



From Bologna to Melbourne: a research experience in modern statistical genetics

Anna Quaglieri

University of Bologna

22nd January 2016

Walter and Eliza Hall Institute of Medical Research

15 different scientific divisions

Wehi divisions





Welcome to the Walter and Eliza Hall Institute website for local bioinformatic resources.

WEHI bioinformatic resources

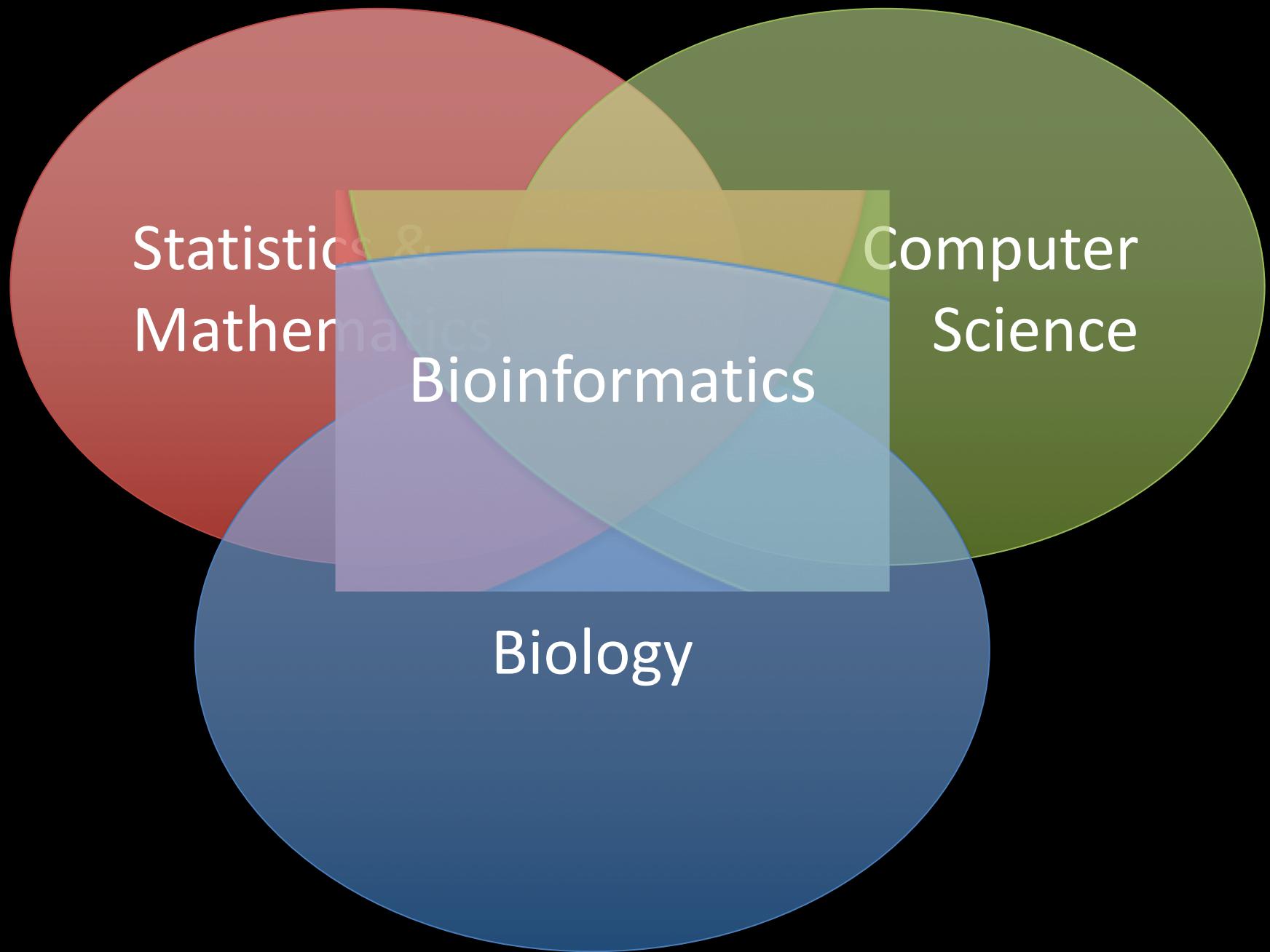
- Software projects
- Supplementary materials for published articles

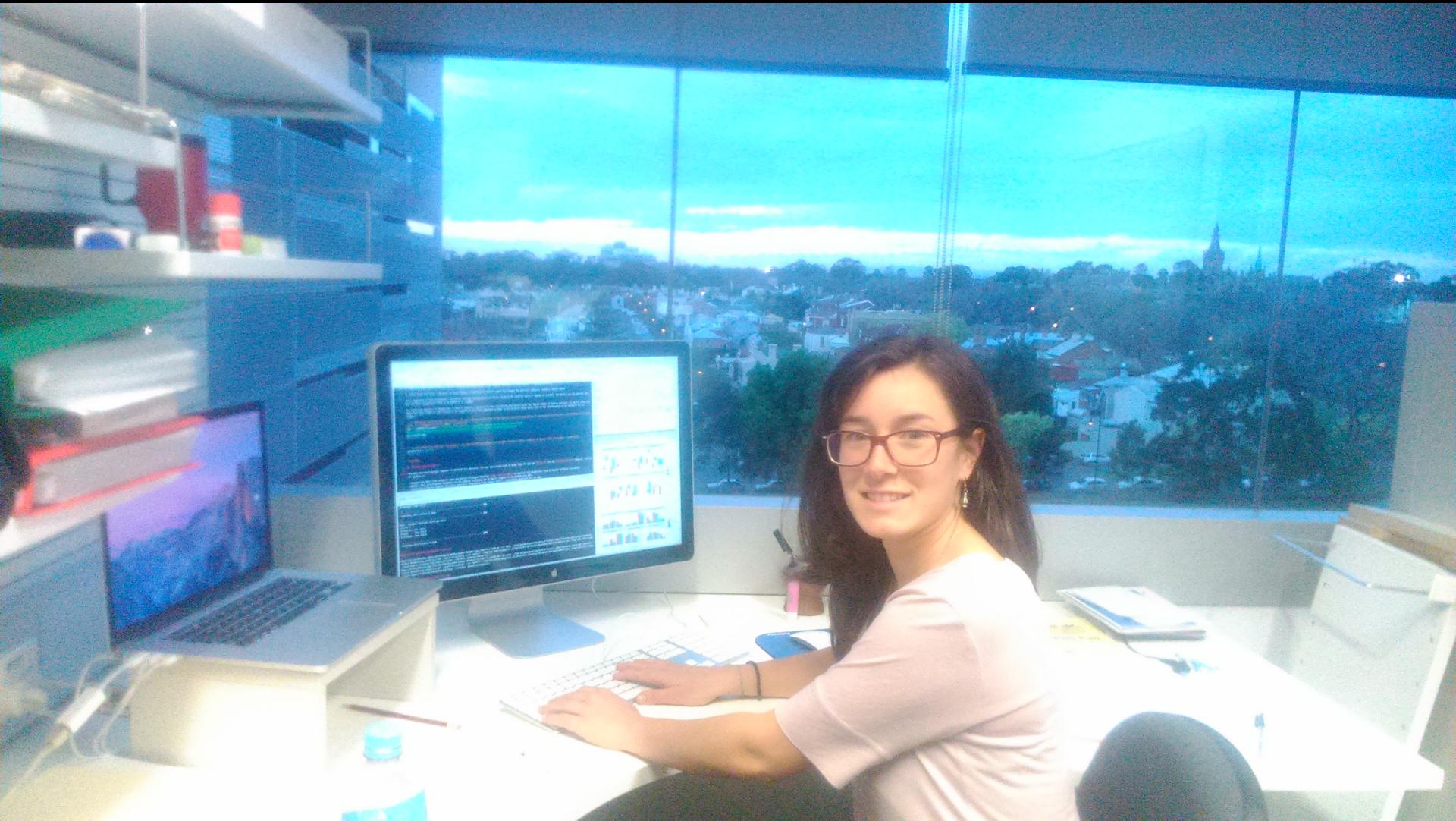
Bioinformatics division seminars

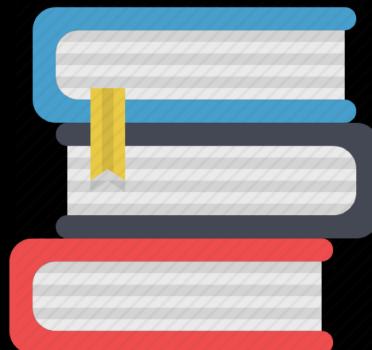
- Current seminar
- All seminars for this year
- Search past seminars

Bioinformatics-orientated research laboratories

- Bahlo Lab
- Papenfuss Lab
- Ritchie Lab
- Smyth Lab
- Speed Lab



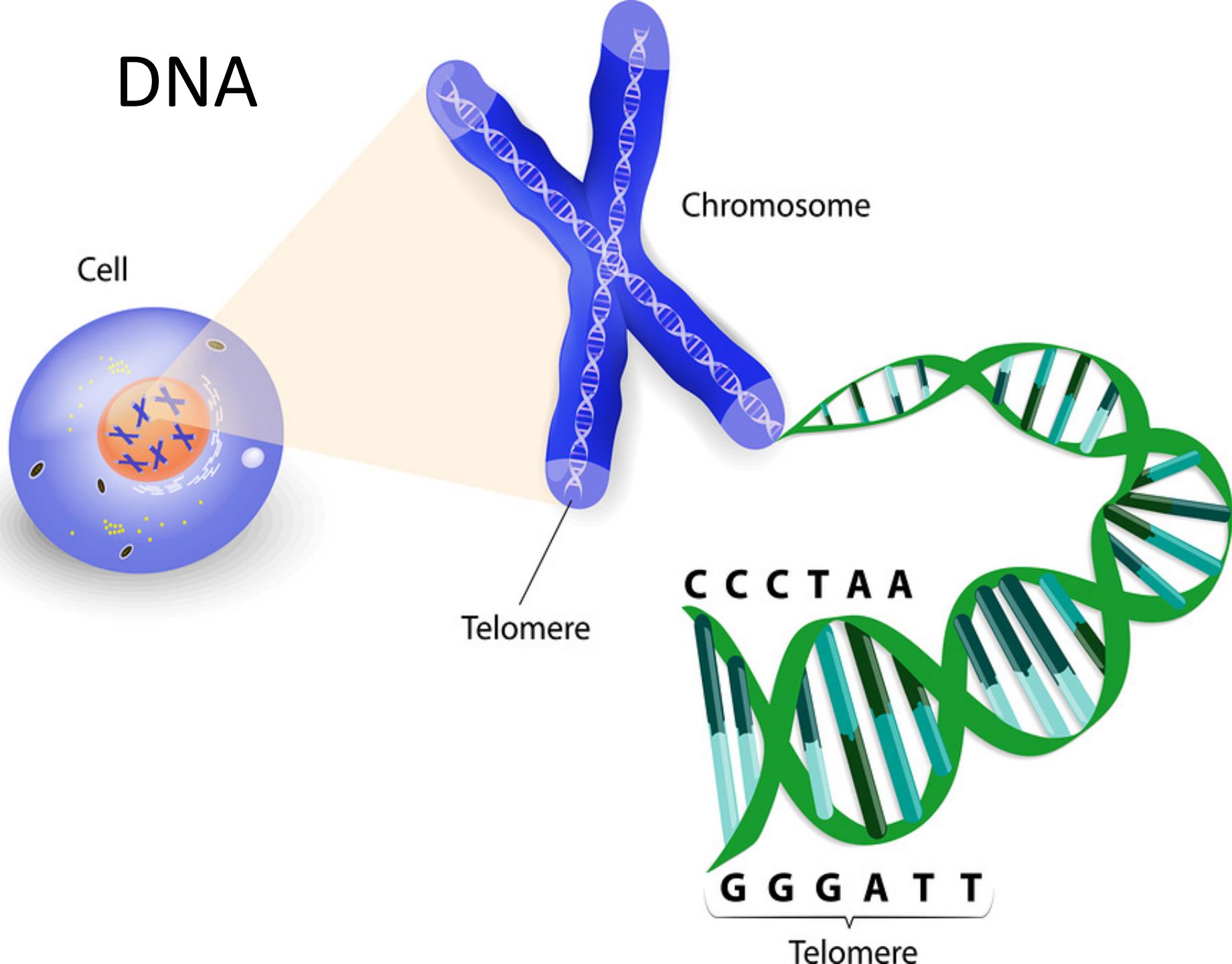




My main projects

- DNA – GWAS on MacTel – an eye disease
- Protein-DNA interaction data: ChIP-Seq

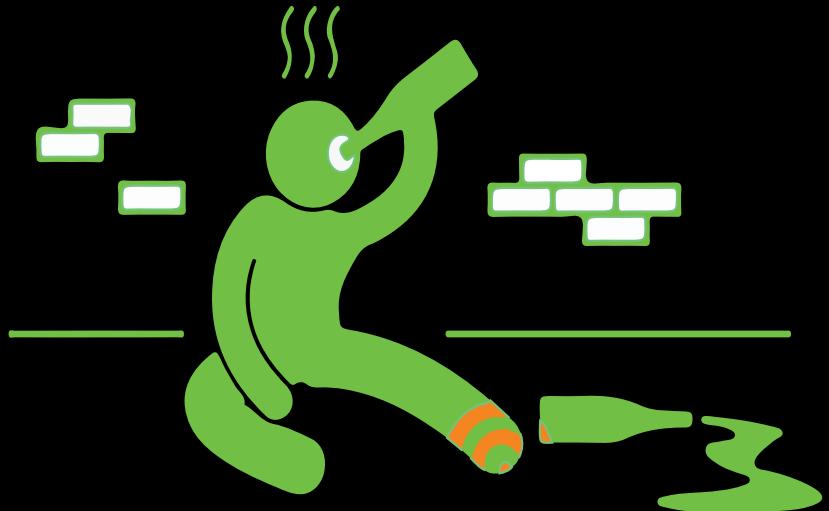
DNA



DNA analysis

- Comparing *people's code*
- 3 billion letters each
- Extremely intensive/expensive task

