --out
Sheep05_afterQC_pruned_structure")
```

Note: For the admixture, you can still use the .bed format

## Tips

#To generate other population genetics parameters:

#To generate the MAF from bed file:

```
system("./plink --bfile Sheep05_afterQC --chr-set 26 --freq --out Sheep05_afterQC_maf")
```

# To generated expected and observed heterozygosity, and hwe:

```
system("./plink --bfile Sheep05_afterQC --hardy --chr-set 26 --out Sheep05_afterQC_het")
```

#To generate coefficient of inbreeding:

```
system("./plink --bfile Sheep05_afterQC --chr-set 26 --het --out Sheep05_afterQC_Fx")
```

## Why checking the PCA is important before proceeding to genomic evaluation?

- It helps to know the structure of the target population so that it gives idea whether we need to implement stratification in the model we will be fitting
- It also helps to provide feedback to the farmers about the structure of the animals they own

# PCA

---

```

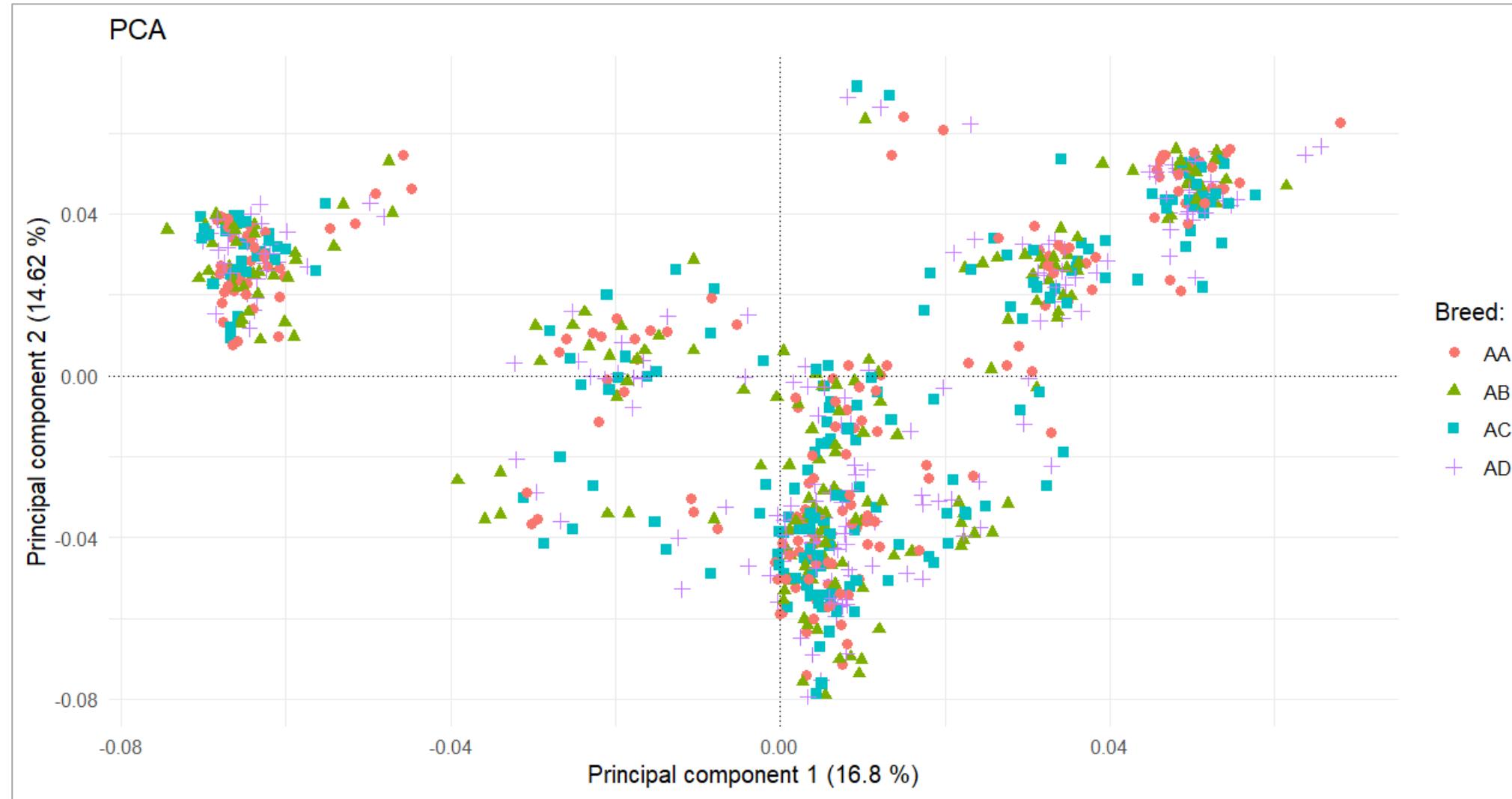
install.packages("tidyverse")
library(tidyverse)
# read in result files
eigenValues <- read_delim("PATH/Sheep05_afterQC_pruned_pca.eigenval", delim = " ", col_names = F)
eigenVectors <- read_delim("PATH/Sheep05_afterQC_pruned_pca.eigenvec", delim = " ", col_names = F)
eigenVectors
eigenValues
## Proportion of variation captured by each vector
eigen_percent <- round((eigenValues / (sum(eigenValues)))*100), 2)
# PCA plot
ggplot(data = eigenVectors) +
  geom_point(mapping = aes(x = X3, y = X4, color = X1, shape = X1), size = 1, show.legend = F ) +
  geom_hline(yintercept = 0, linetype="dotted") +
  geom_vline(xintercept = 0, linetype="dotted") +
  labs(title = "PCA",
       x = paste0("Principal component 1 (",eigen_percent[1,1]," %)"),
       y = paste0("Principal component 2 (",eigen_percent[2,1]," %)"),
       colour = "Breed:", shape = "Breed:") +
  theme_minimal()

```

