

Factors Influencing accuracy of genomic selection with sequence Information

Ignacy Misztal, Ivan Pocrnic*, Daniela Lourenco

University of Georgia

*now Roslin Institute



Genomic selection, effective population size and sequence data

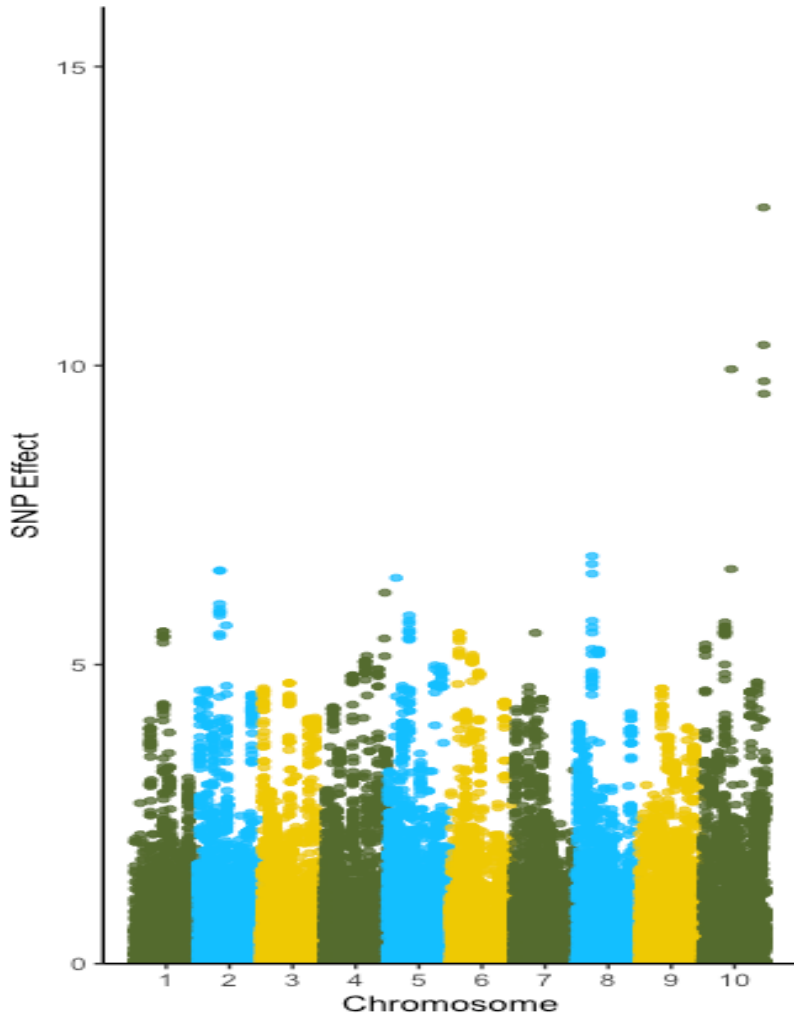
- Efficient genomic selection with sequence data - need QTN location and QTN variance (Fragomeni et al., 2017)
- GWAS developed mostly for human data – large effective population size ($N_e=3000-10000$)
- Farm populations have small N_e : 40 to 150
 - 5k to 15k independent chromosome segments (Stam, 1980; Pocrnic et al., 2016)
 - Each segment size 200k to 600k
- Questions
 - Does small effective population size affect GWAS?
 - What is Manhattan plot composed of?
 - Why small benefits of using sequence data for prediction?

Simulation study

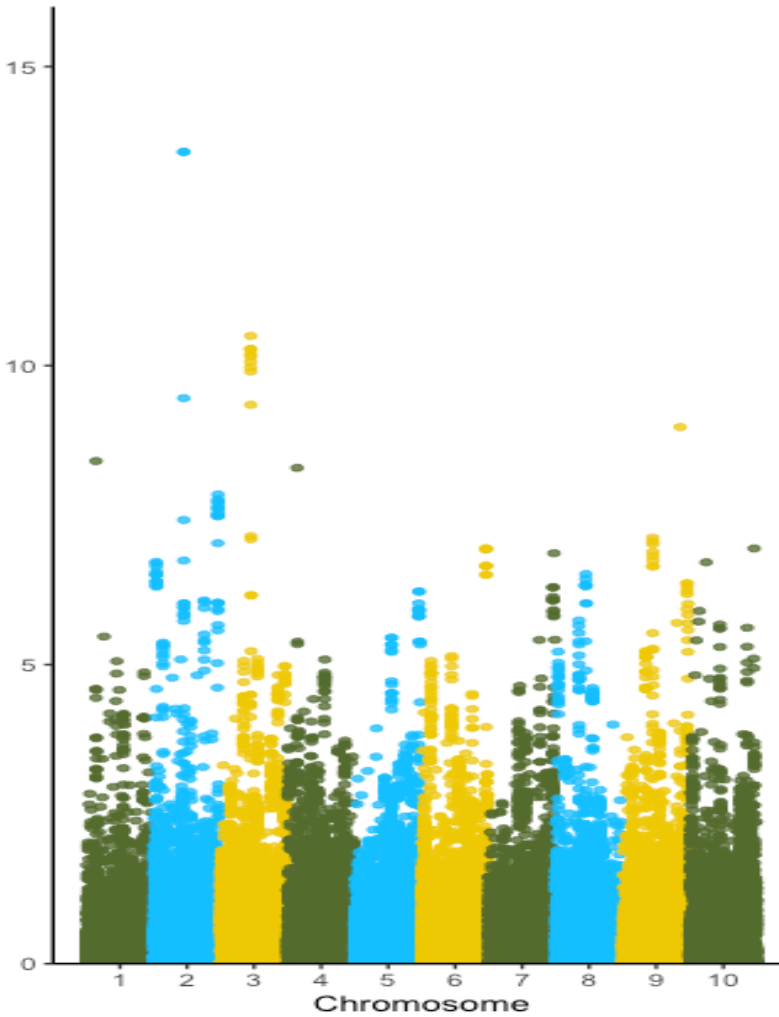
- AlphaSimR (Faux et al., 2016)
- 100 QTN with same effect and equally spaced
- 10 generations with 2000 animals in each
- Last 3 generations genotyped for 50k SNP – causative SNP included
- 3 data sets
 - Ne=60
 - Ne=600
 - Ne=60_3x --3 times more data (6000 per generation)
- Analyzes by ssGBLUP with p-value option