Select representative letters



99% the
same





GTACTGGTCATTGAG TCTTTGCA AC
CATG ACCAG TAACTCA AAGAAACGTTGAC

T

G

T

T



GTACTGGTCATTGAG TCTTTGCA AC
CATTGACCAGTAACTCA AAGAAACGTTGAC

GGTCATCGTATGGCAT ACAAAATGCCAT
CCA GTAGATAAC CG TAATGTTA CGGTA



SNP T/G

T 90%

G 10%

SNPs as
representatives of
people's genome

Associate a SNP to the
occurrence of a disease

Vanilla GWAS

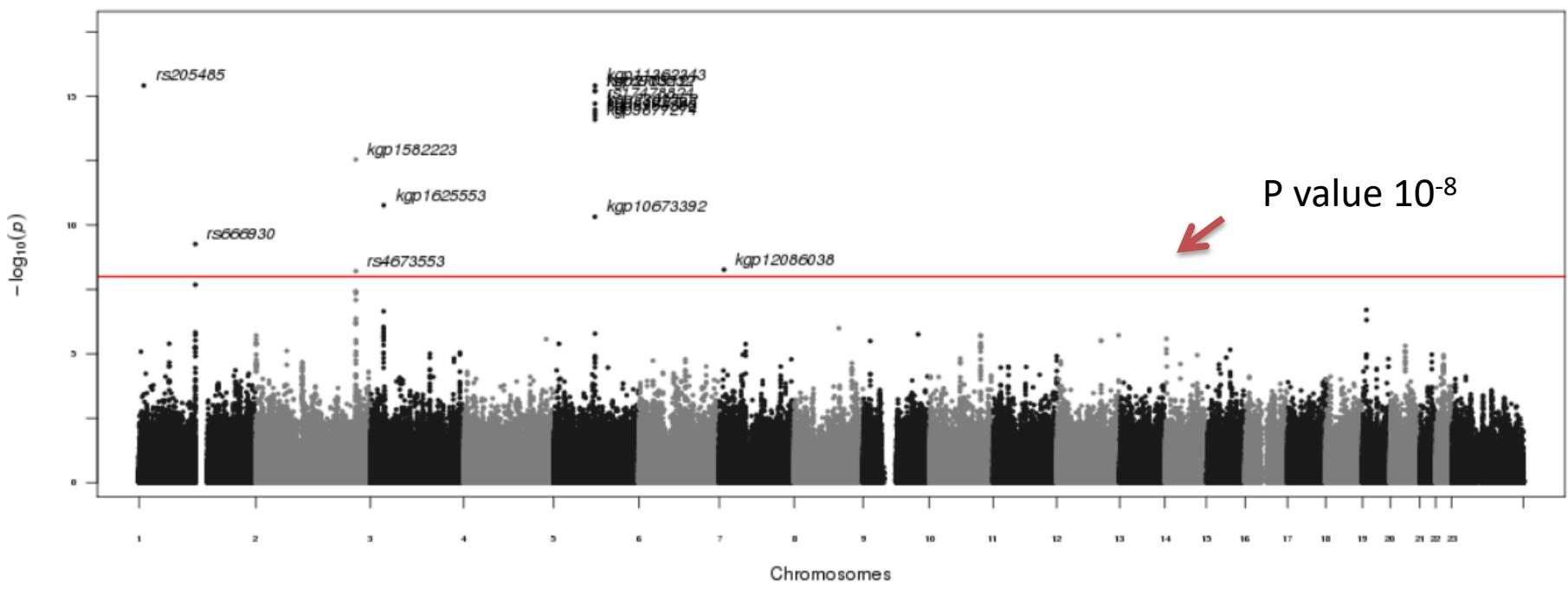
SNP_{A/B}

Allele	A	B	Total
Case	n ₁₁	n ₁₂	N ₁₀
Control	n ₂₁	n ₂₂	N ₂₀
Total	N ₀₁	N ₀₂	N

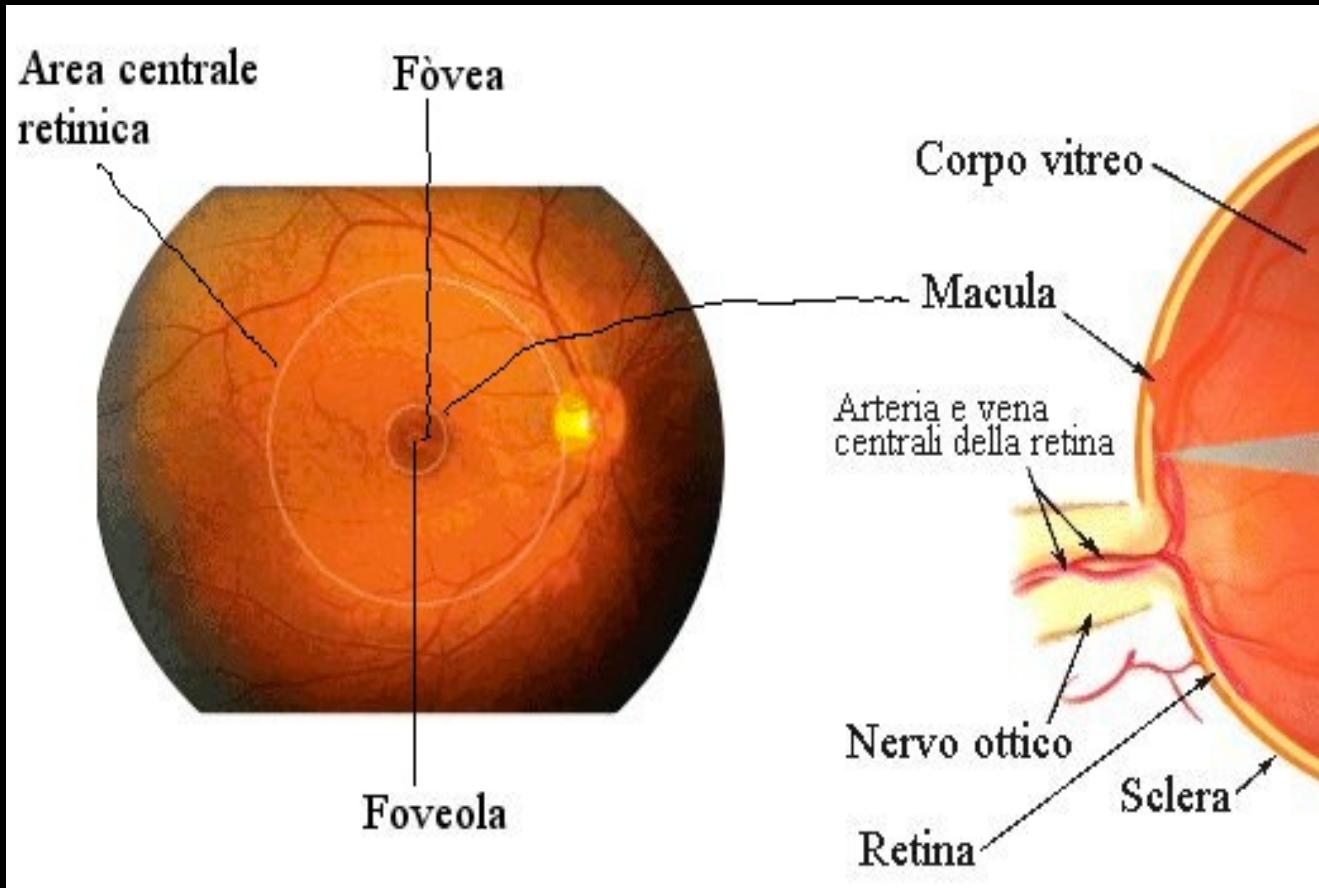
Odds ratios of association SNP by SNP

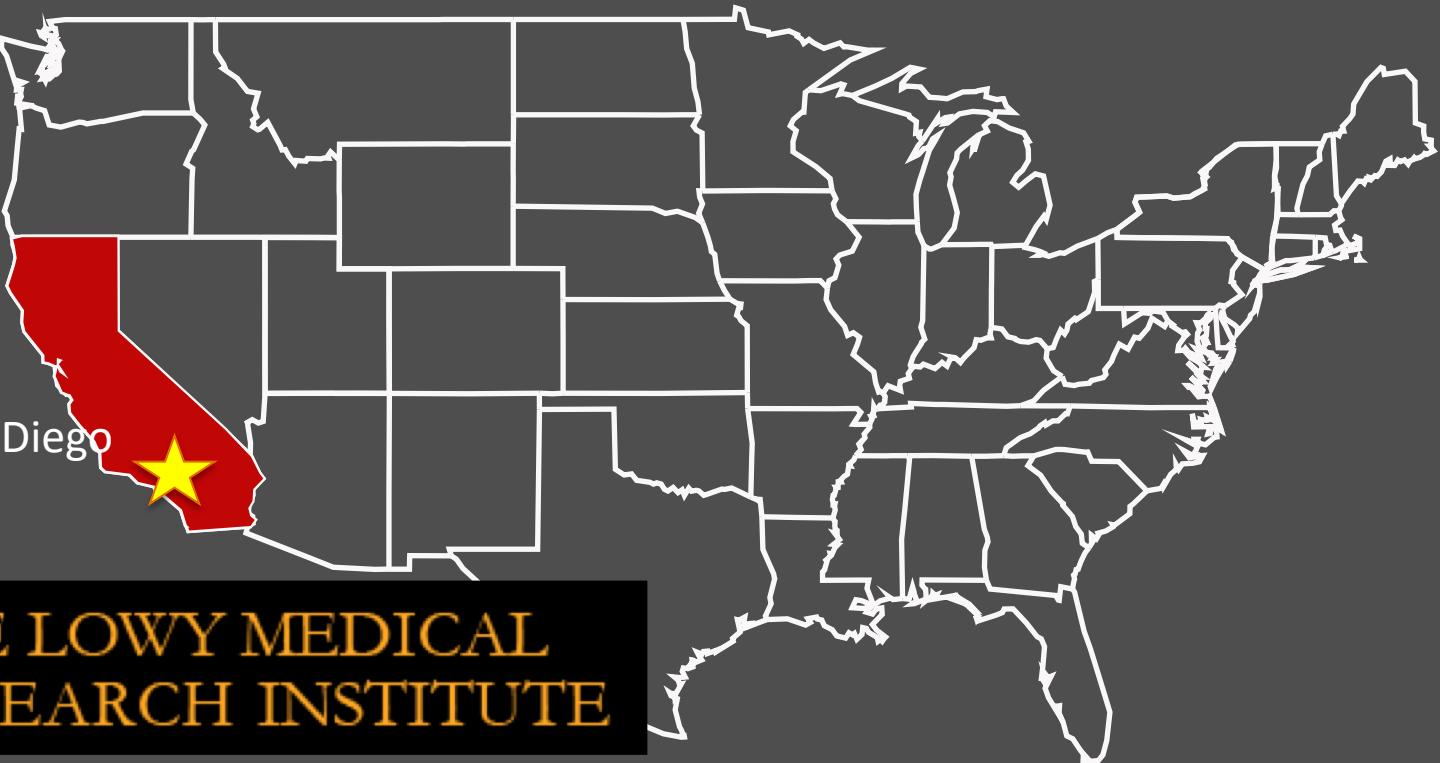
GWAS are **association** studies. Any discovered SNPs **not** be thought of as **causative**, until investigated further.