# PCA plot generated from the example data (*blackface*)



Centre for  
Tropical Livestock  
Genetics and Health



# Admixture package

---

## Why admixture is preferred?

([Alexander and Lange., 2012. Enhancements to the ADMIXTURE algorithm for individual ancestry estimation. Read](#))

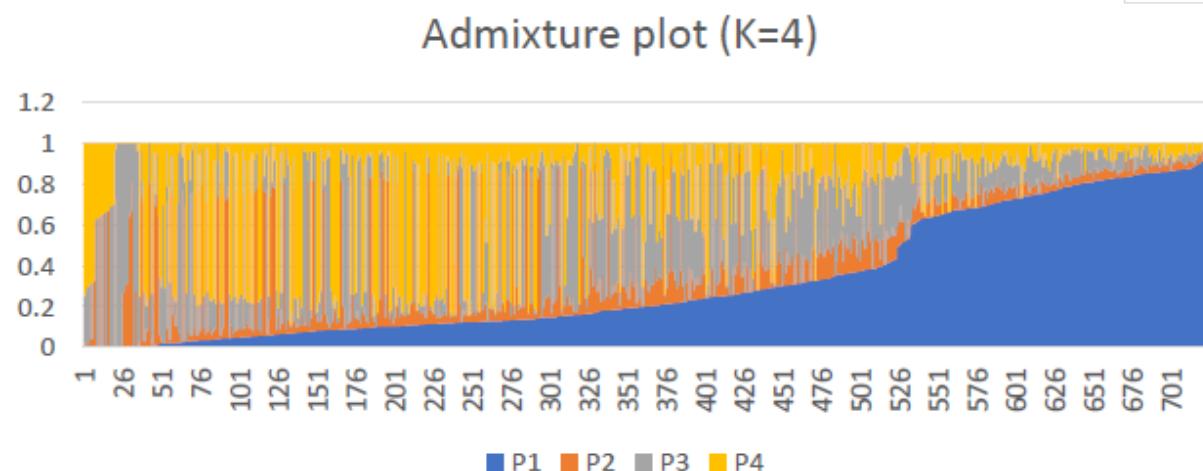
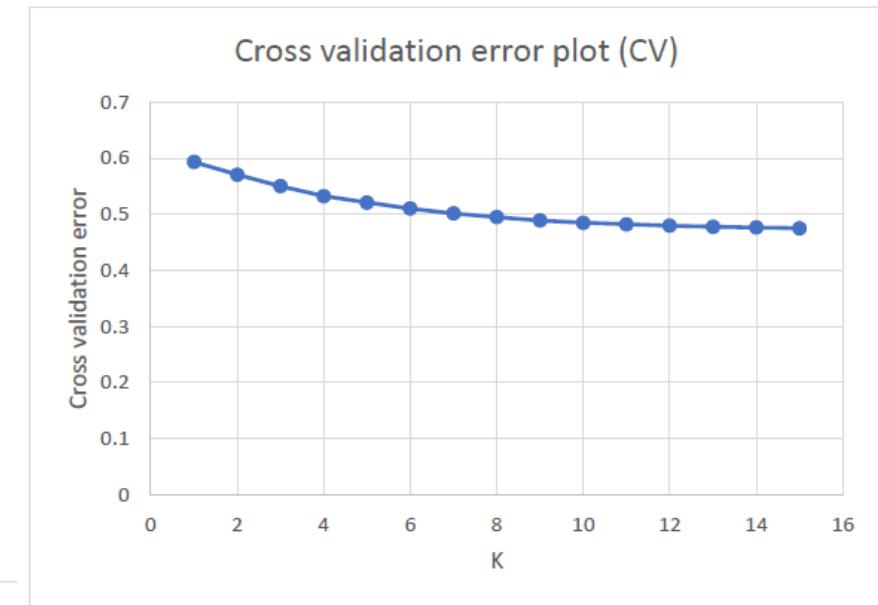
- *Helps to generate the breed proportion of large dataset*
- ADMIXTURE estimates individual ancestries by efficiently computing maximum likelihood estimates in a parametric model.
  1. ADMIXTURE can be used to estimate the number of underlying populations through **cross-validation**.
  2. Individuals of known ancestry can be exploited in supervised learning to yield more precise ancestry estimates.
  3. By penalizing small admixture coefficients for each individual, one can encourage model parsimony, often yielding more interpretable results for small data sets or data sets with large numbers of ancestral populations.
  4. By exploiting multiple processors, large data sets can be analyzed even more rapidly.

# Running admixture

- **Generating the input file:**
- ADMIXTURE requires unlinked (i.e.LD-pruned) SNPs in plink format.
- It is very easy to generate the input file from a VCF containing such SNPs.
- Plink helps to generate the .bed file which can be read by ADMIXTURE
- The default cross-validation 5-fold CV (you can change as the K value you expect to be) and starts as K=2.
- **ADMIXTURE produced 2 files:** *Q* which contains cluster assignments for each individual and *P* which contains for each SNP the population allele frequencies.
- The default cross-validation 5-fold CV (you can change as the K-value you expect to be) and starts as K = 2.
- Grep the results to access the CV error log Files as: `grep 'CV error' log_*`

- Running population admixture in using admixture package