Manhattan plots for SNP effects

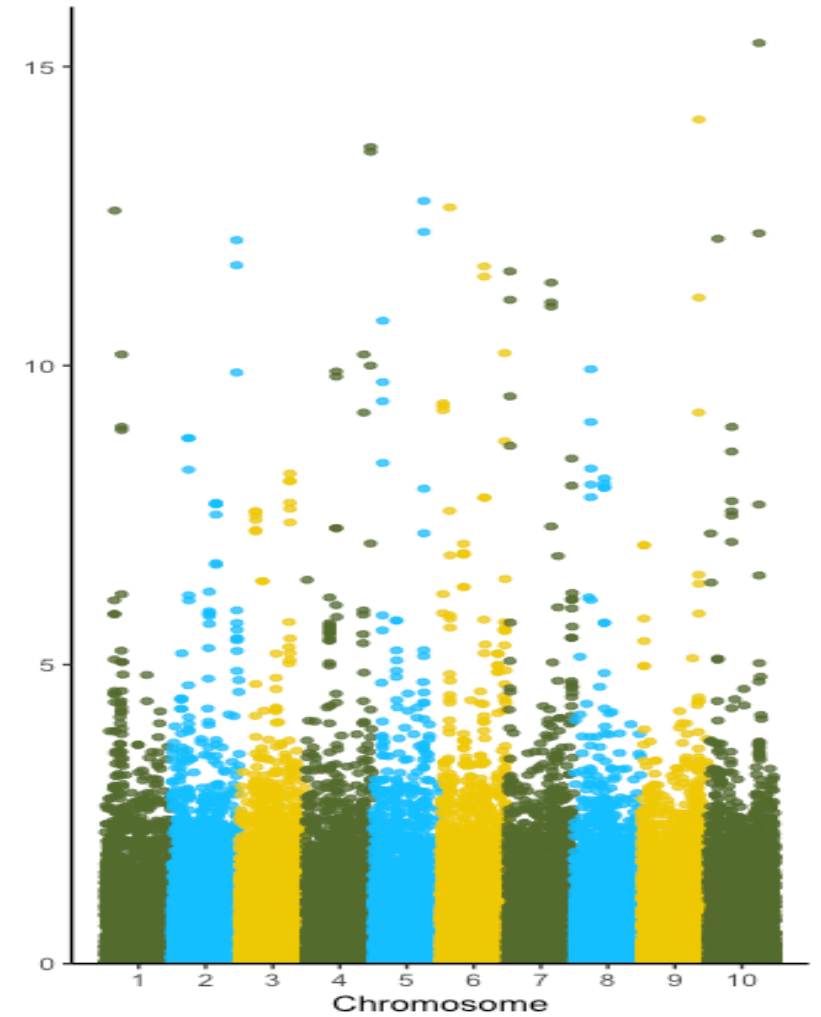
Ne=60



Ne=60 3X

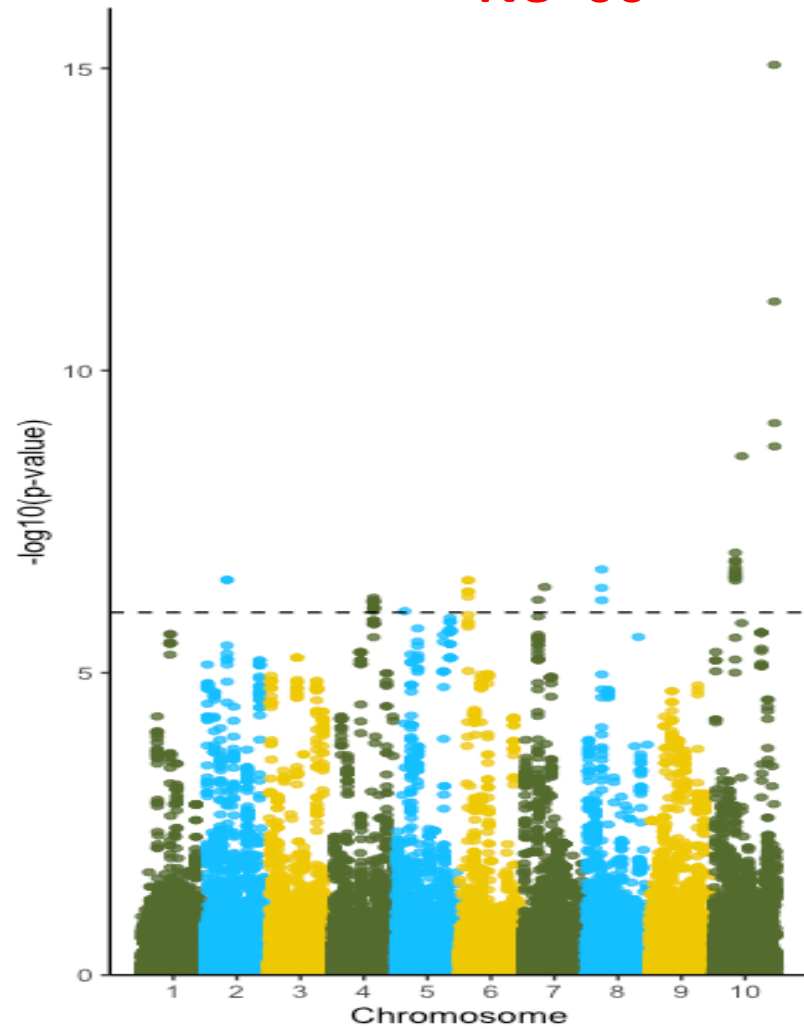