~1million SNPs



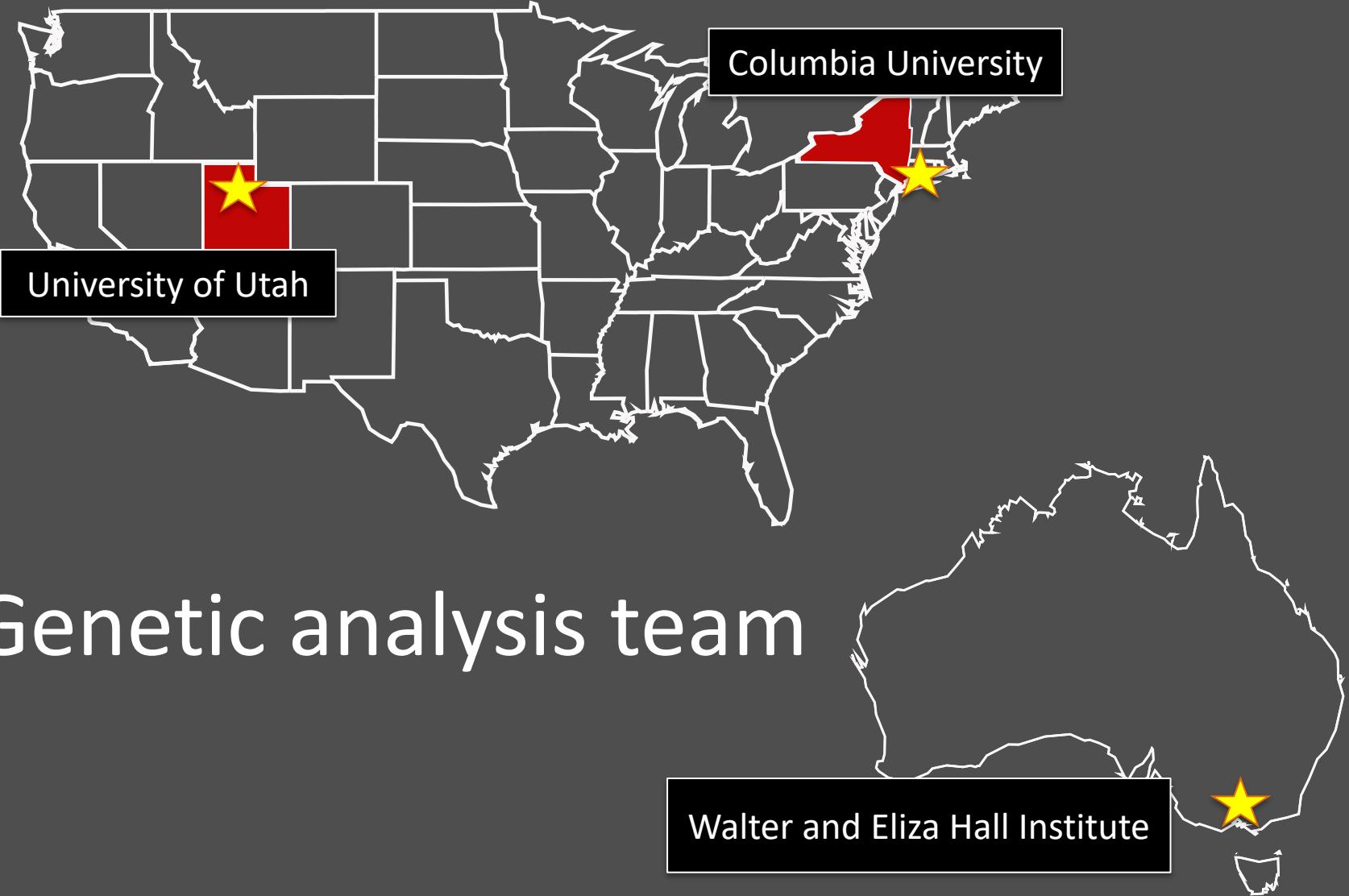
Case Study: MacTel





San Diego

THE LOWY MEDICAL
RESEARCH INSTITUTE

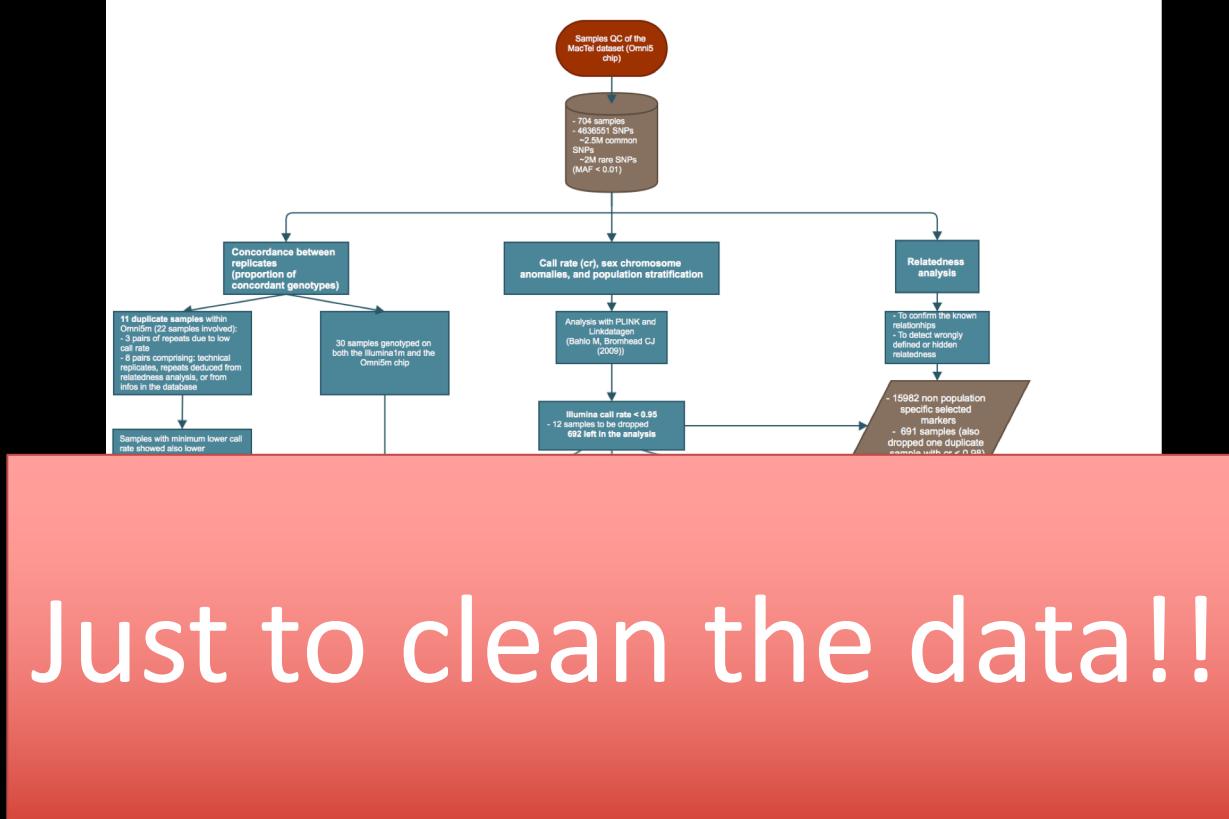


WEHI - Statistical genetic group

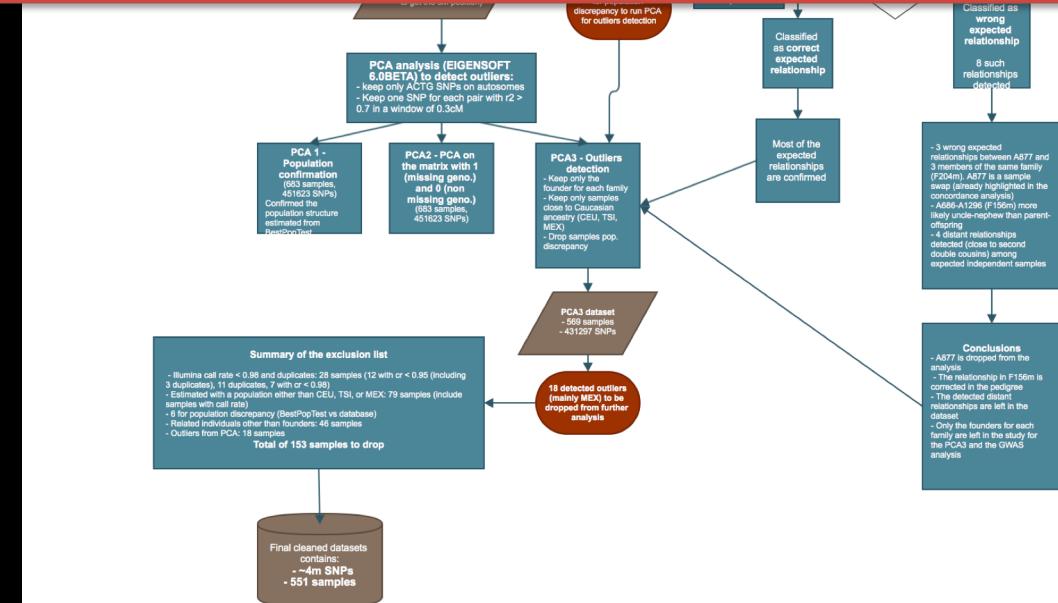


Affected people's DNA





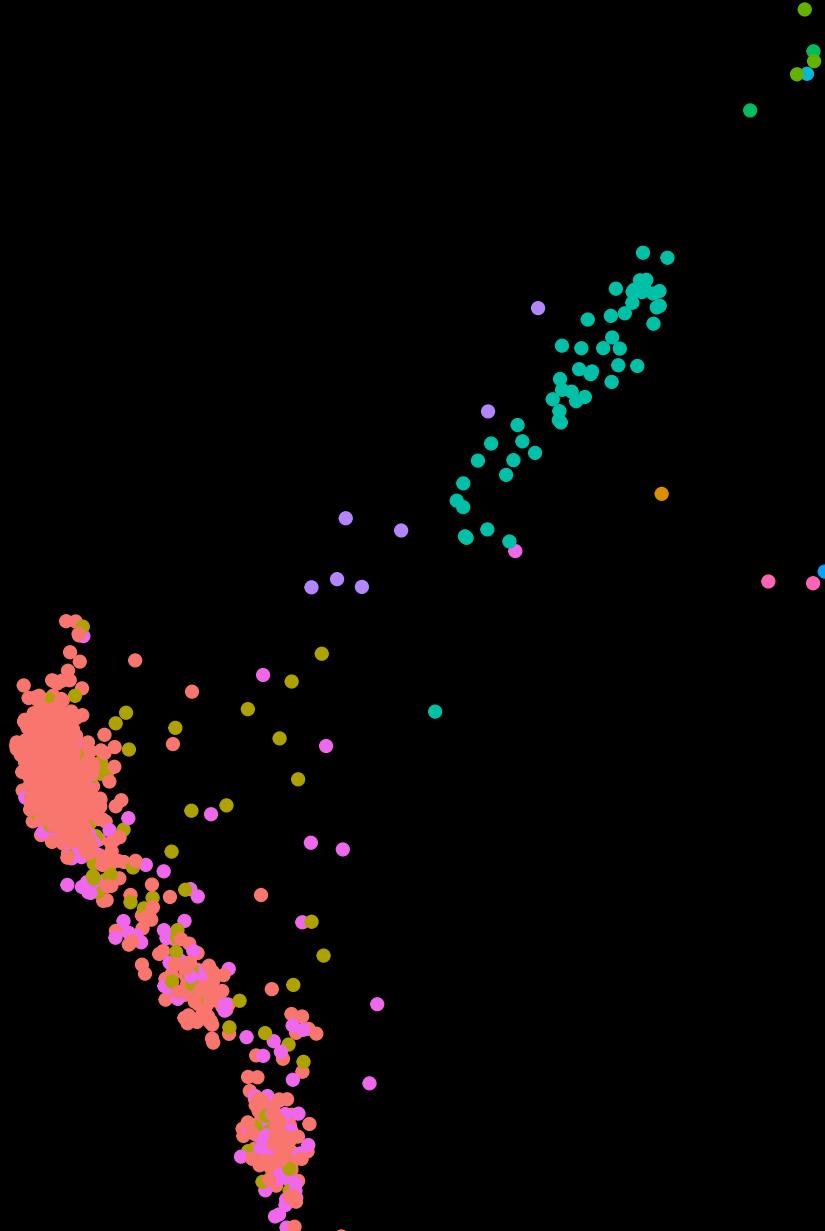
Just to clean the data!!



PC2

PC1

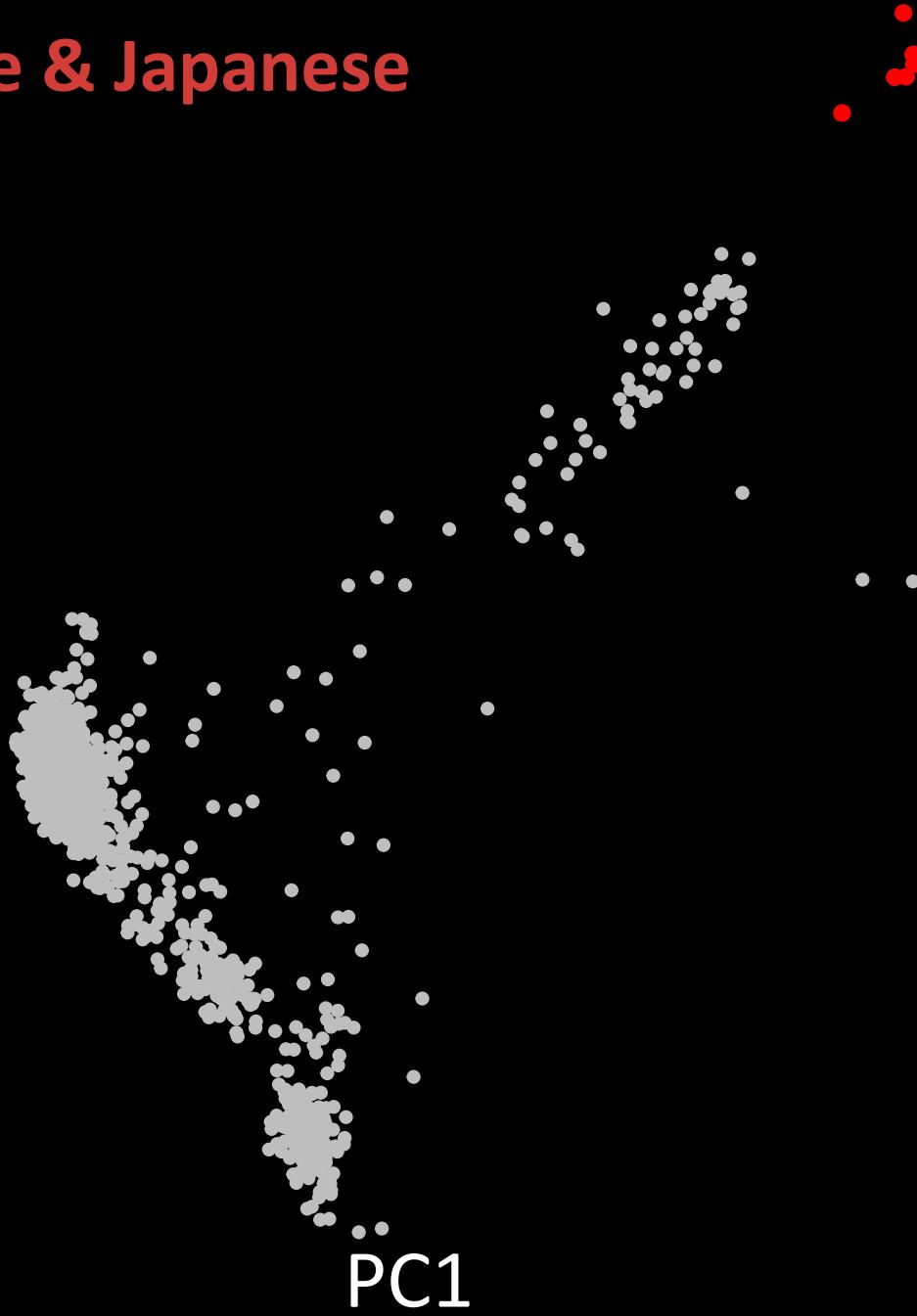
- areds
- ASW
- CEU
- CHB
- CHD
- GIH
- JPT
- LWK
- MEX
- TSI
- YRI



