

(Artificial) intelligent and non-intelligent methods to explore multidimensional genomic data

Joanna Szyda & The THETA Biostatistic Group



Different worlds



Small datasets

Huge datasets



Different worlds



- Assymptotic test properties (type I, type II errors)
- Missing data patterns
- Do not remove observations



- Redundant, correlated information
- Fitting patterns to „noise”
- Input, interim, final data storage problem
- Statistical modeling problem
- Computationally intensive
- Extensive data preprocessing



Outline

multidimensional

genomic

AI

data



Outline

multidimensional

genomic

AI

data

1. cattle

2. humans



Outline

multidimensional

genomic

1. SNPs from WGS

AI

data

1. cattle
2. humans



Outline

multidimensional

genomic

1. SNPs from WGS

AI

1. DL based classification

data

1. cattle
2. humans



Outline

Dimensionality reduction by feature selection

1. basic
2. tagSNP
3. 1D-SRA
4. MD-SRA

genomic

1. SNPs from WGS

AI

1. DL based classification

data

1. cattle
2. humans



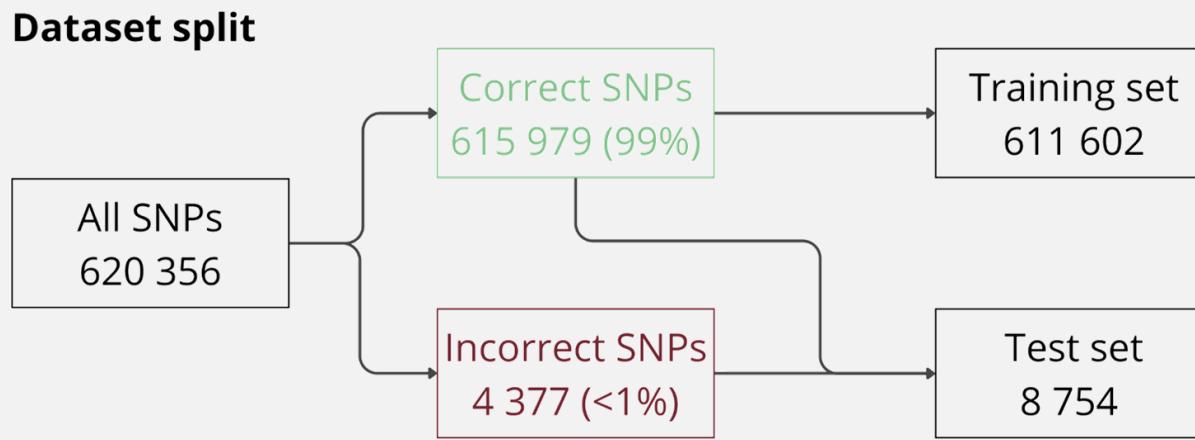
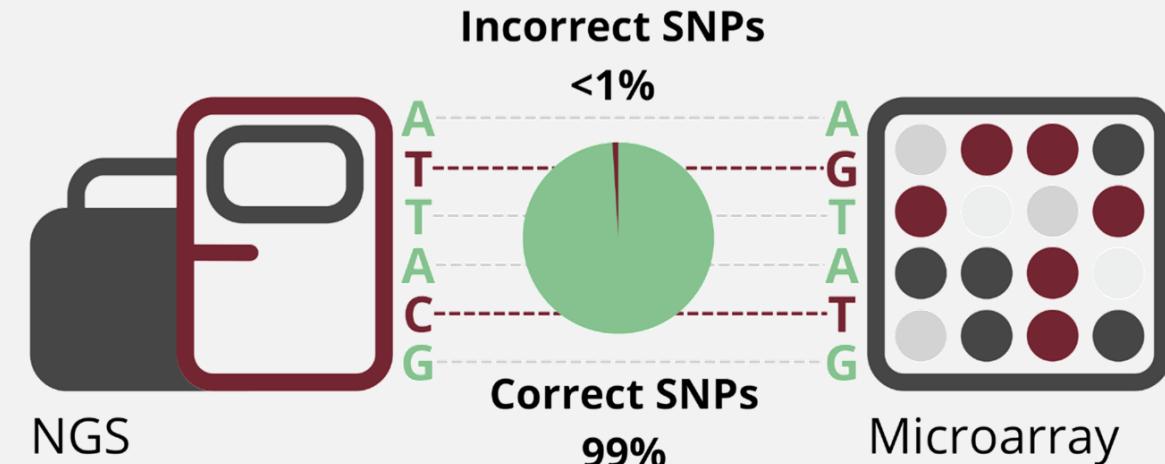
basic approach to feature selection