Ne=600

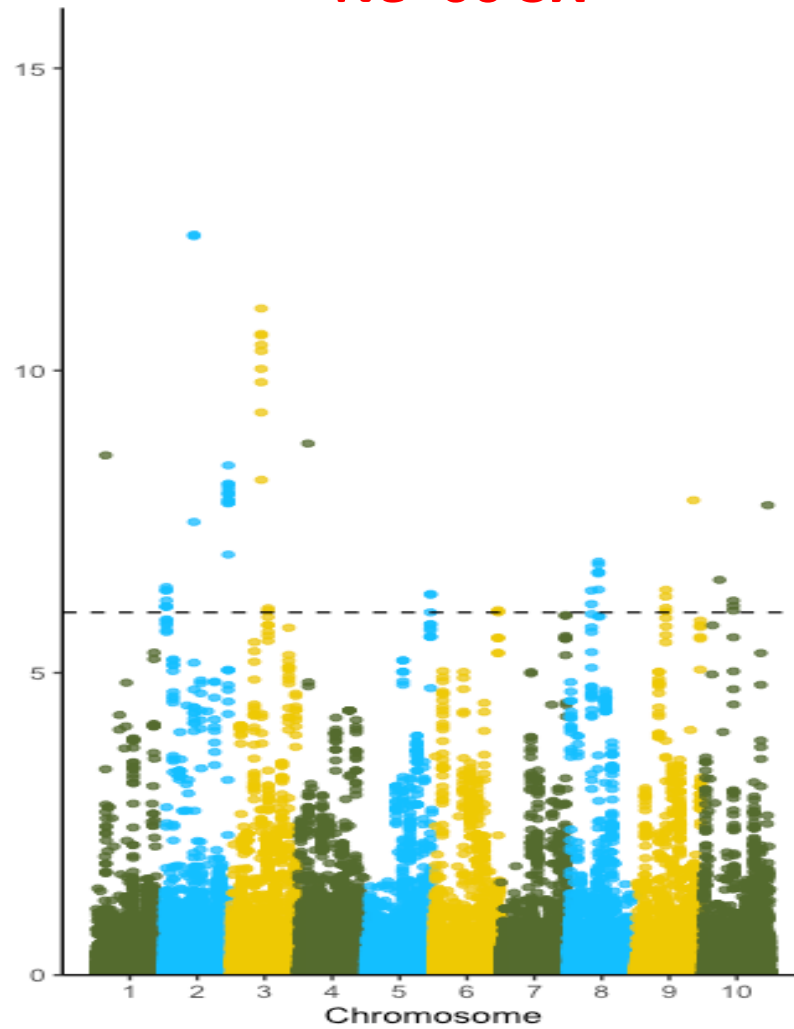


Manhattan plots for P-values

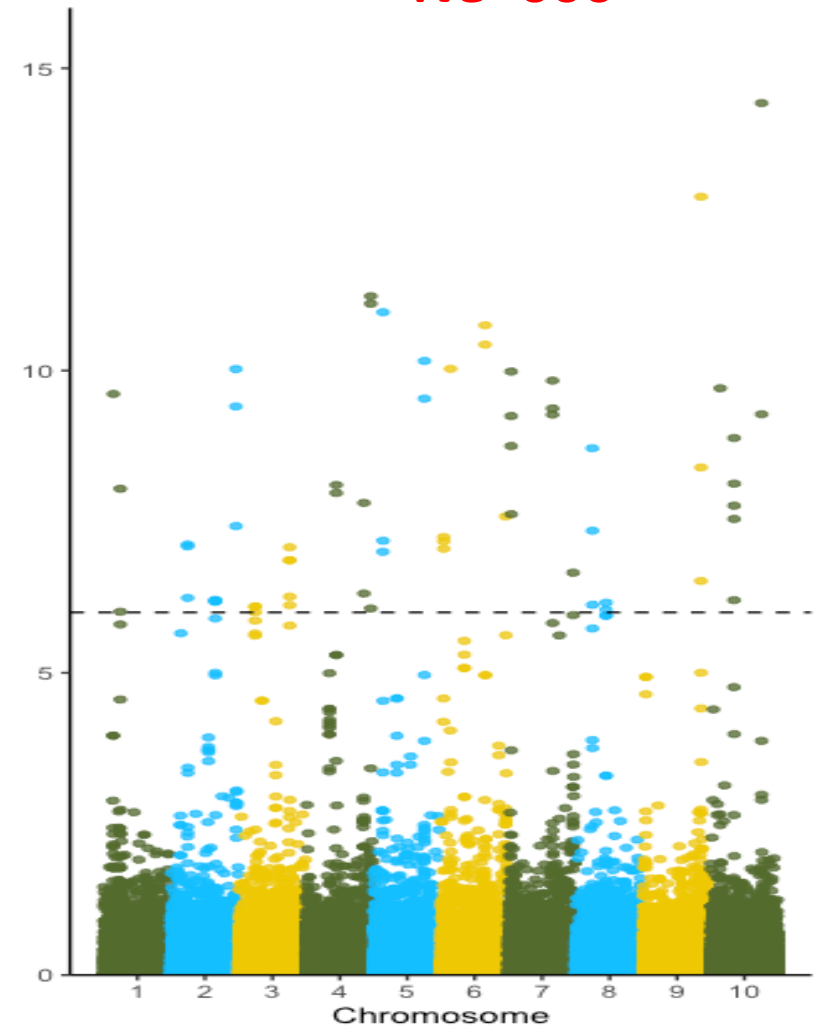
Ne=60



Ne=60 3X