Chinese & Japanese

PC2

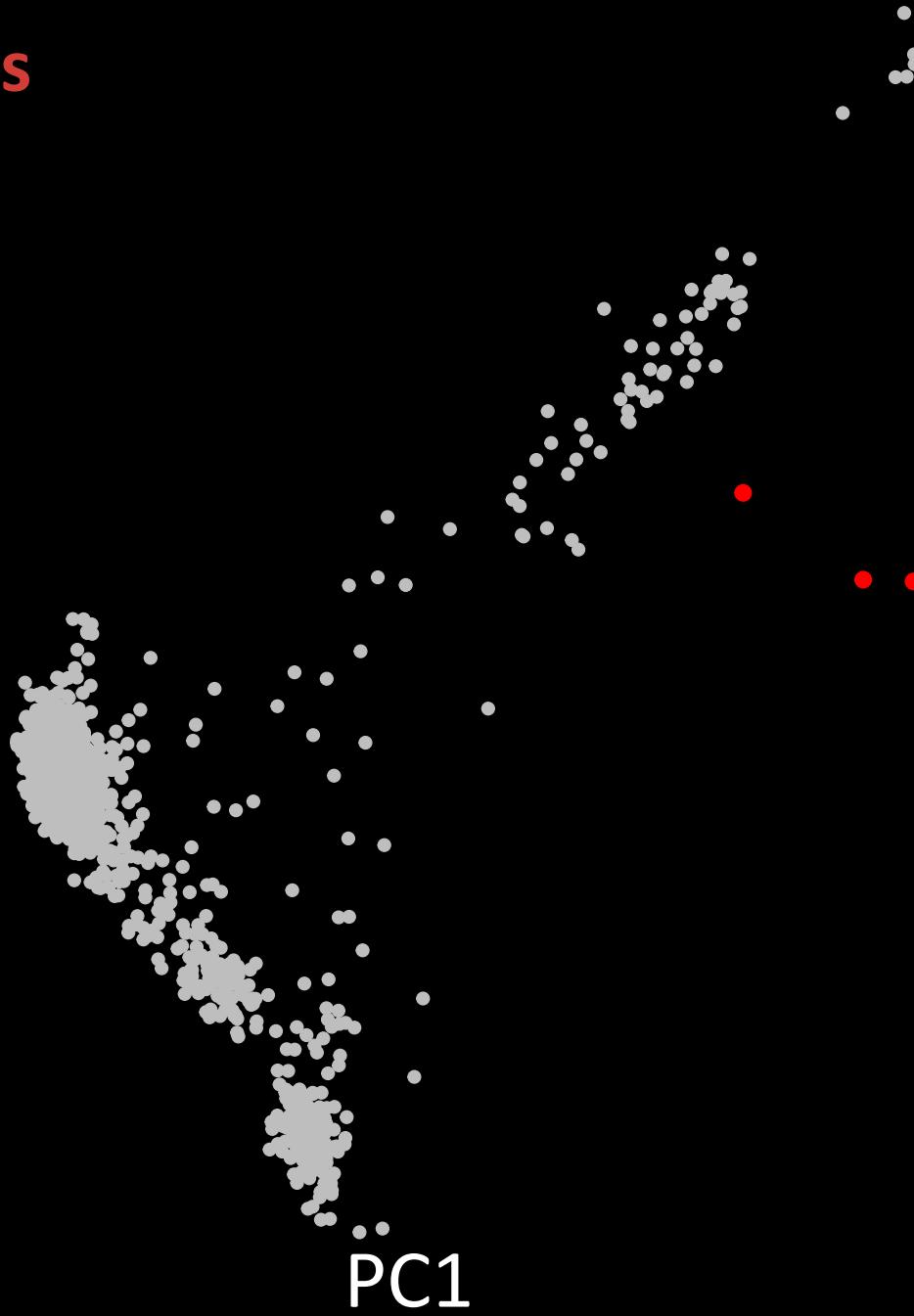
PC1



Africans

PC2

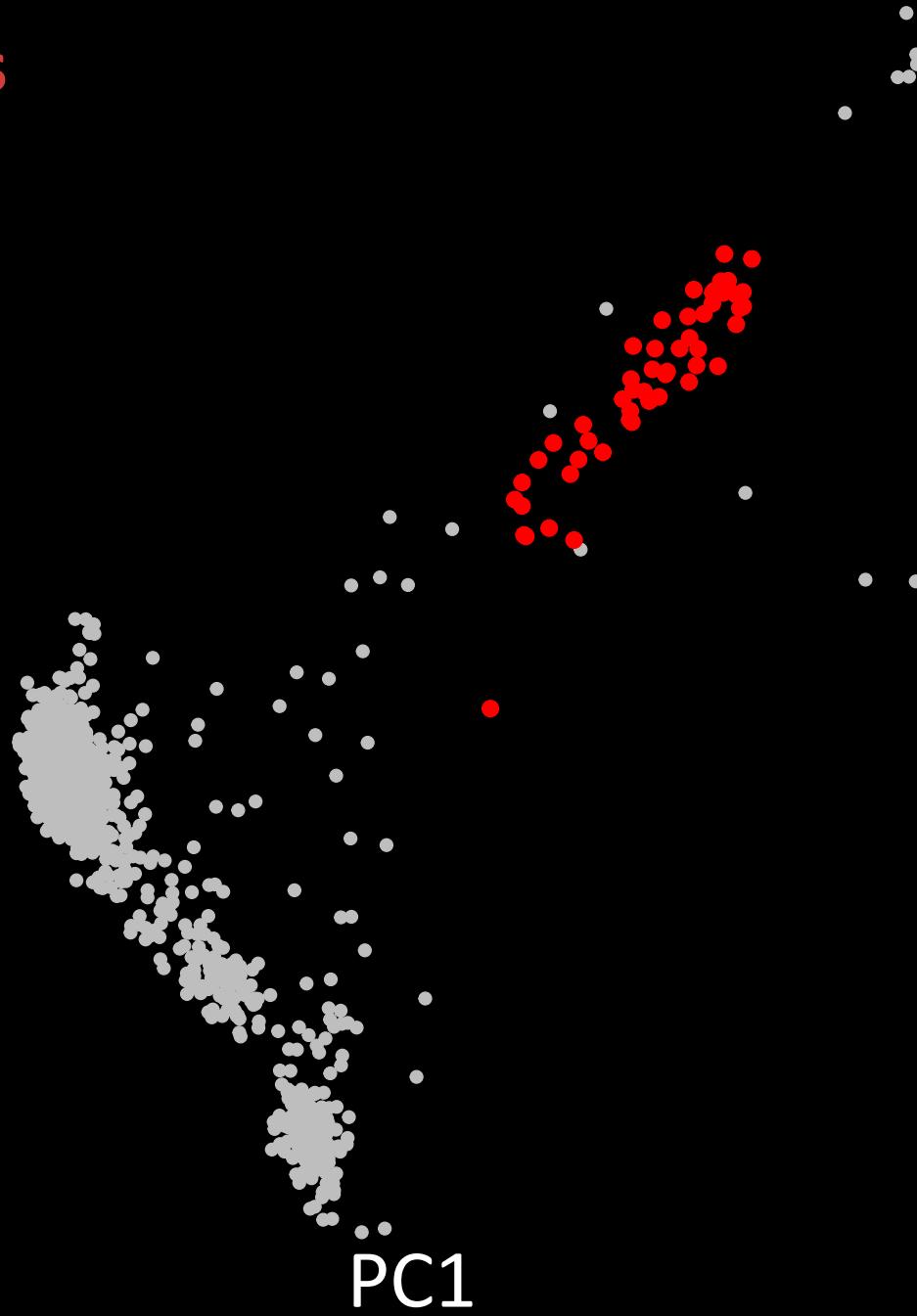
PC1



Indians

PC2

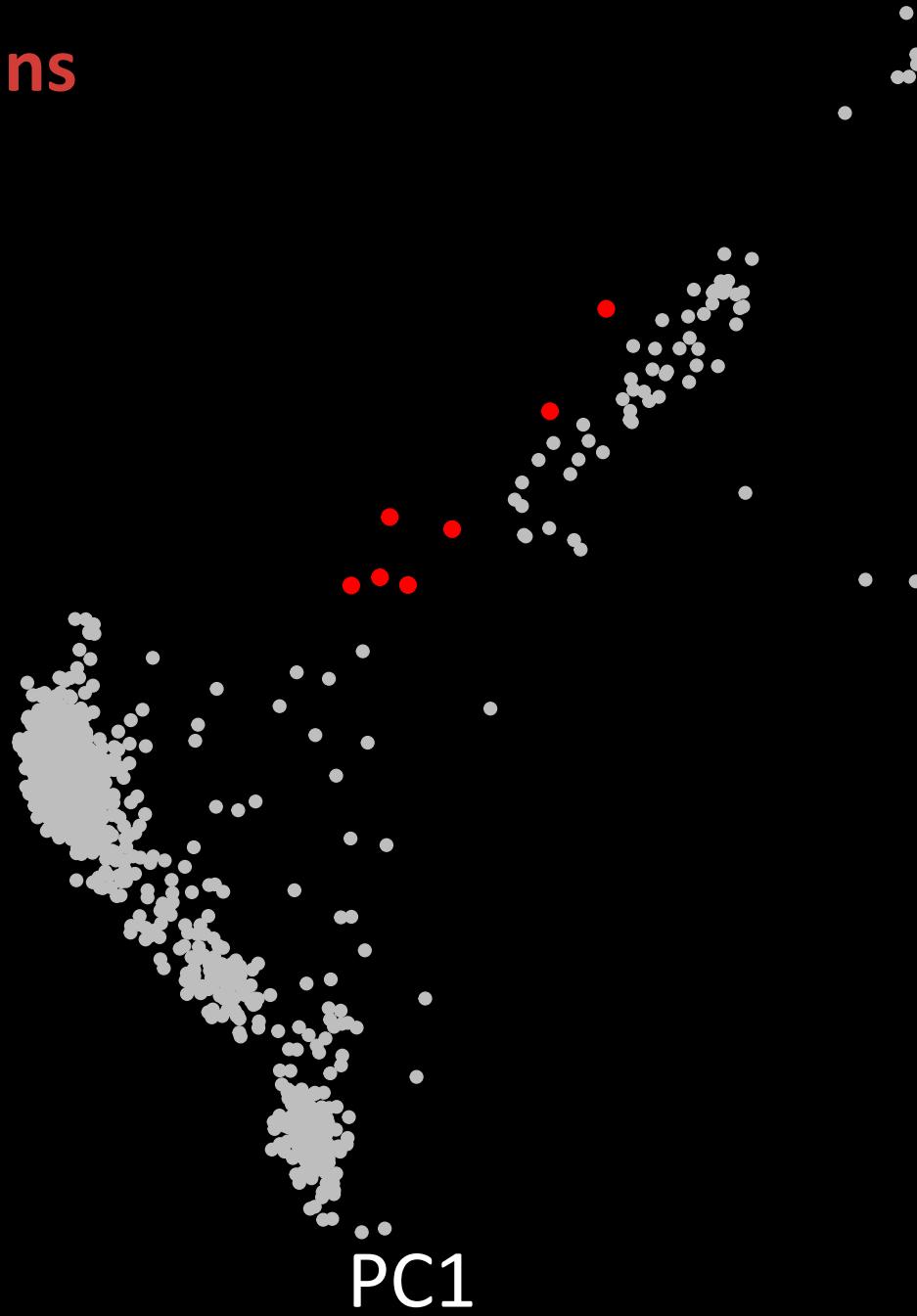
PC1



Mexicans

PC2

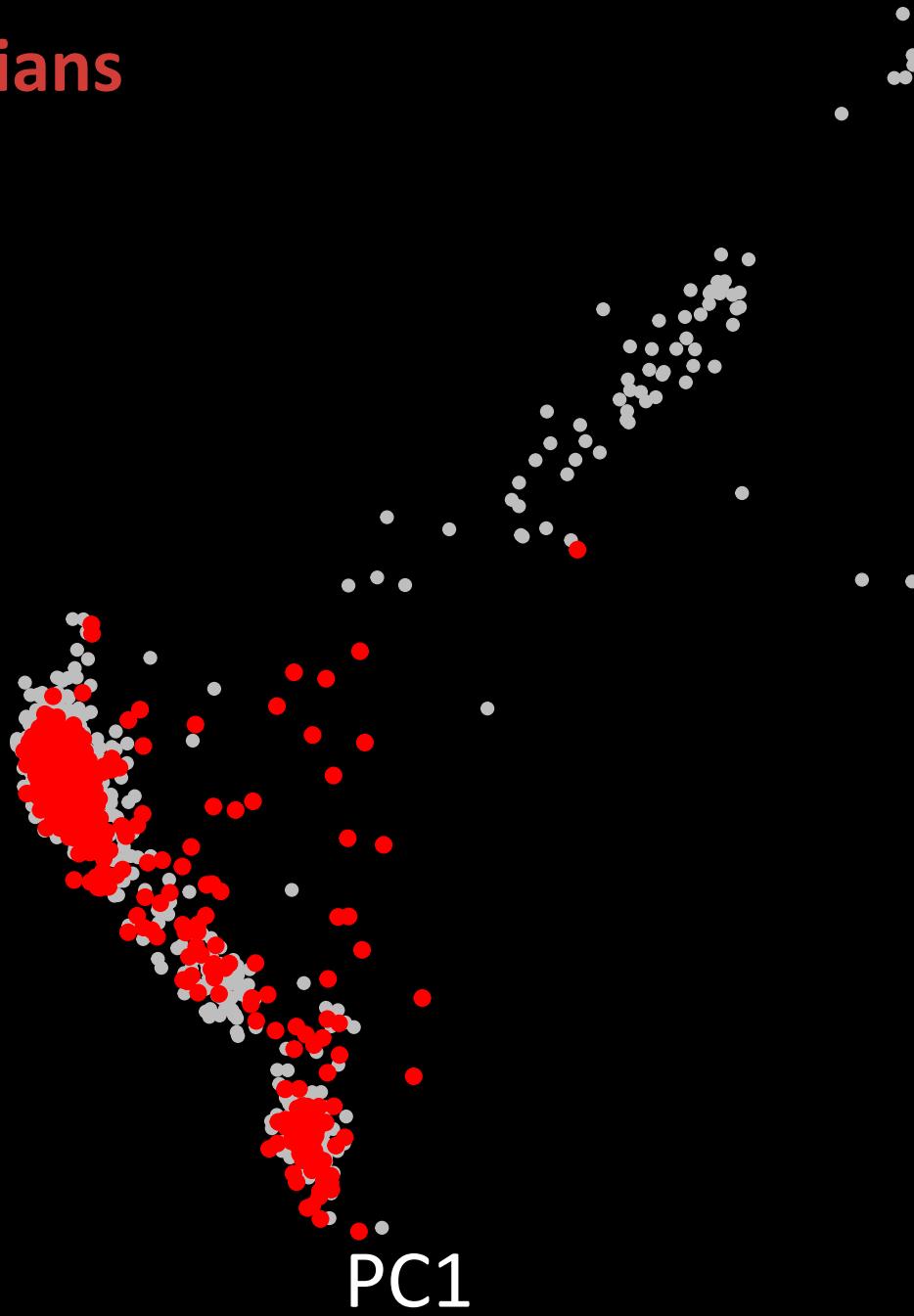
PC1



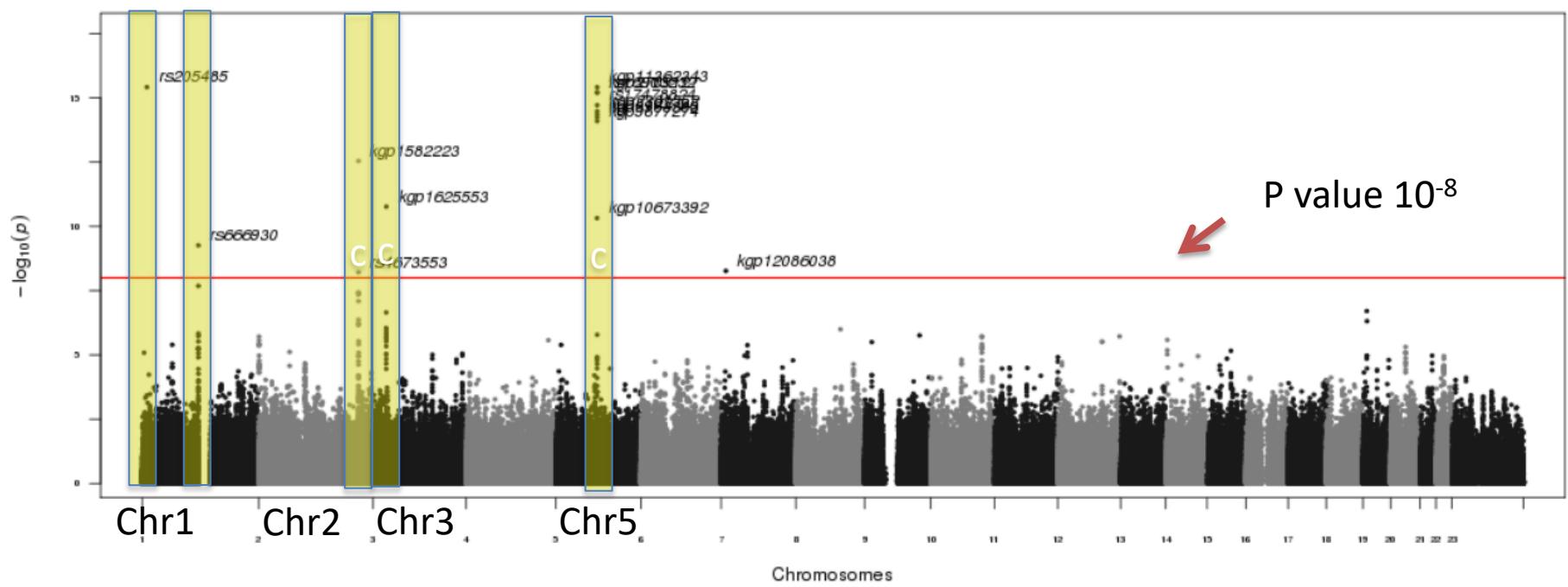
Caucasians

PC2

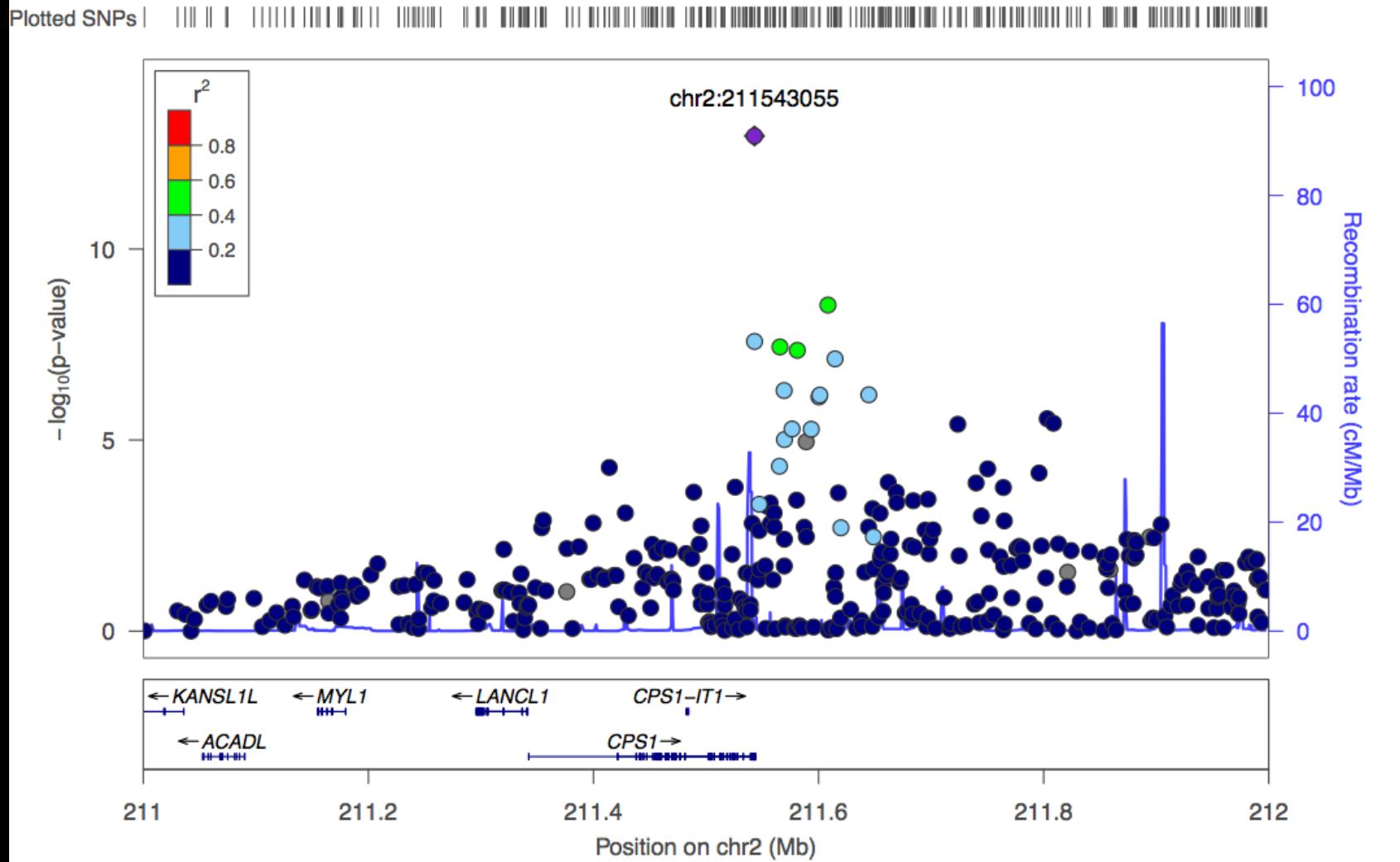
PC1



**Logistic regression (Mt + Areds)
(2 e.vectors correction)**



kgp1582223 – MacTel + AREDS

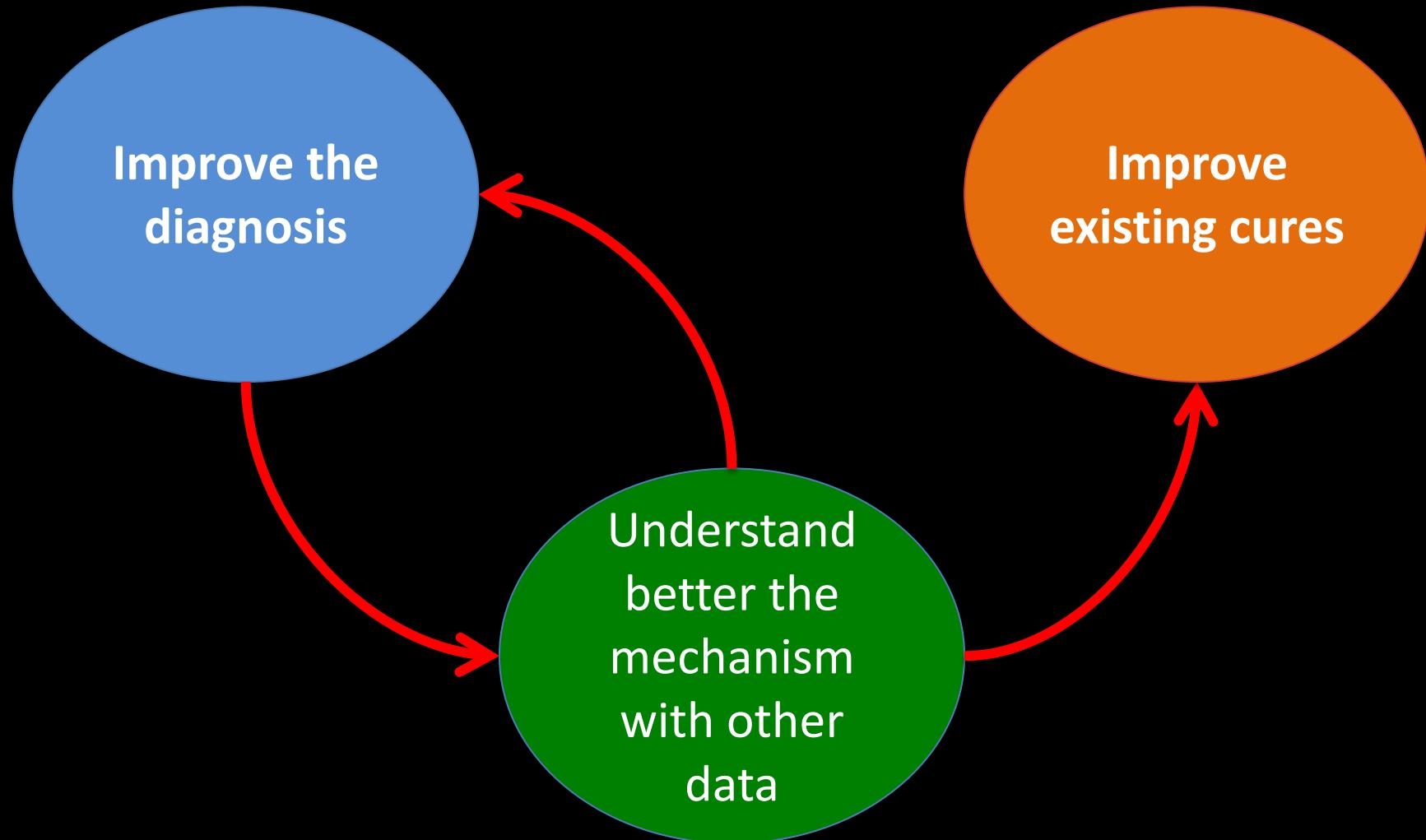




Interpretation to understand causation