```
bed=/PATH/shp724afterQCFinal.bed
mkdir -p / PATH /newFolder
out=/ PATH / newFolder
cd / PATH / newFolder
for k in $(echo {1..15..1});
do echo ${k};
./admixture --cv=10 ${bed} ${k} > ${out}/log_${k}.txt;
done
grep 'CV error' log_*
```



# Formatting the genotype data for GWAS

- # pick chromosomes of interest for GWAS:
- ./plink --file Sheep05\_724 --recode --chr 1-26 --out Sheep05\_afterQC
- #converting from Allele form to numeric for blupf90 package
- ./plink --bfile Sheep05\_afterQC --chr-set 26 --recode A --out shp\_blupf90
- # removing first line or header
- cat shp\_blupf90.raw | sed 1d > snpblup.txt
- #rmoving column
- awk '{\$1=\$3=\$4=\$5=\$6="";print \$0}' snpblup.txt > snpblup1.txt
- #removing space from column
- awk '{s=\$1;gsub(\$1 FS,x);\$1=\$1;print s FS \$0}' OFS= snpblup1.txt > snpblup2.txt
  - OR use this command: sed -i "s/ //g" snpblup1.txt
- #adding and removing space between column 1 and 2
- awk '{print ""\$1" "\$2 }' snpblup2.txt >snpblupclean3.txt
- #slecting map file
- awk '{print ""\$2" "\$1" "\$4 }' shp724afterQCFinal.map > gwasmapxxx.txt



Centre for  
Tropical Livestock  
Genetics and Health

---

## CTLGH Funders

BILL & MELINDA  
GATES foundation



Biotechnology and  
Biological Sciences  
Research Council