Ne=600

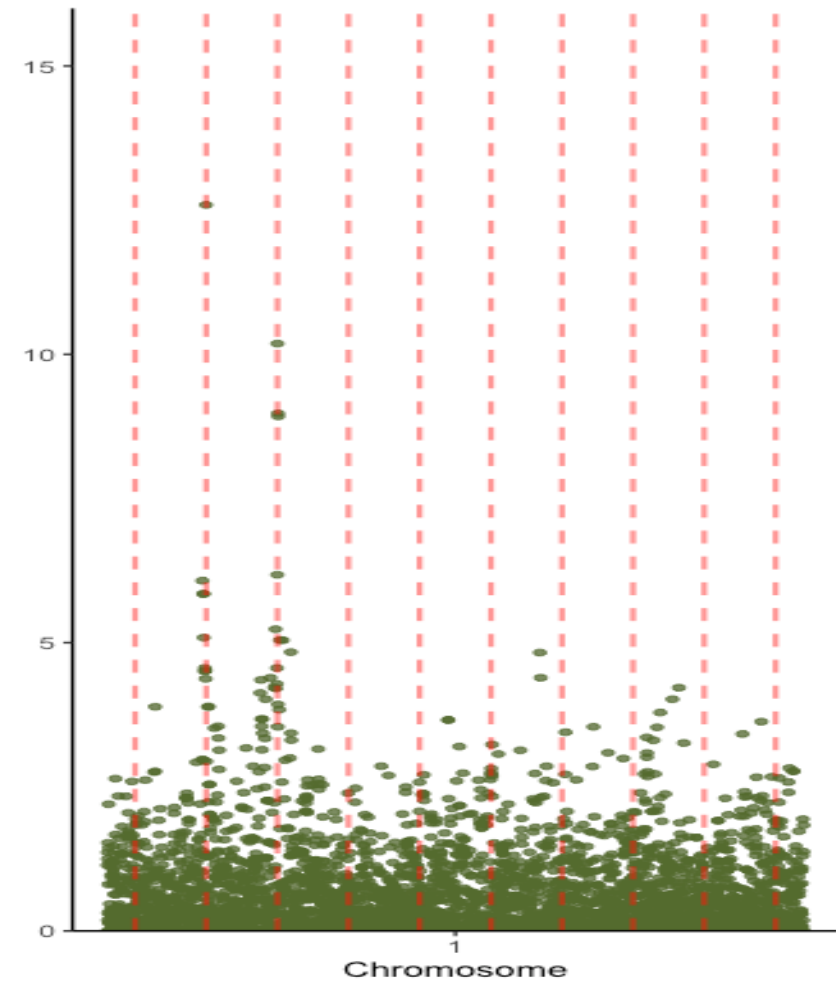
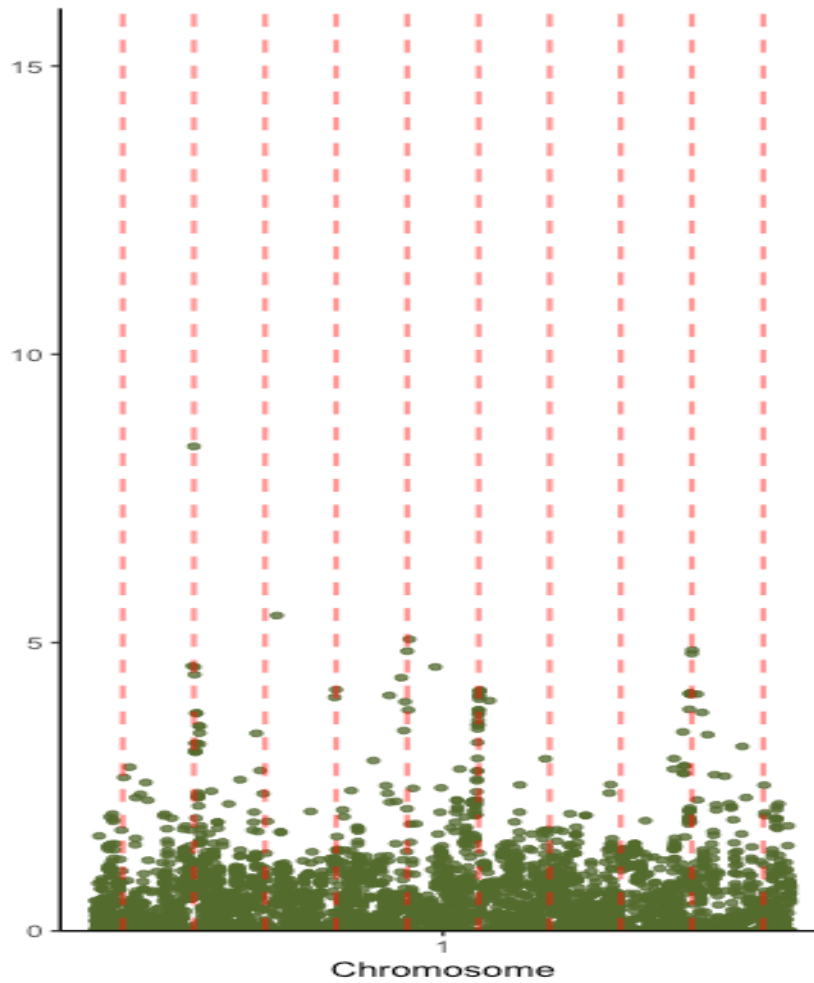
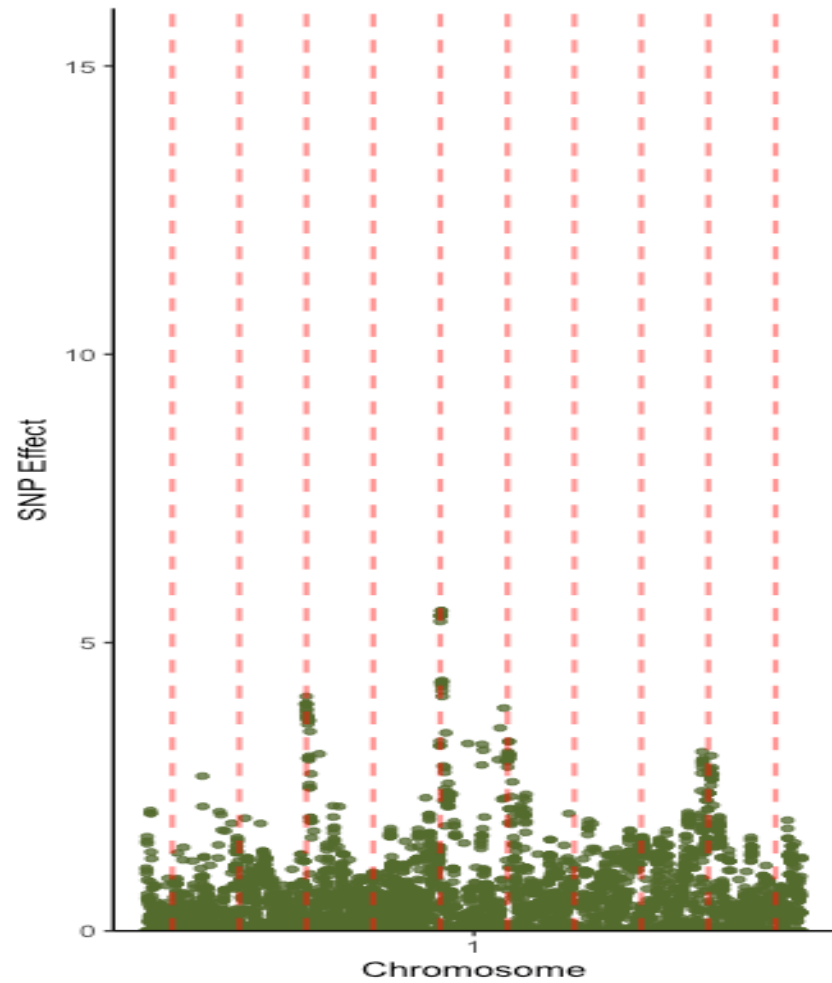