- Search through literature
- More experiments related to the specific genes highlighted by GWAS
- Understand the function of the genes involved

Next important aims



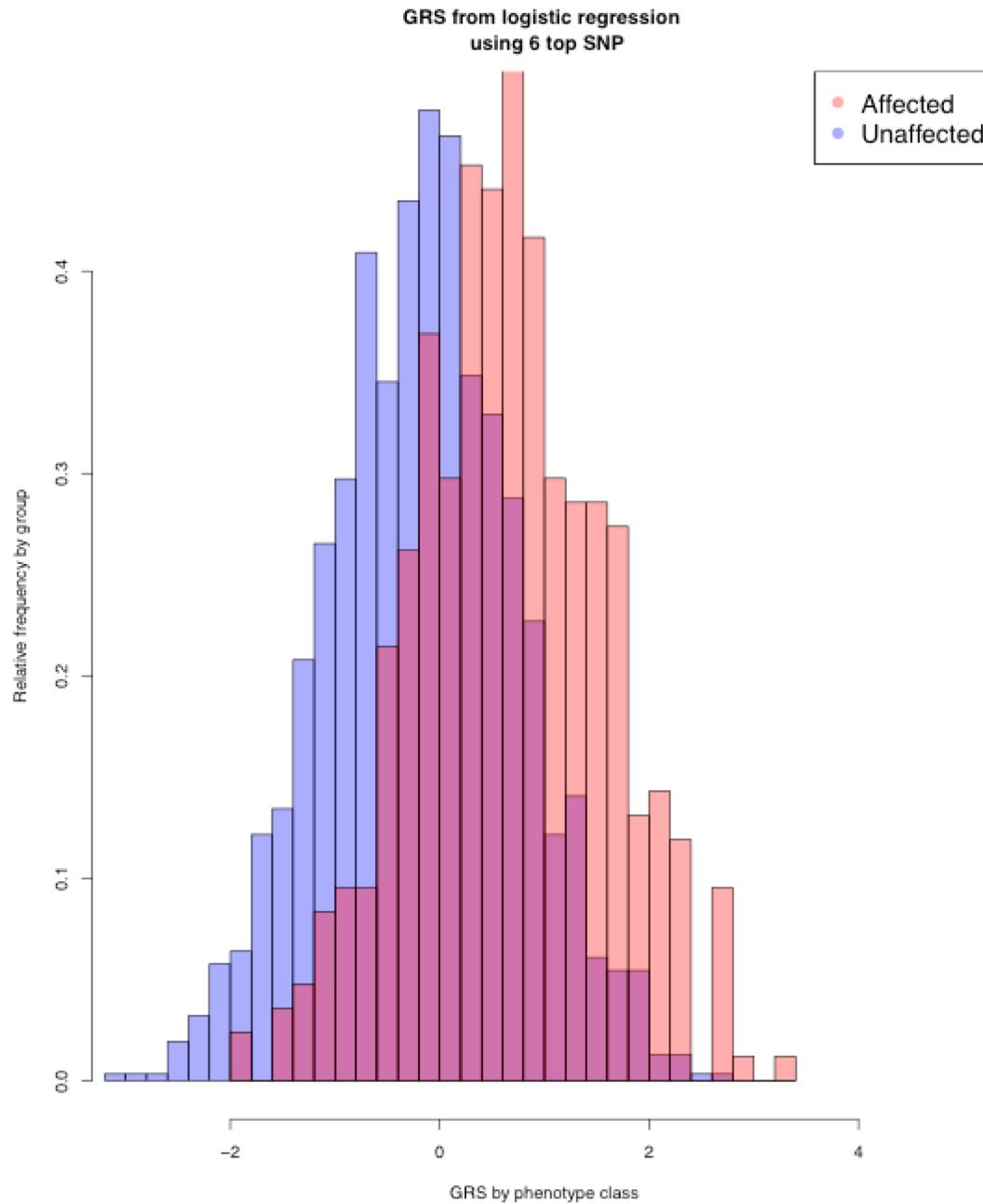
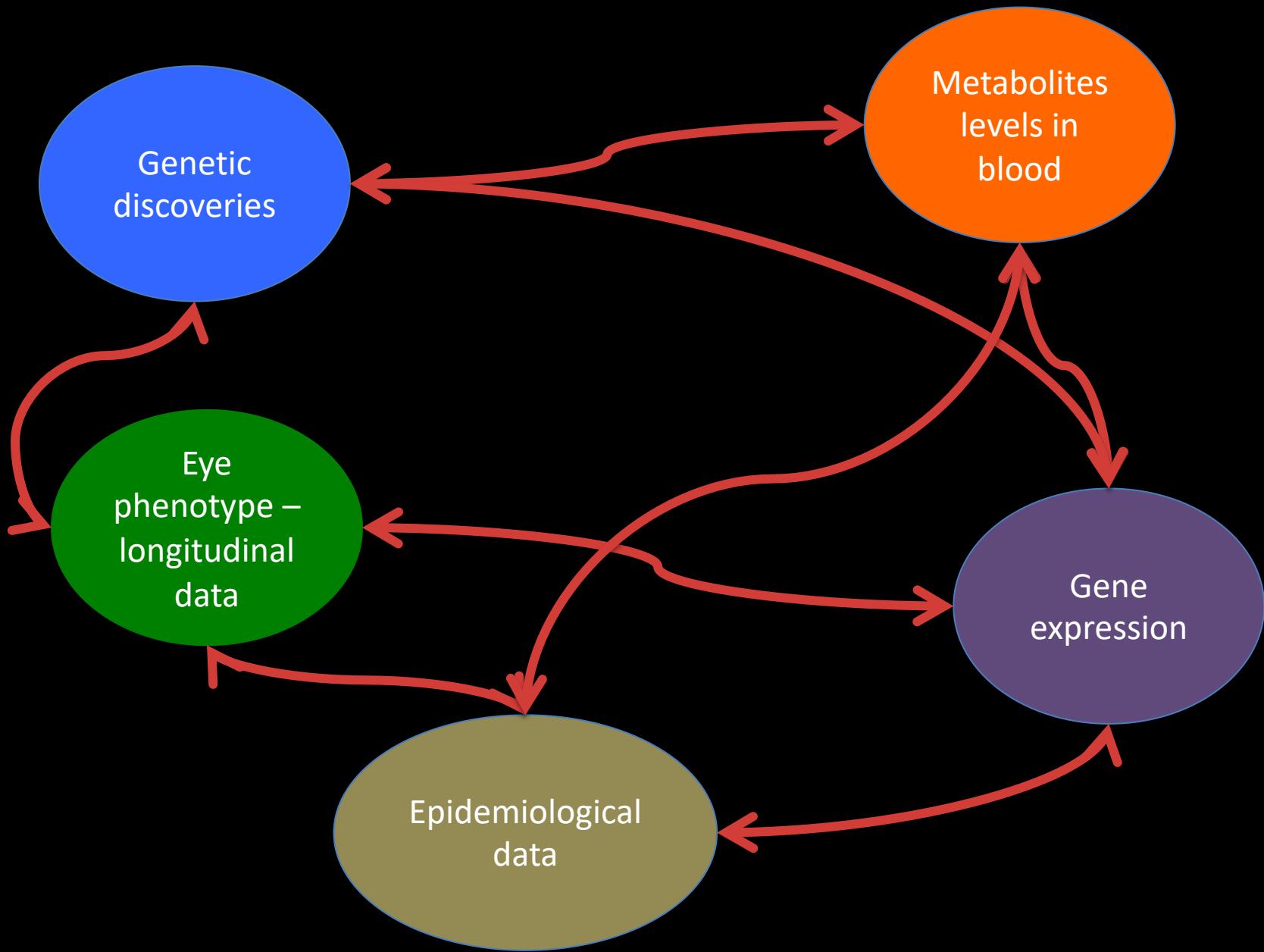


Figure 12: Relative frequency by affection group.



Conclusions

Conclusions

- Never be afraid of asking for something

Conclusions

- Never be afraid of asking for something
- Study as better as you can

Conclusions

- Never be afraid of asking for something
- Study as better as you can
- Learn R!