PCA

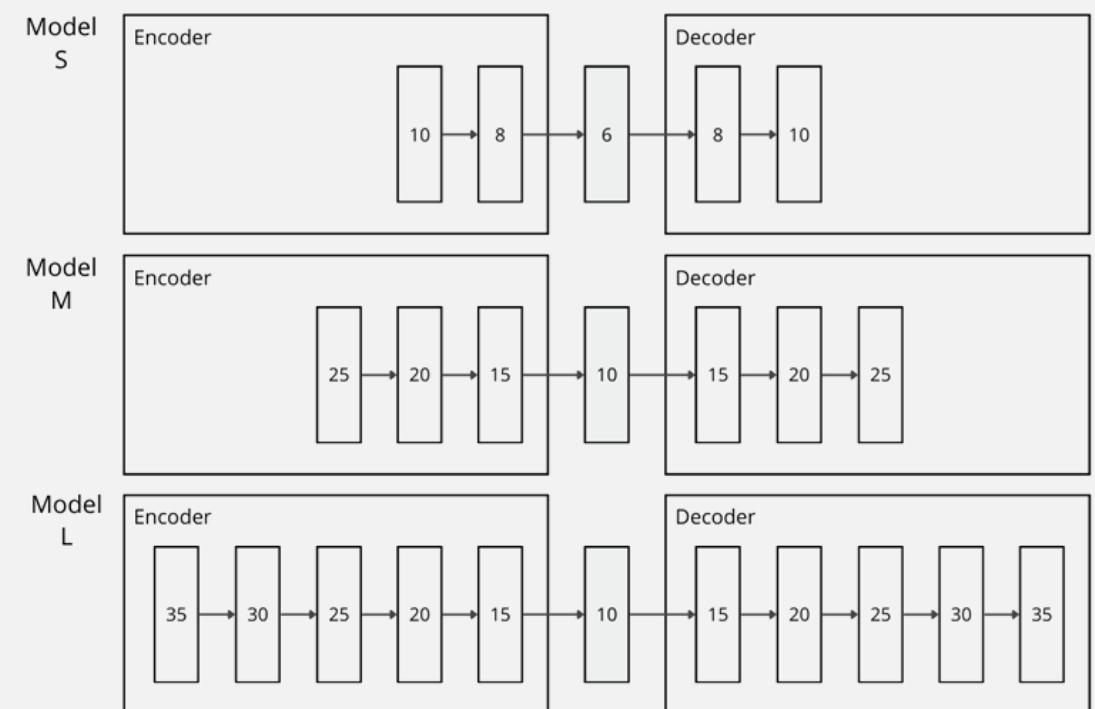


PCA

- SNP genotypes of 20 HF cows

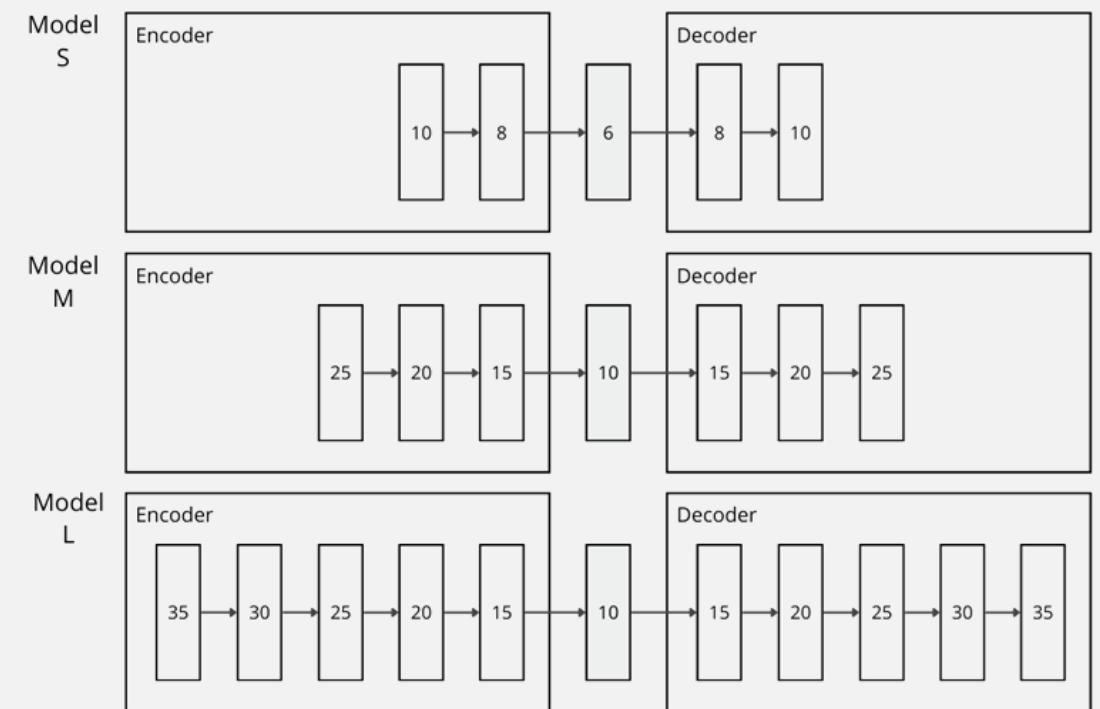


- 3 AutoEncoder architectures
- Classification of SNPs
 - correct
 - incorrect



PCA

- Considered features
 - Reference allele
 - Alternative allele
 - Sequencing depth
 - SNP genotype quality
 - Sequence context
 - 4 nucleotides downstream
 - 4 nucleotides upstream
- 3 AutoEncoder architectures
- Classification of SNPs
 - correct
 - incorrect



PCA

- Considered features
 - Reference allele
 - Alternative allele
 - Sequencing depth
 - SNP genotype quality
 - Sequence context
 - 4 nucleotides downstream
 - 4 nucleotides upstream
- Features measured on different scales
 1. Original encoding → poor classification
 2. PCA

23	57.69% ±0.46%	55.06% ±0.33%	58.64% ±0.32%
21	58.09% ±0.34%	56.46% ±0.61%	57.38% ±0.47%
19	59.40% ±0.76%	58.93% ±0.47%	57.89% ±0.28%
17	58.34% ±0.80%	59.05% ±0.72%	58.23% ±0.92%
15	55.55% ±0.71%	57.18% ±0.92%	57.37% ±0.94%
13	54.60% ±0.84%	51.82% ±0.42%	51.54% ±0.45%
11	54.30% ±0.74%	53.15% ±0.42%	54.29% ±0.95%
9	53.71% ±0.34%	51.77% ±0.73%	51.87% ±0.51%
7	52.30% ±0.71%	53.84% ±0.67%	53.99% ±0.83%
5	50.44% ±0.93%	52.56% ±1.05%	53.66% ±0.70%
3	51.51% ±0.62%	51.70% ±0.40%	51.94% ±0.53%