Manhattan plots for first chromosome

Ne=60

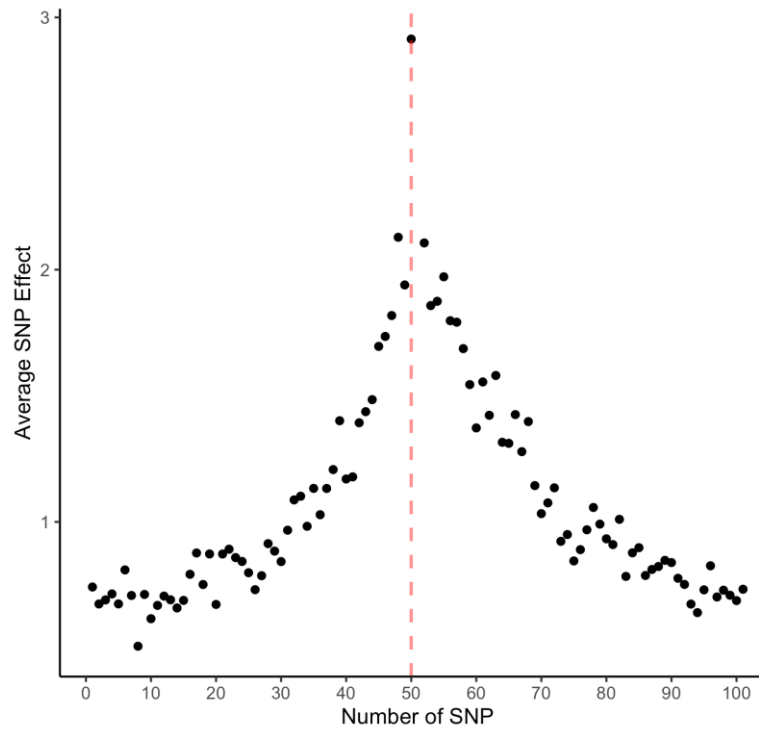
Ne=60 3X

Ne=600

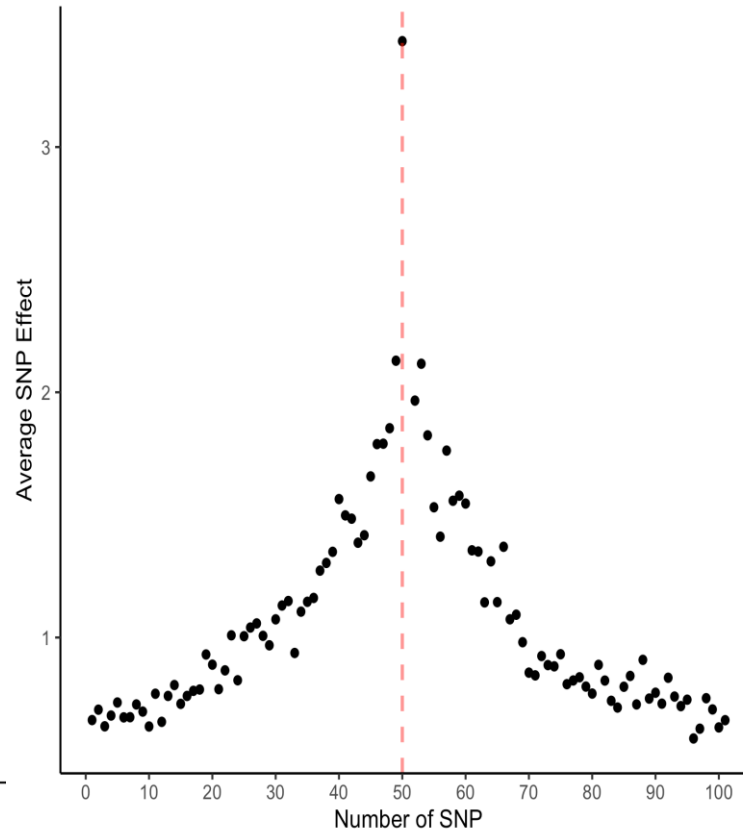


Plots averaged for all QTNs

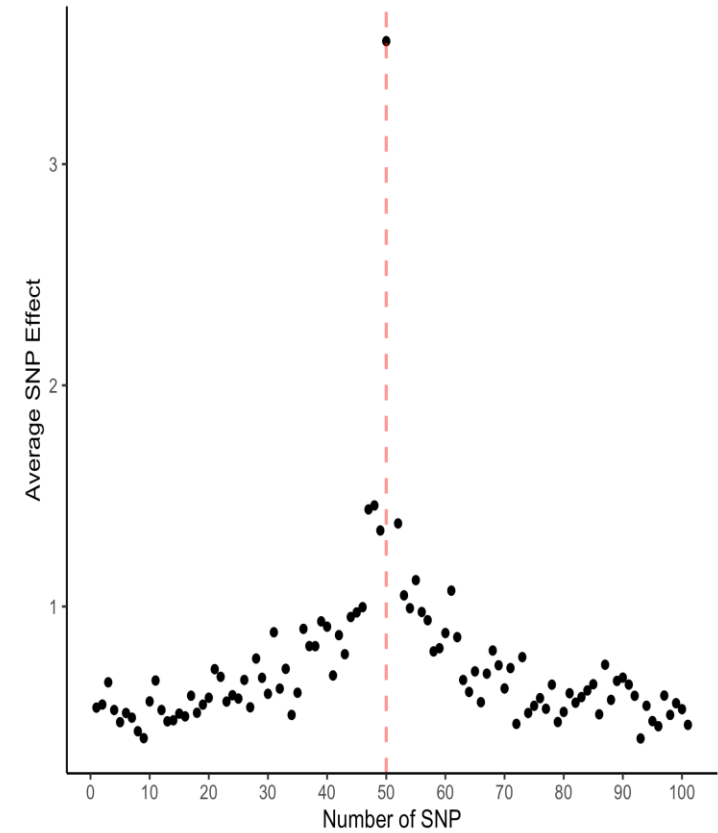
Ne=60



Ne=60 3X



Ne=600



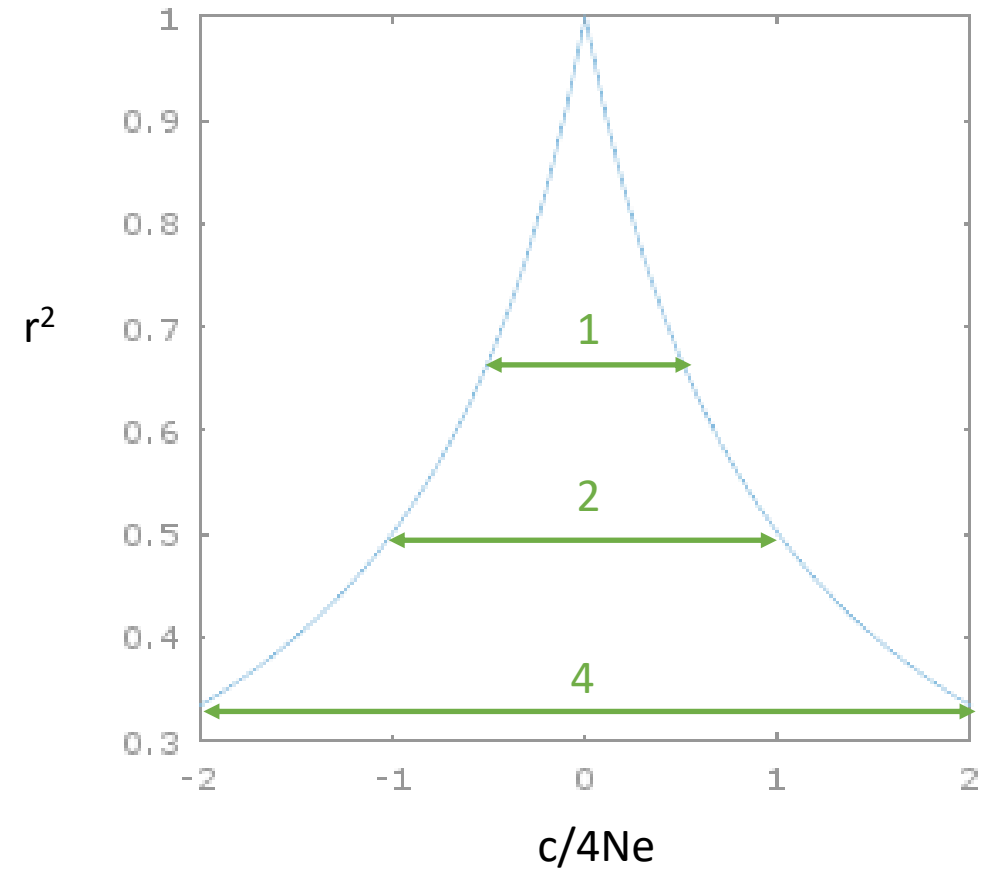
Pairwise linkage disequilibrium (r^2)