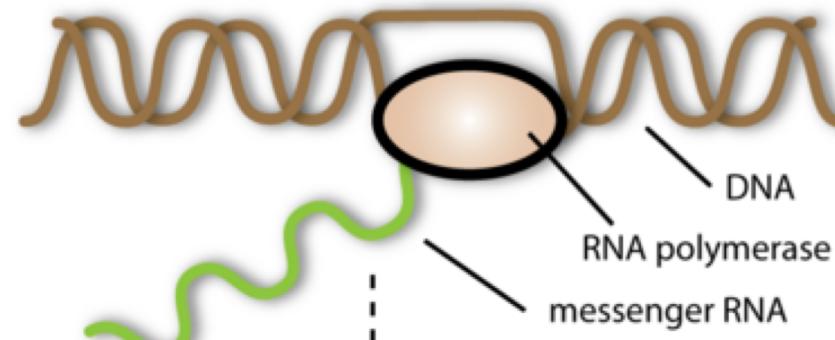
Conclusions

- Never be afraid of asking for something
- Study as better as you can
- Learn R!
- Be curious

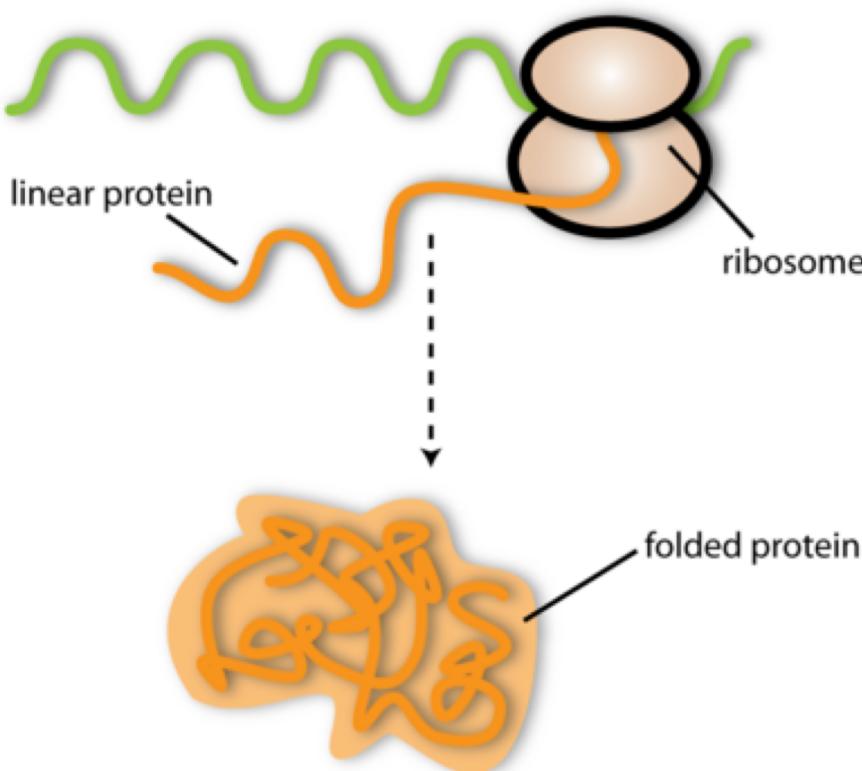
Thanks for listening!

Any question about
food and kangaroos?

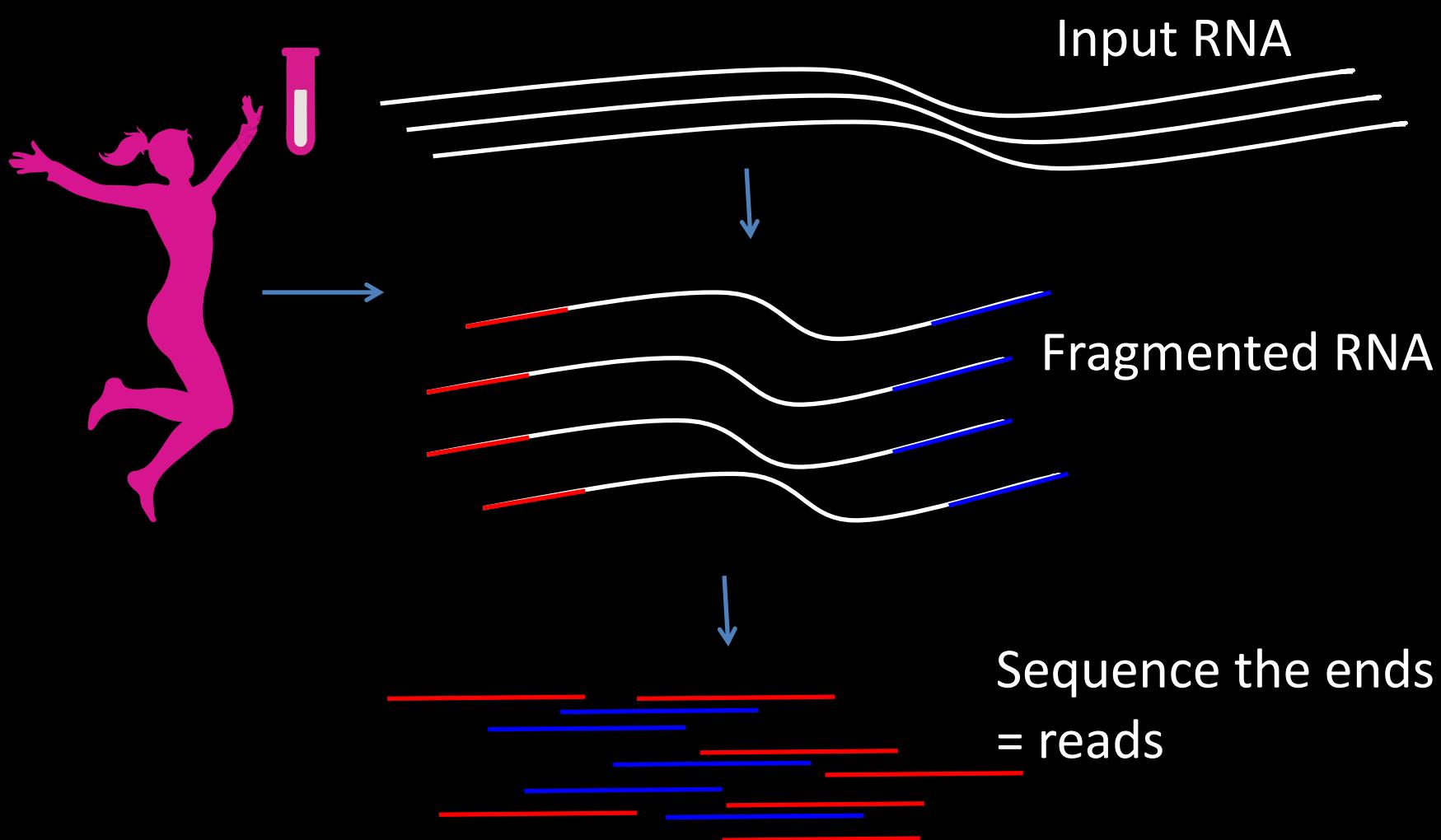
1. Transcription



2. Translation



RNA-sequencing





1. mRNA Isolation

2. Illumina Sequencing

3. Align Sequences against Genome

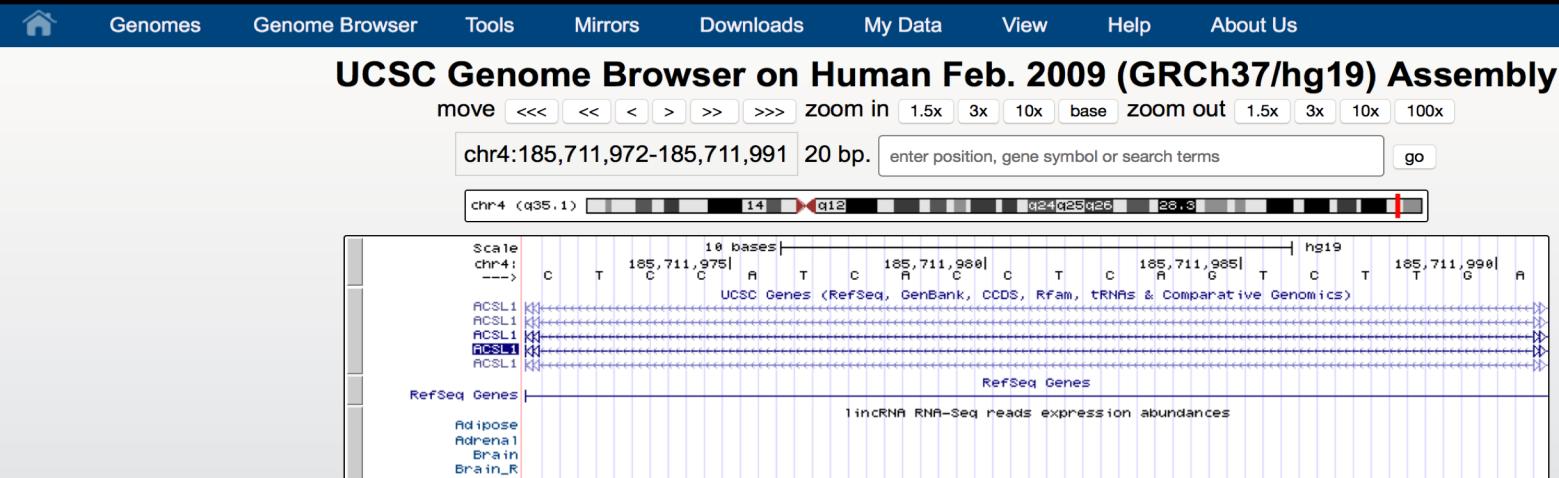


4. Generate Sequence Counts for all Genes in Genome

Gene A: $\frac{30}{10} = 3 \text{ fold change}$

Gene B: $\frac{10}{5} = 2 \text{ fold change}$

Reference genomes



[genome browser](#)

RNA analysis

	Treat1	Treat2	Control1	Control2
Gene1	100	120	20	15
Gene2	150	80	15	15
Gene3	0	15	80	100
Gene4	0	0	0	0
Gene5	180	200	250	200

Table of analysis

	Cases1	Cases2	Control1	Control2
Gene1	100	120	20	15
Gene2	150	80	15	15
Gene3	0	15	80	100
Gene4	0	0	0	0
Gene5	180	200	250	200

Table of analysis

	Treat1	Treat2	Control1	Control2
Gene1	100	120	20	15
Gene2	150	80	15	15
Gene3	0	15	80	100
Gene4	0	0	0	0
Gene5	180	200	250	200

Table of analysis

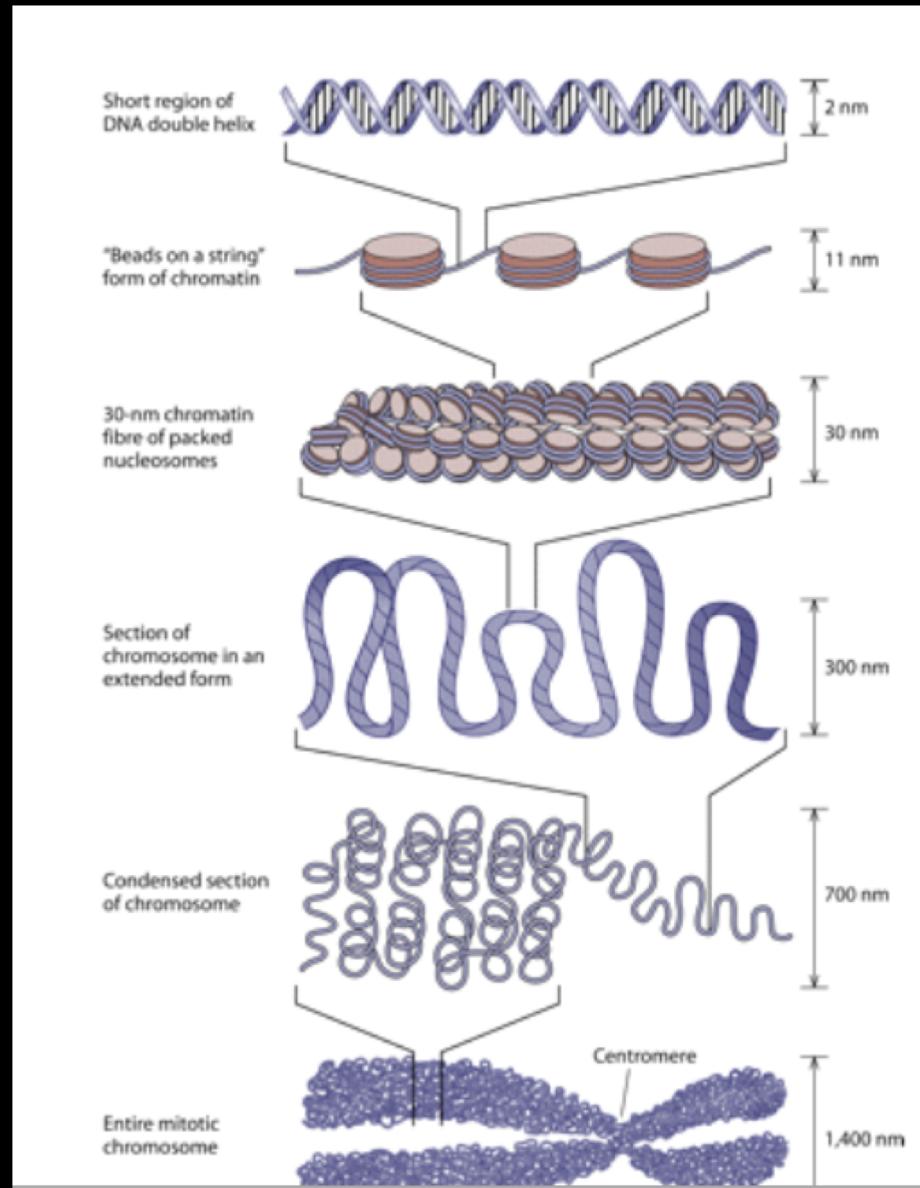
	Treat1	Treat2	Control1	Control2
Gene1	100	120	20	15
Gene2	150	80	15	15
Gene3	0	15	80	100
Gene4	0	0	0	0
Gene5	180	200	250	200

Protein-DNA interaction

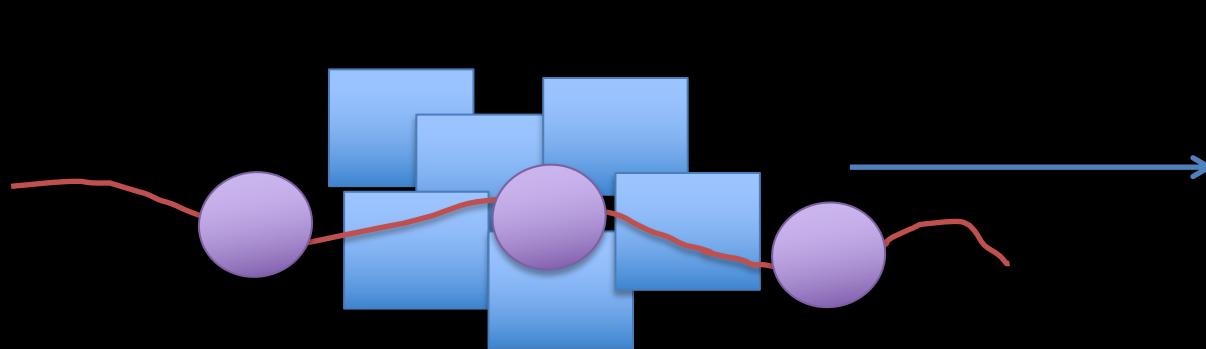
[Epigenetics - Ted Talk1](#)

[Epigenetic - Ted Talk2](#)