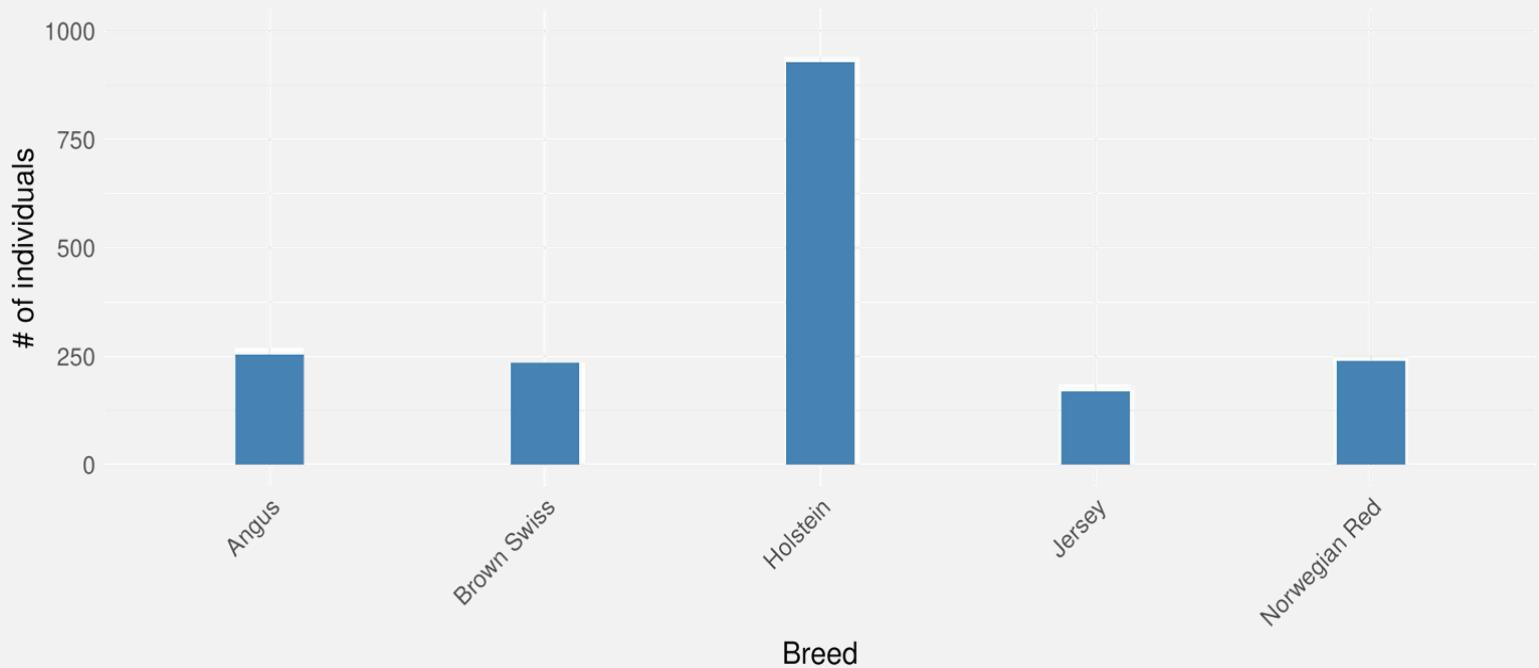


**feature selection based on
SNP tagging, 1D-SRA, MD-SRA**



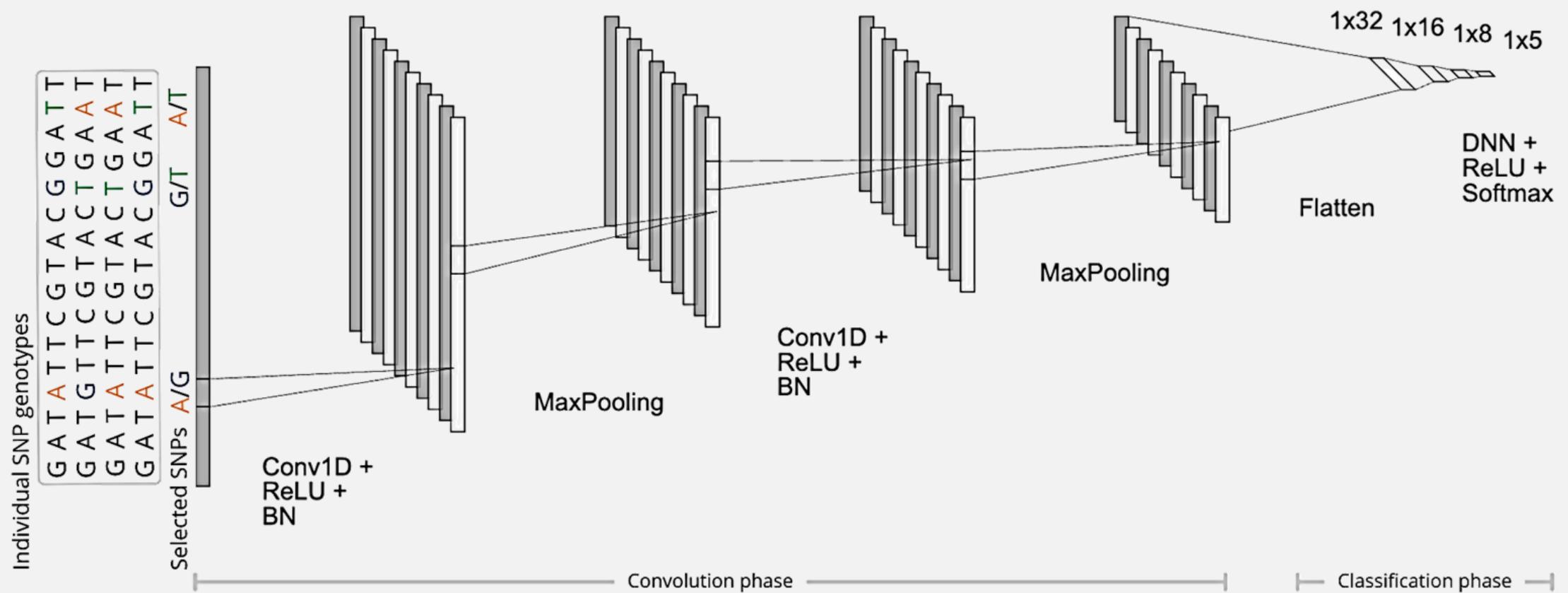
material

- 1000 Bull Genomes Project, Run 9
- Full data → 5,063 bulls with WGS → 33,595,340 SNPs
- Edited data → 1,825 bulls, 5 breeds → 11,915,233 SNPs = features