$$E(r^2) = 1 / (4 c N_e + 1) \quad \text{Sved (1971)}$$

c – distance between genes

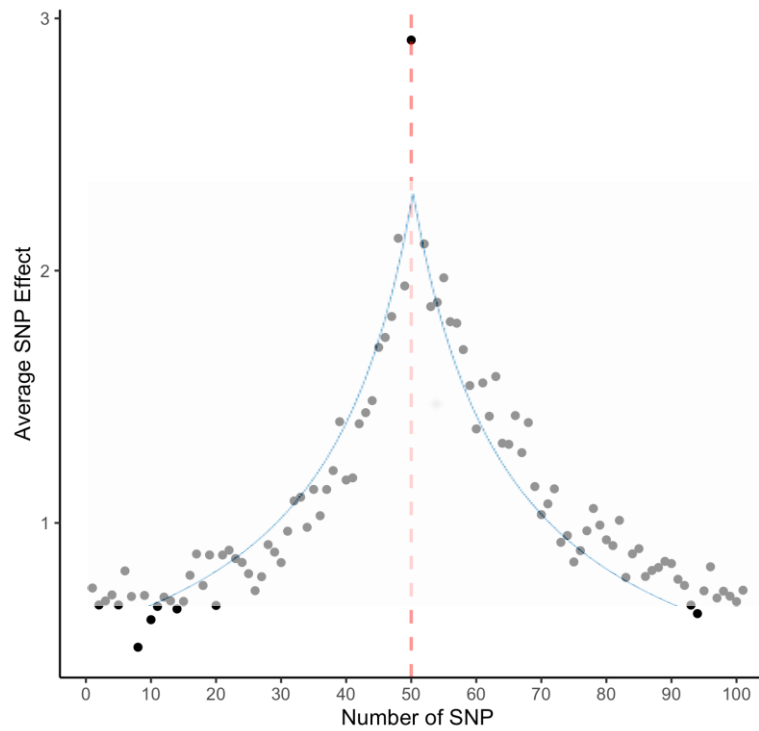
N_e – effective population size

One Stam segment - $(1/4N_e \text{ Morgan})$

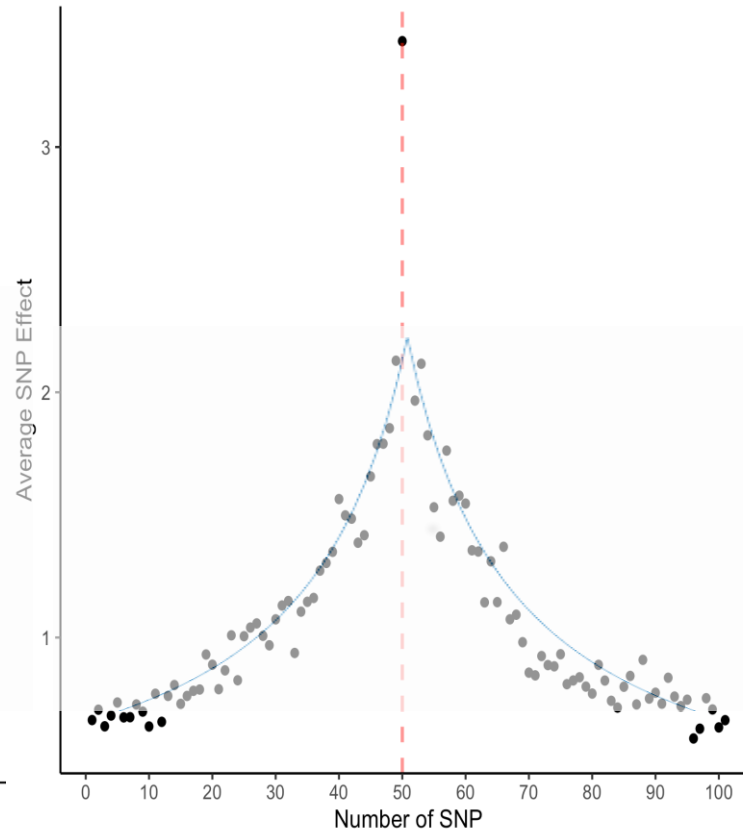


LD fitting....

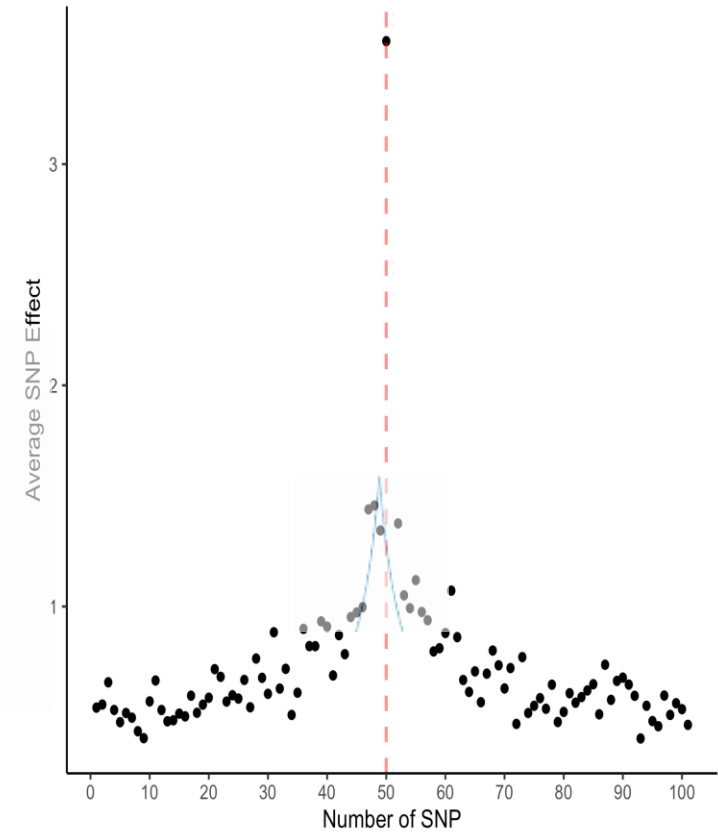
Ne=60



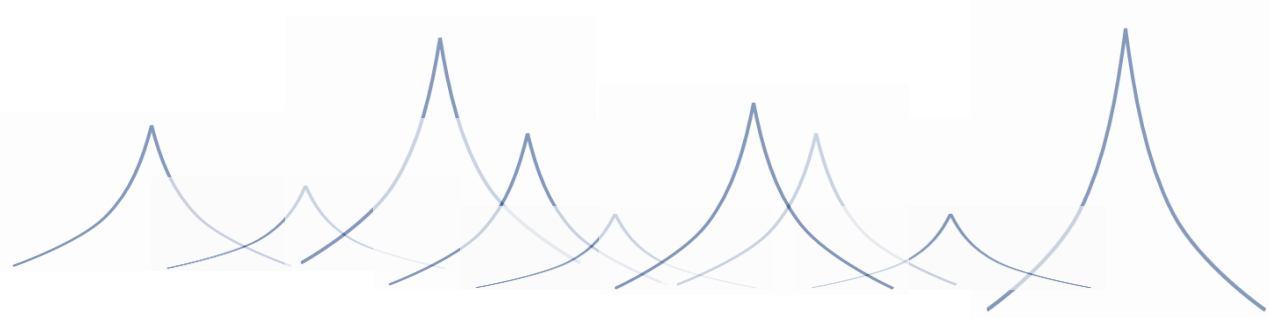
Ne=60 3X



Ne=600



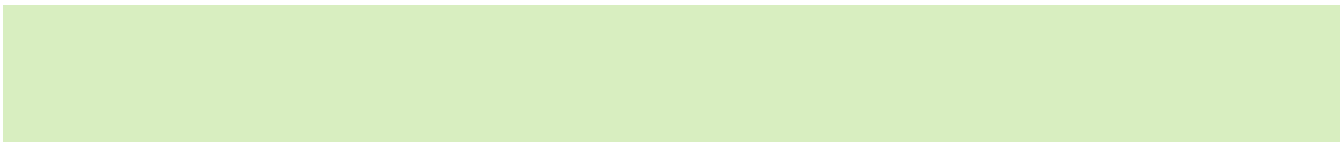
What is Manhattan plot composed of?