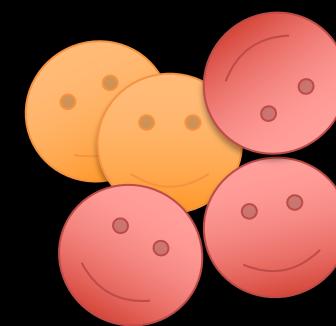
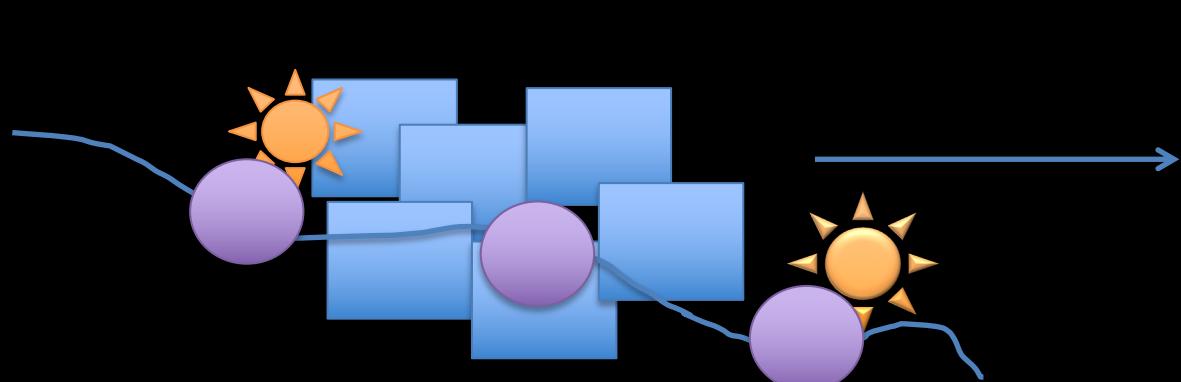
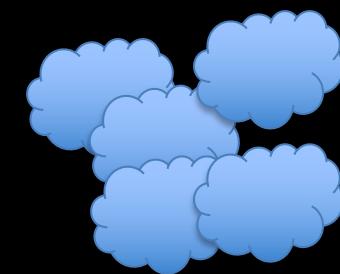
Protein-DNA interaction



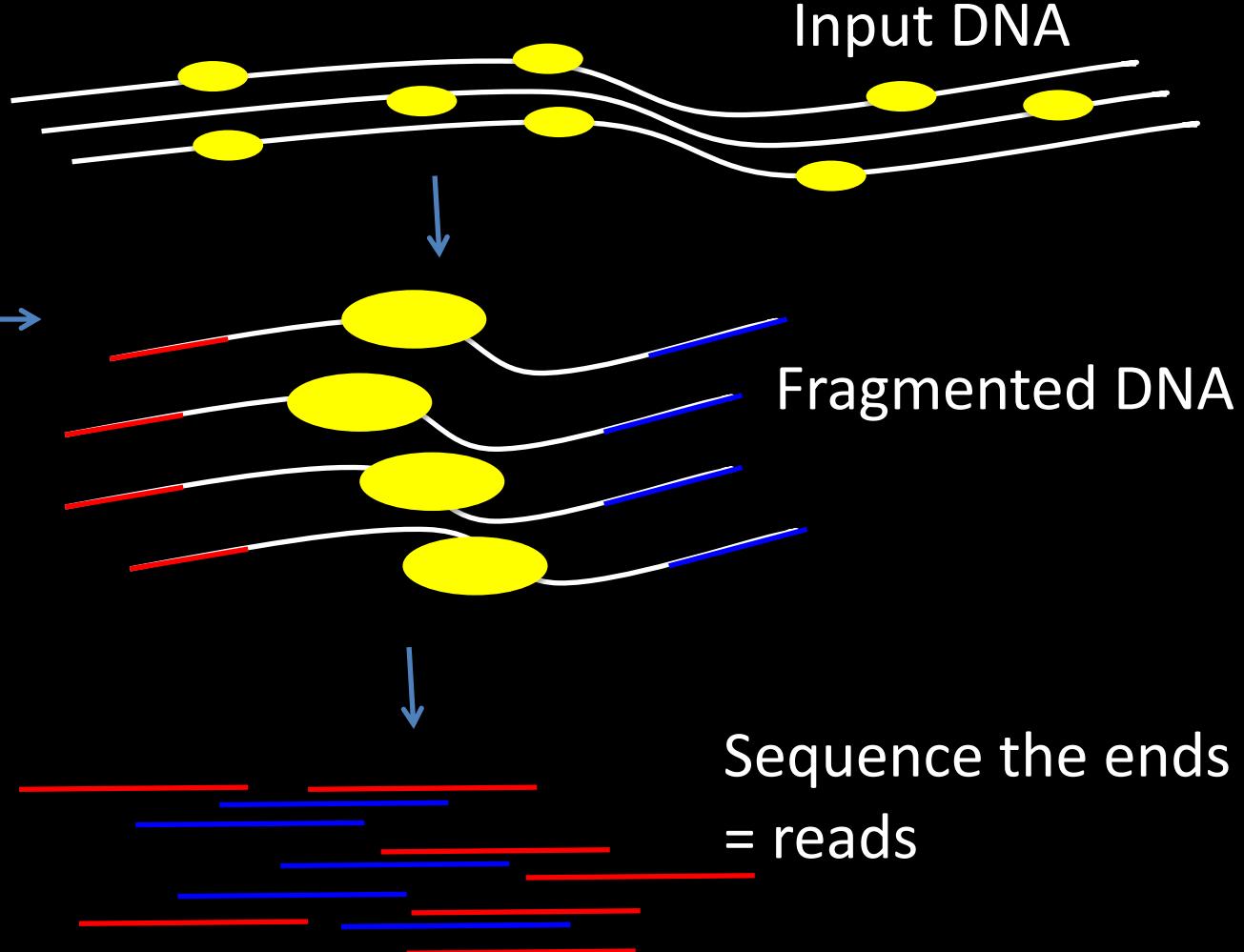
Initial cells



Differentiated cells

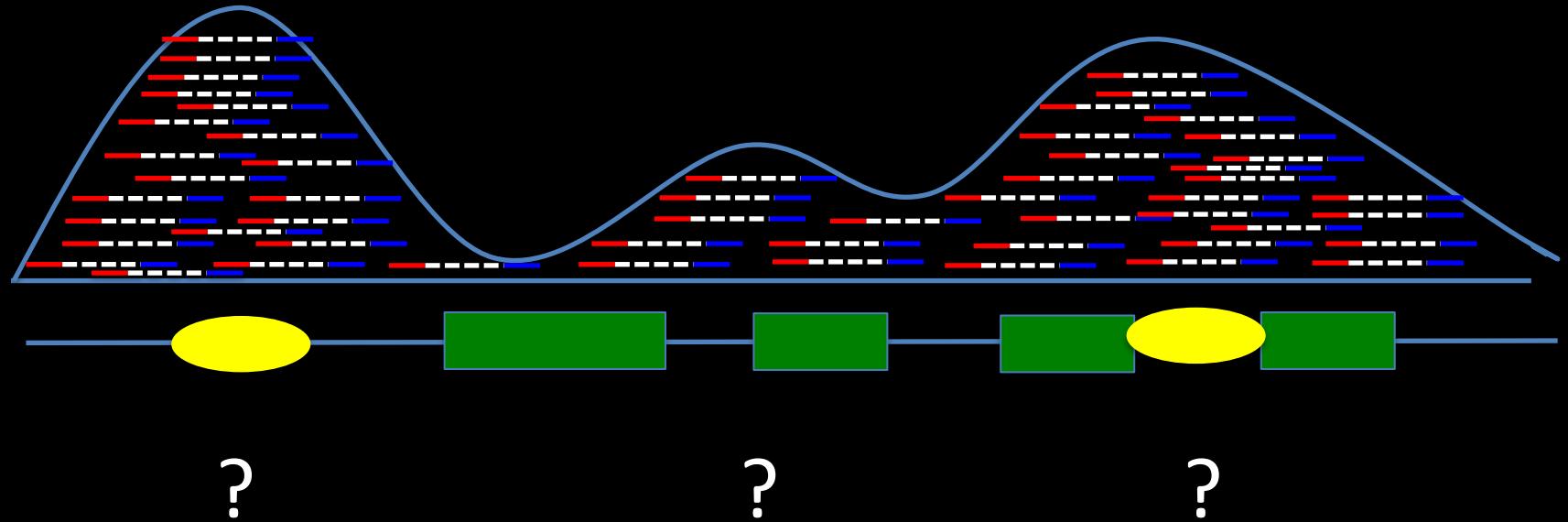


ChIP-sequencing



ChIP-Seq data

Objective: separate noise from true signal

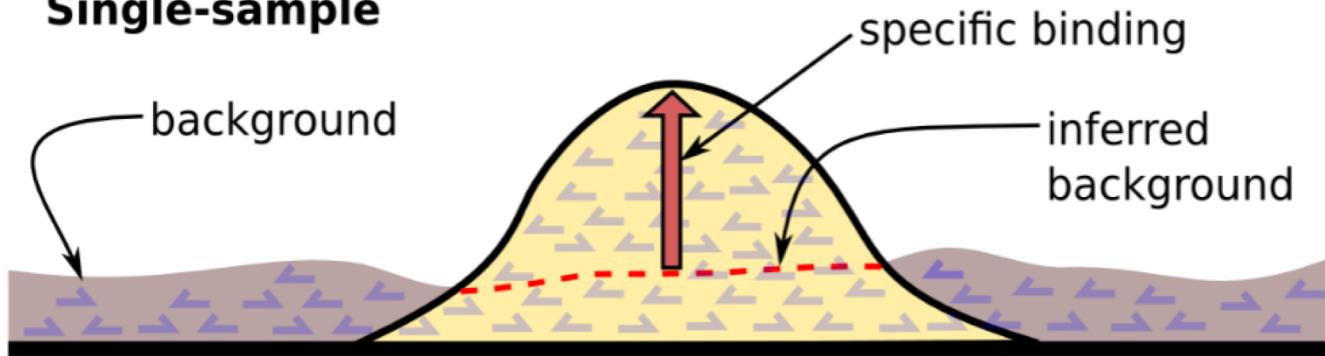


Noise also induced by the biological steps for preparation

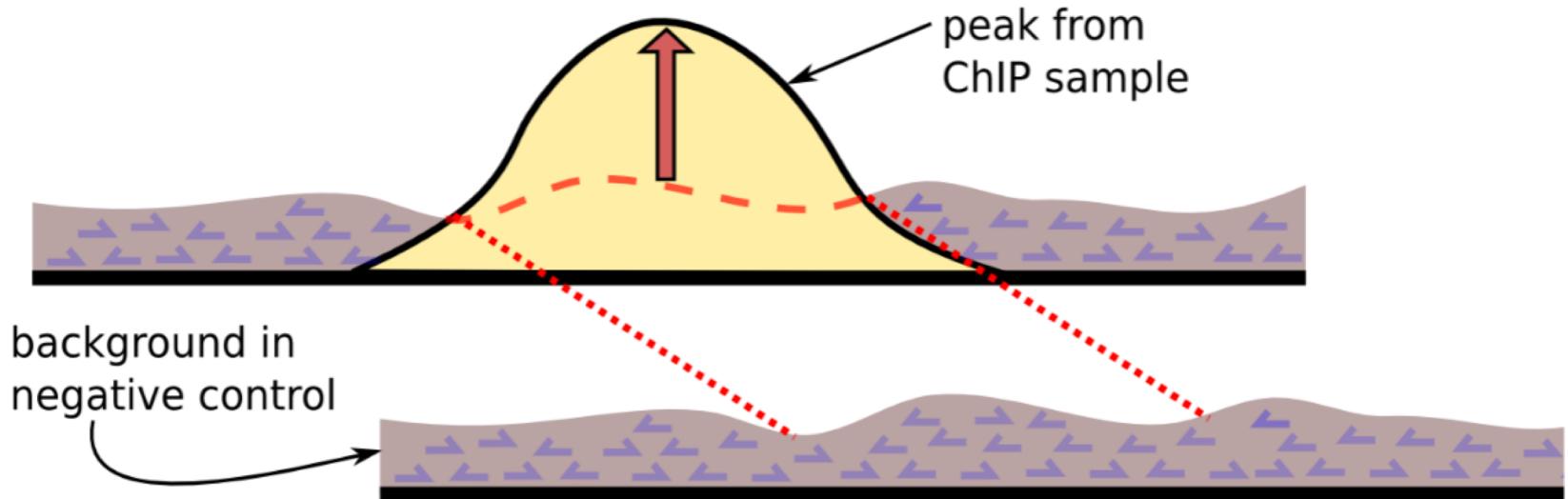
Conventional analyses for ChIP-seq

Detect absolute binding by peak calling:

Single-sample



Two-sample



ChIP-Seq most used software

- MACS2 – Liu lab

New from WEHI Bioinformatics

- csaw– Aaron Lun (Smyth lab)

Model counts

- Wide use of statistical model for counts
 - Poisson
 - Negative binomial
 - Binomial
- Multivariate analysis to study the noise
 - PCs
 - SVD

- Remove unwanted variation

Computational & statistical challenges

- Multivariate analysis
- Data mining
- R and C as support
- Communication with biologist with user-friendly form

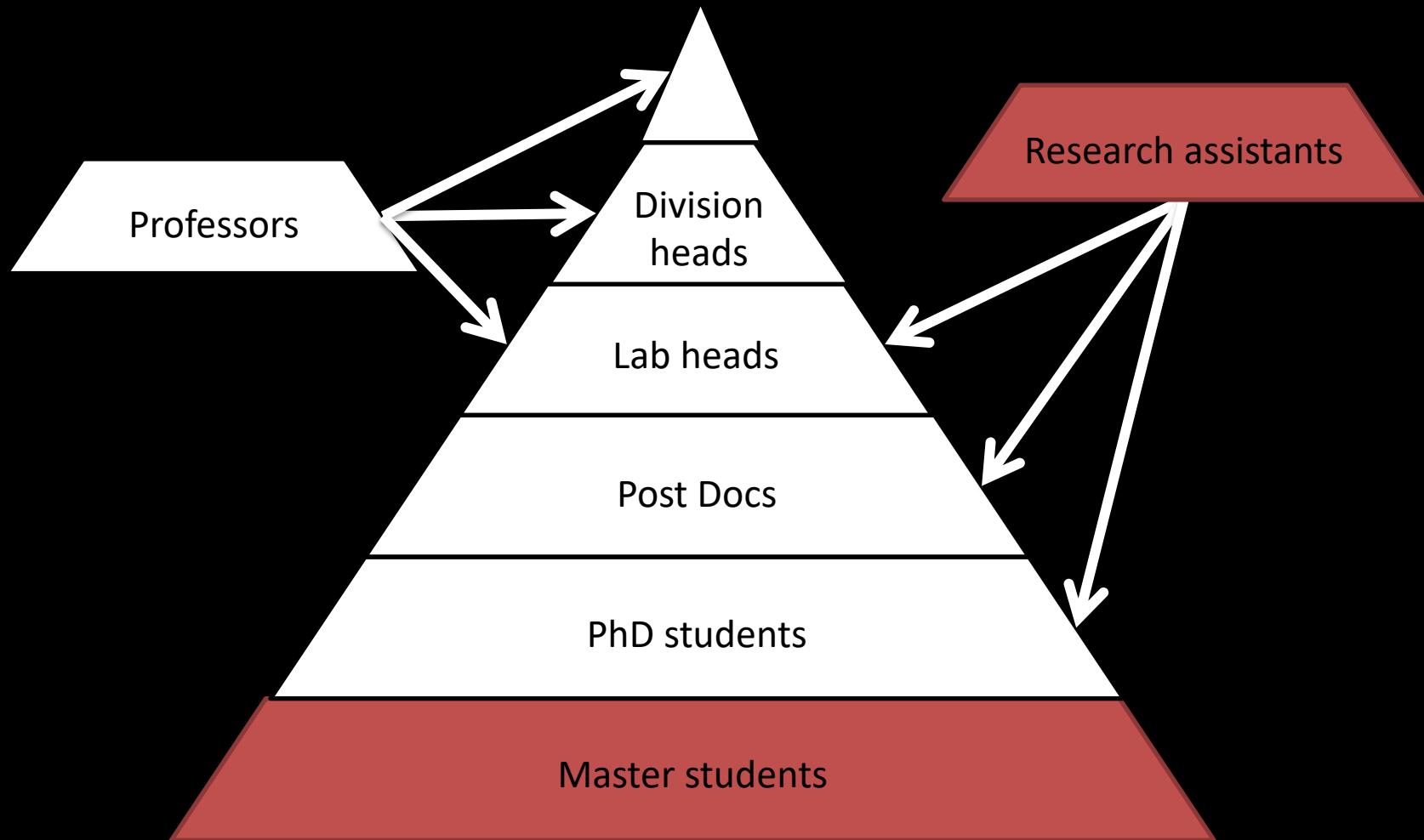
The dream

	Cases1	Cases2	Control1	Control2
Gene1	100	120	20	15
Gene2	150	80	15	15
Gene3	0	15	80	100
Gene4	0	0	0	0
Gene5	180	200	250	200