goal

- Classification of bulls to breeds → multi-class
- DL based on CNNs and DNNs



SNP tagging

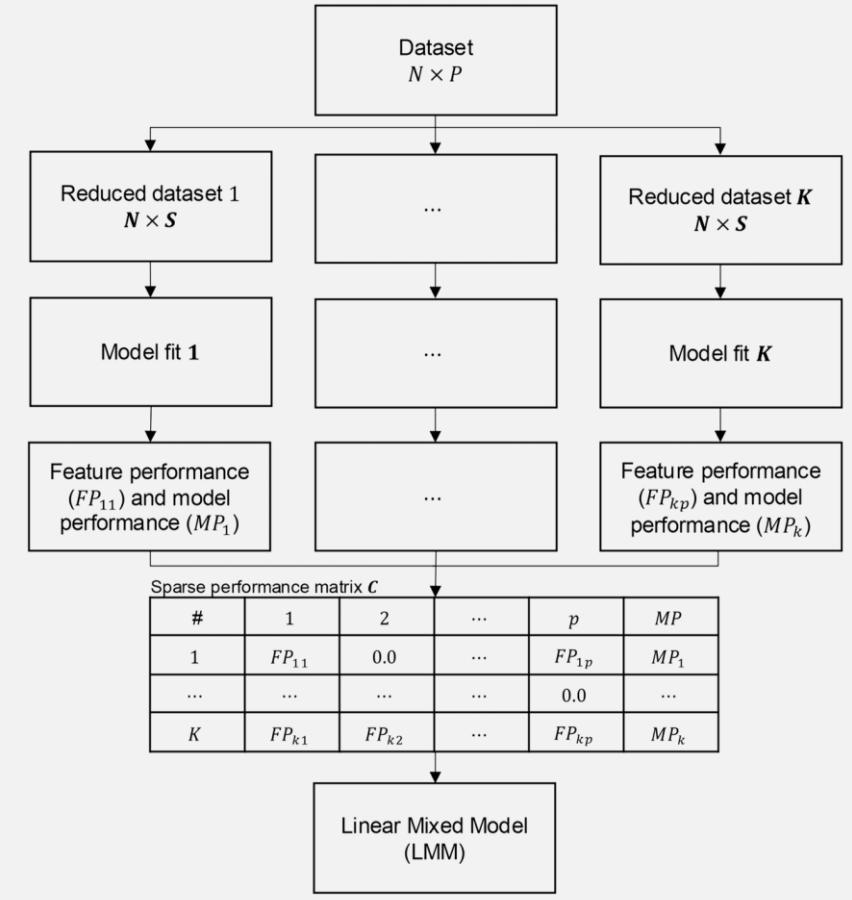
- Selection of representative SNPs based on pairwise LD
- PLINK 1.9
 - LD (R^2) threshold → 0.5
 - windowsize → 100 000 bp
 - stepsize → 1 000 bp