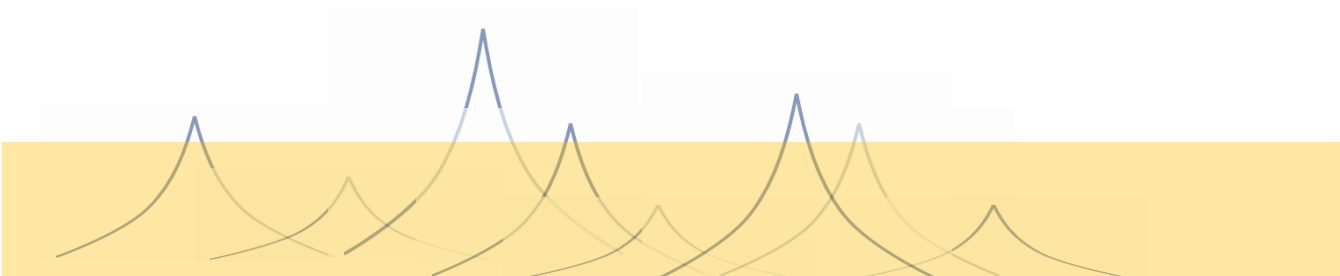
QTNs



Relationships

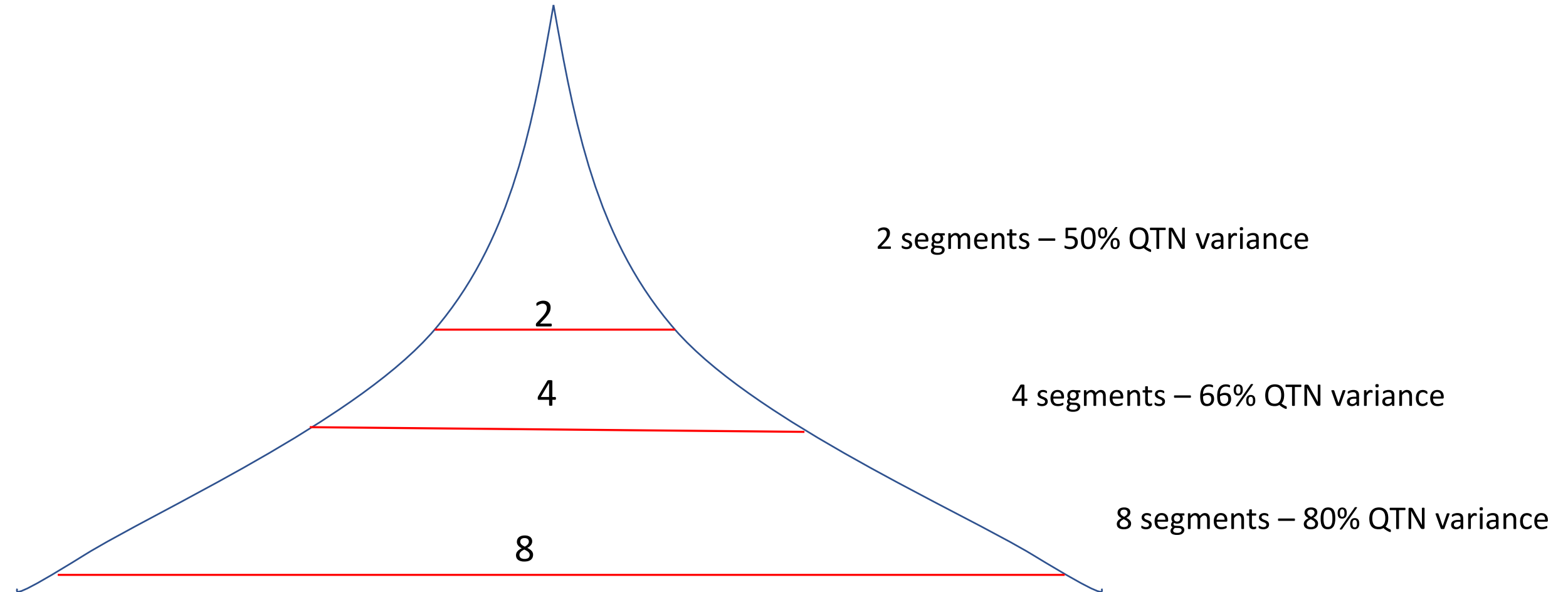


Noise



Composite plot

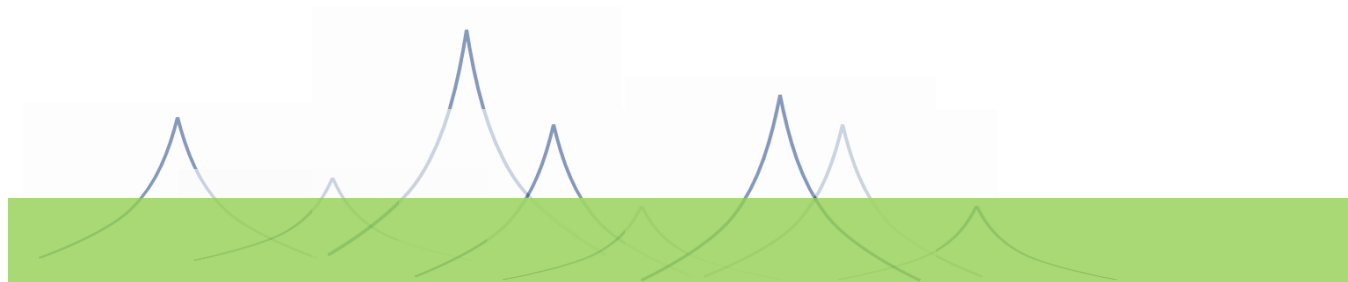
Proportion of QTN variance explained by different number of segments