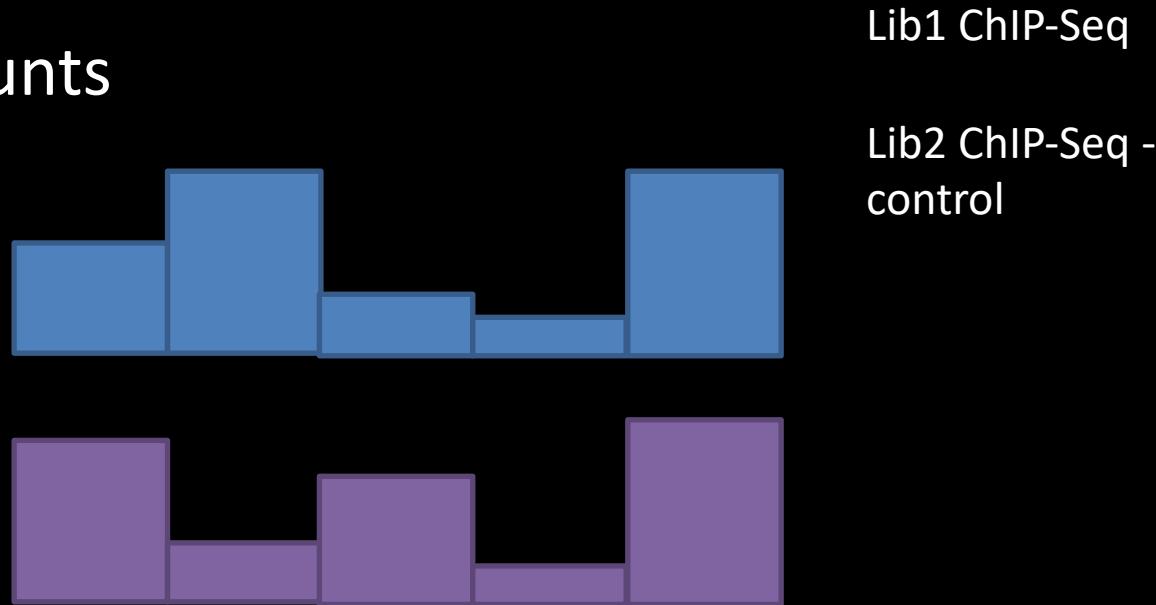
Genes with significant SNPs highlighted in the GWAS

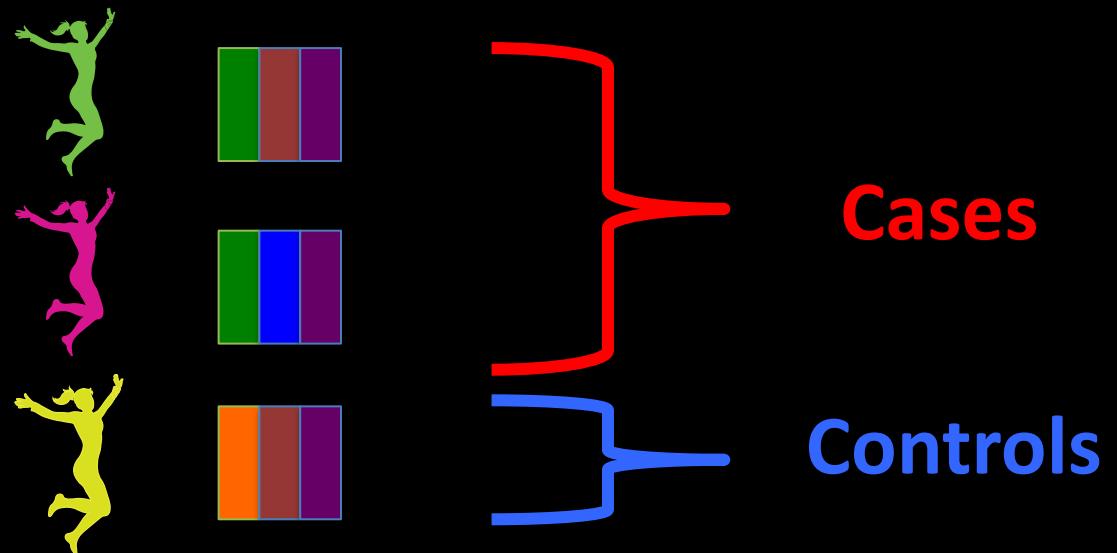
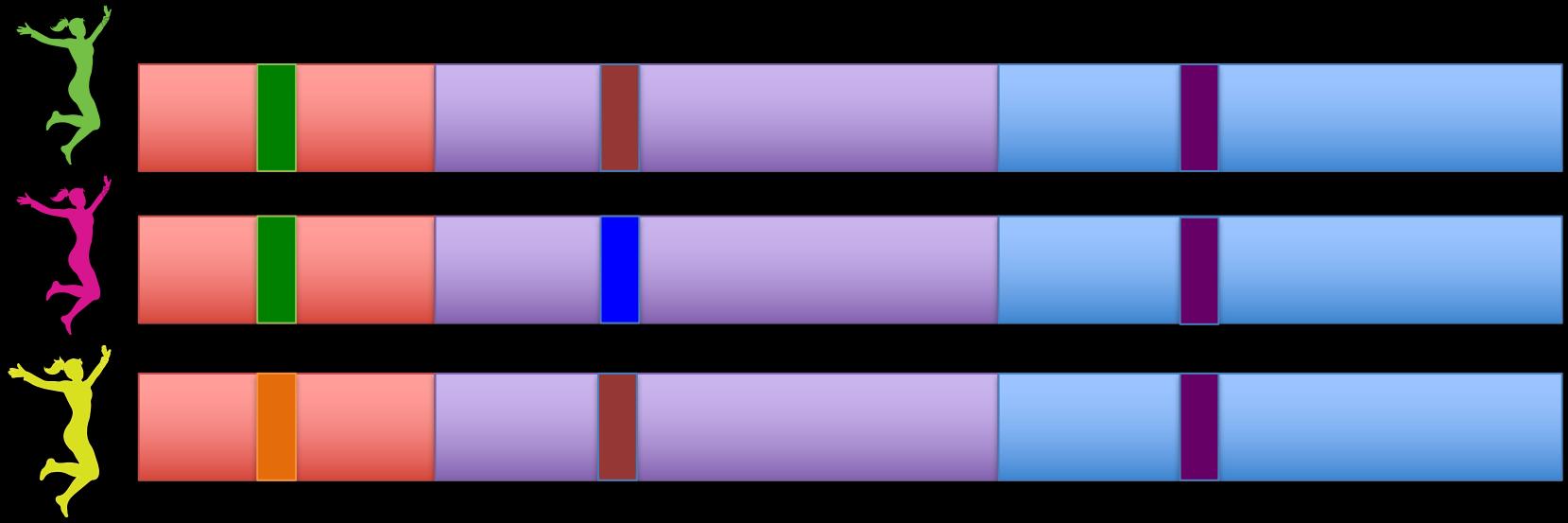
To work in a research institute

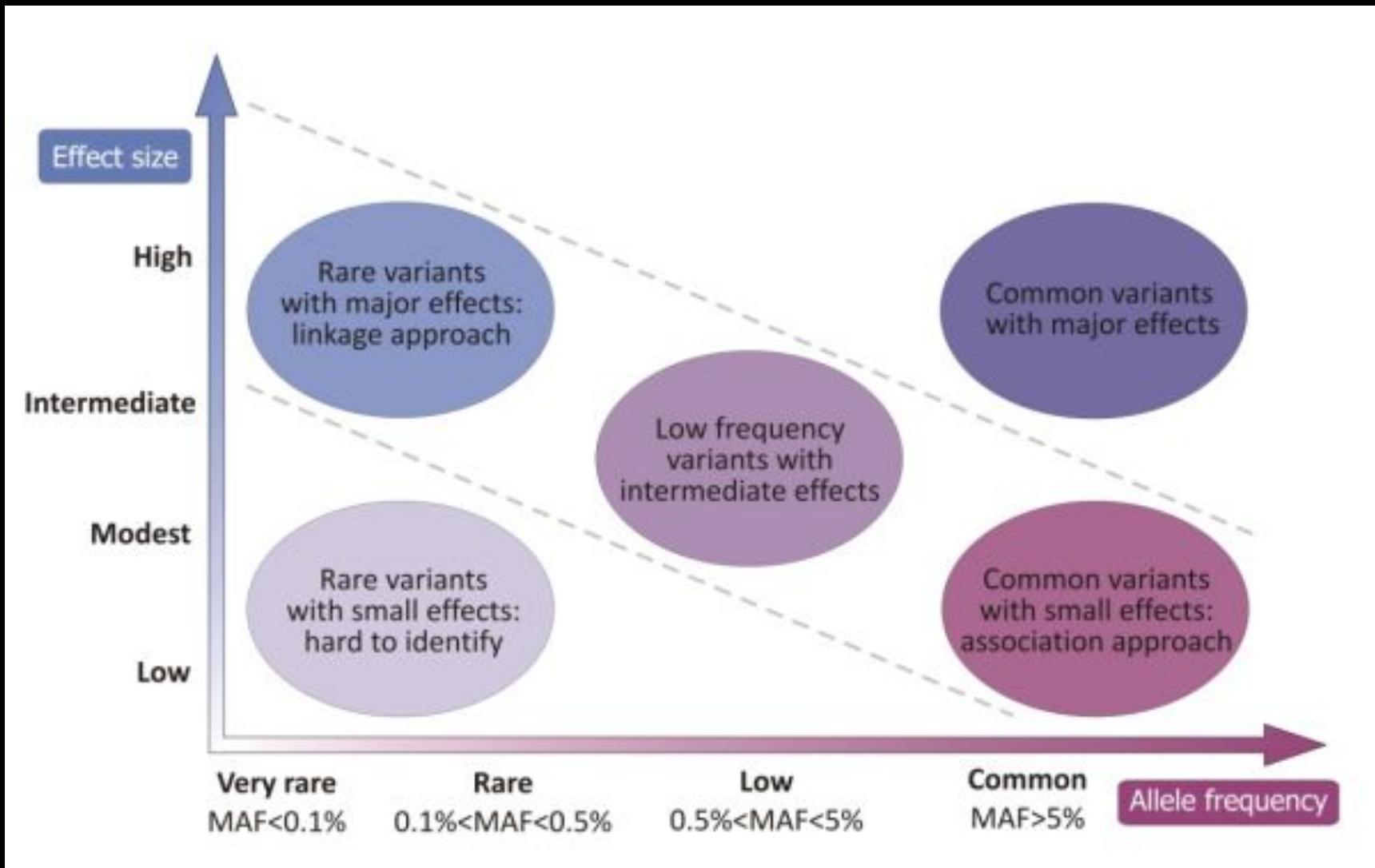


Classic ways of analysis

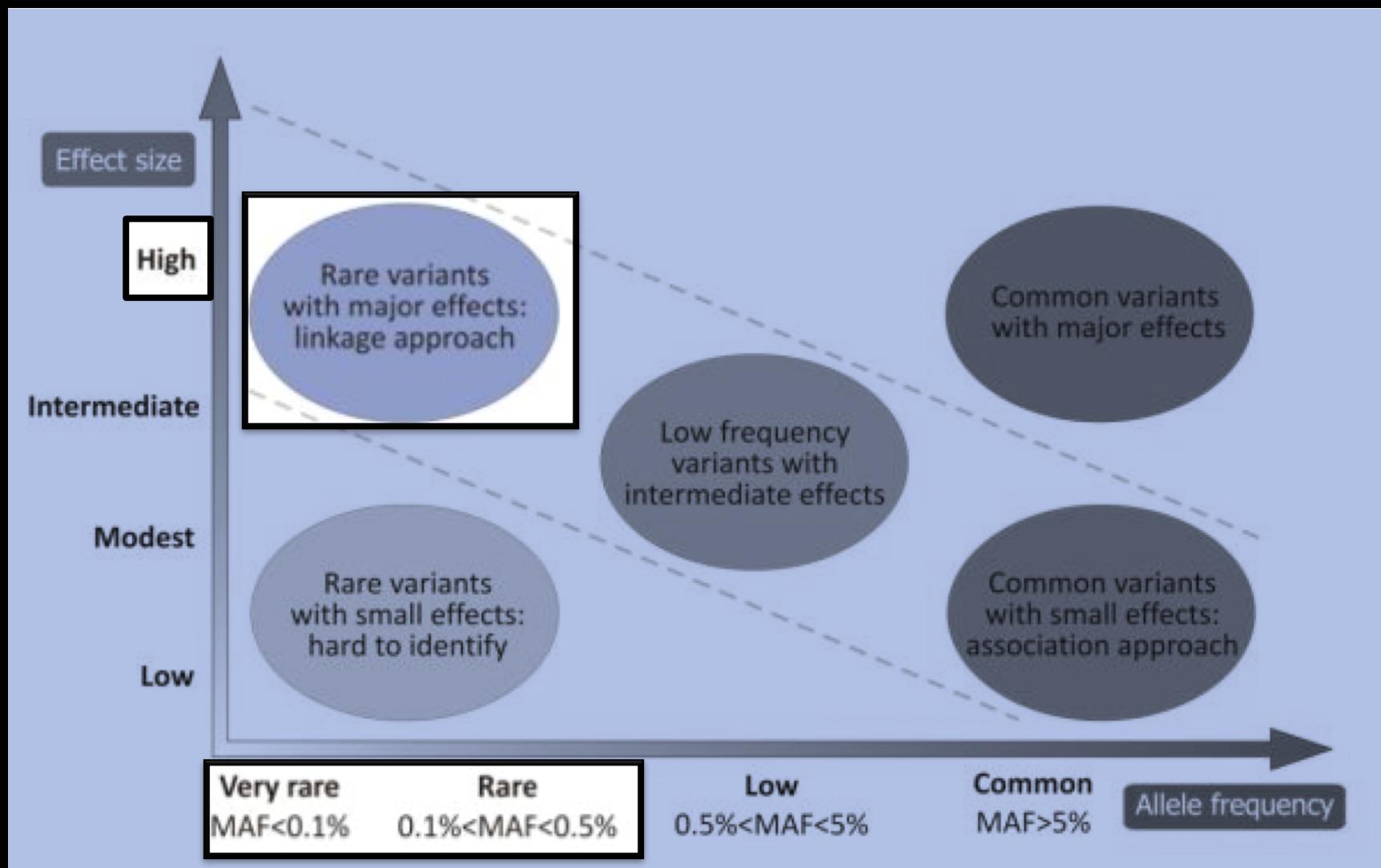
Compare
binned counts







Old approach to MacTel: poor results



New approach to MacTel: GWAS