Size of data



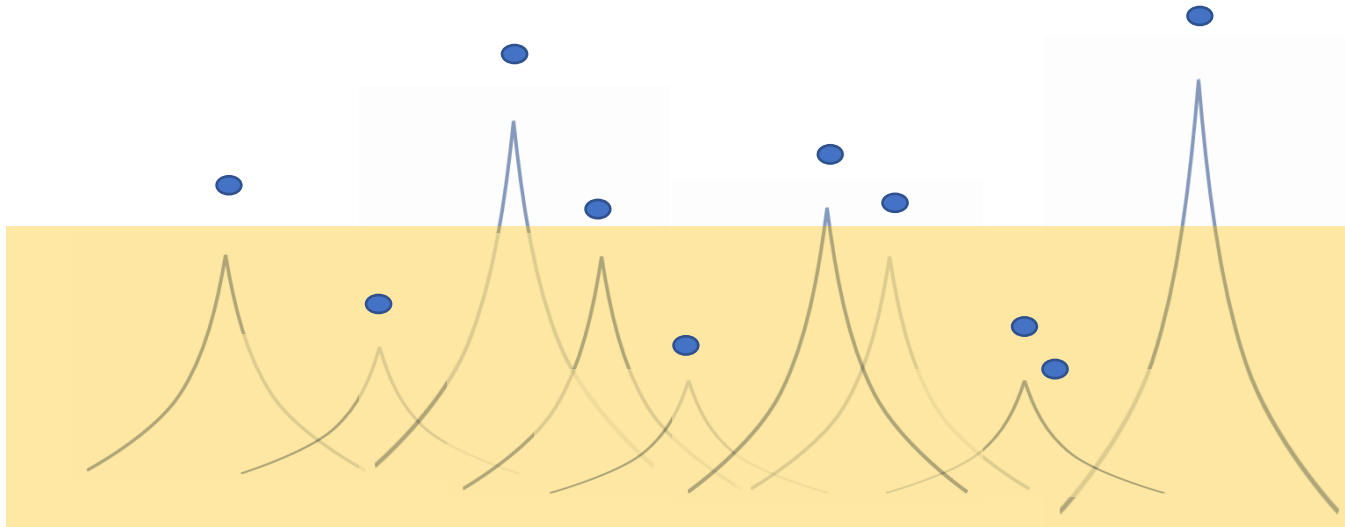
Small data



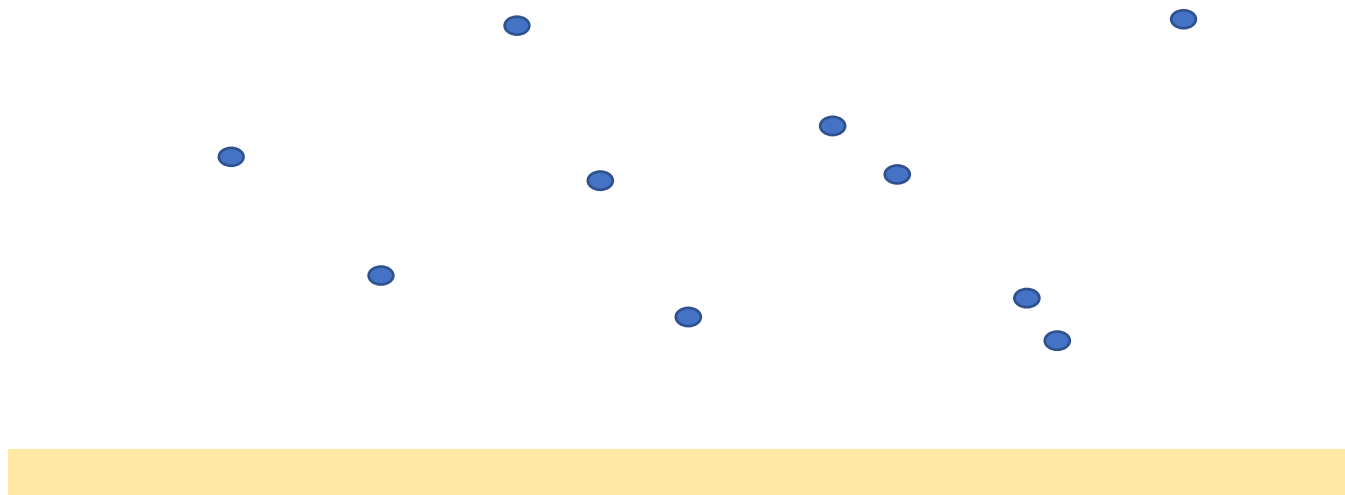
Large data

If GWAS by sliding window, optimal size 4-8 Stam segments
1-2 Mb cattle, 3-6 Mb pigs and chicken

+ QTN in the data

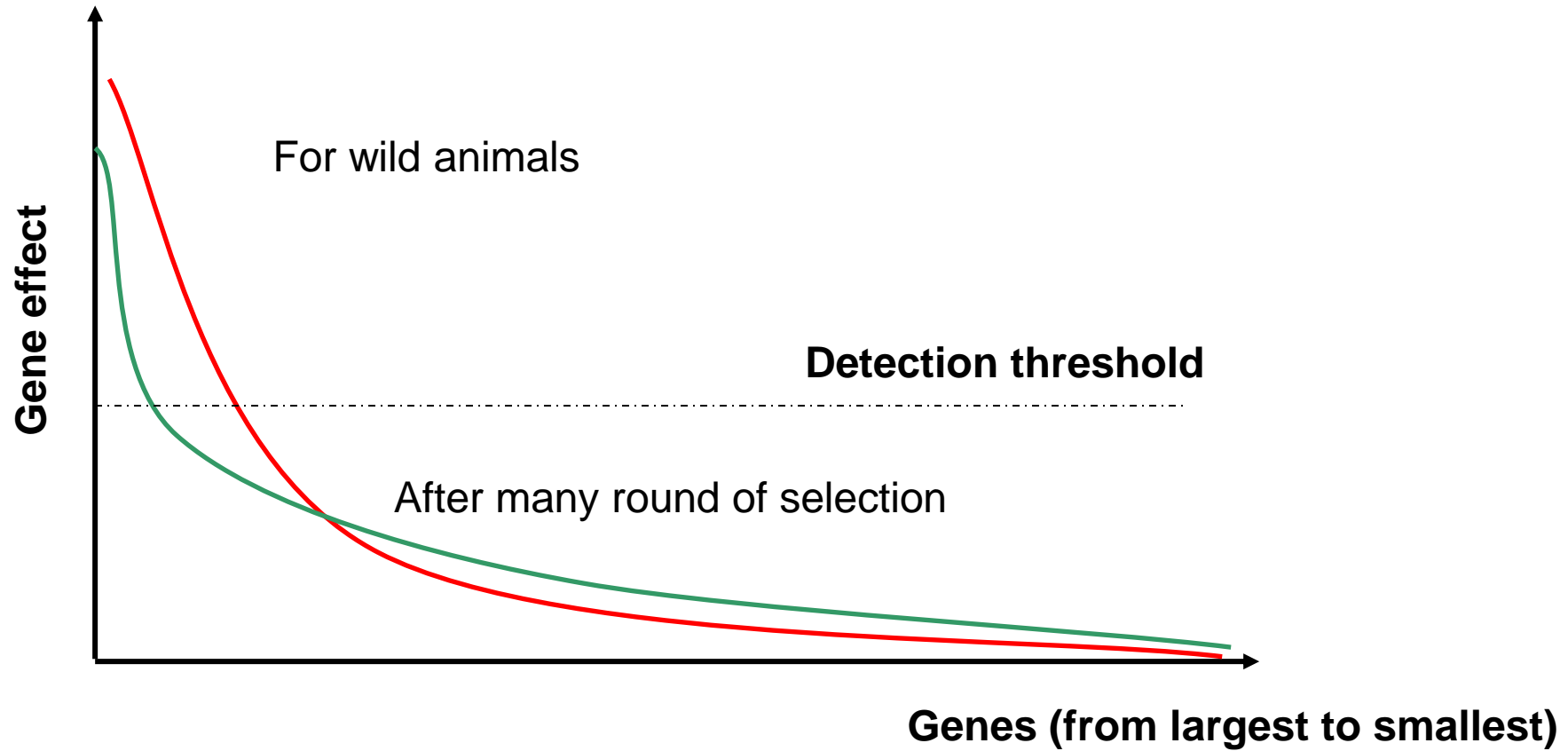


All SNP including QTN unweighted



QTN properly weighted

Distribution of QTL effects



Other factors

- Bayesian methods pick up large signals due to QTL, relationships and noise
 - If small signals undetected, probably large signals inflated
- Large signals without LD curves due to imputation errors (Li Ma, personal communication)

Conclusions

- QTN profile wide with small effective population size
 - Optimal window size 4 to 8 segments
 - About 1-2 Mb at $N_e=100$
- Sharp peak for QTN
- Large signals in GWAS due to QTN, relationships and noise and imputation
- Most QTNs probably below detection limit



Acknowledgements