1D - SRA

Supervised Rank Aggregation

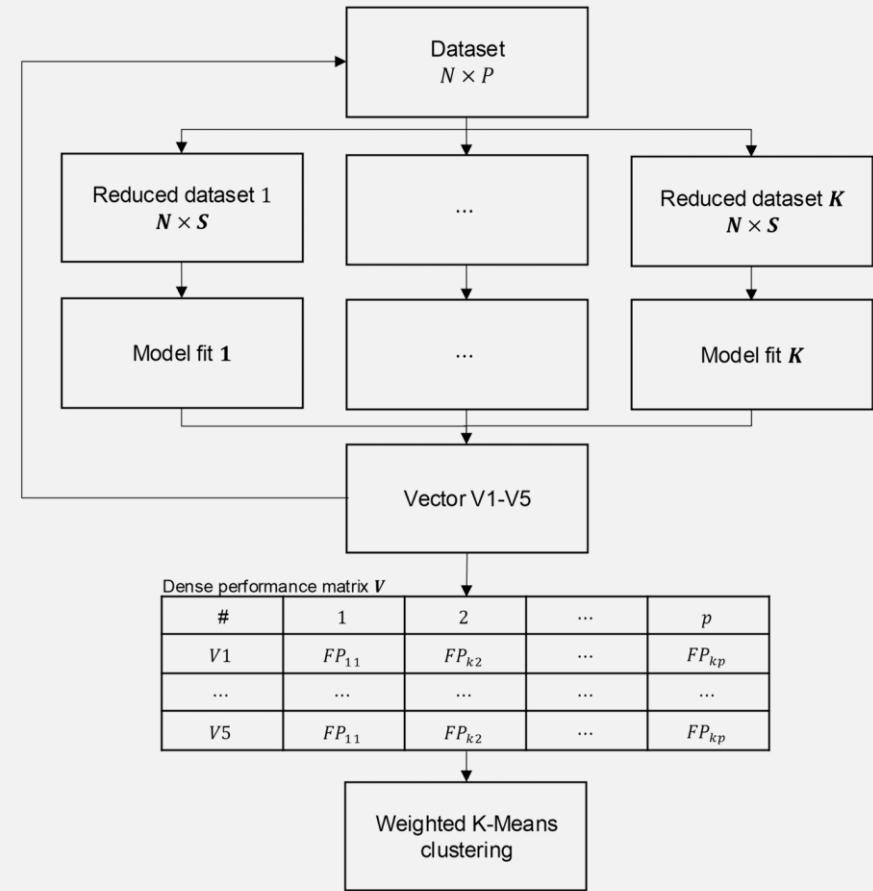
1. Fitting multiple **submodels** to **subsets** of features → multinomial logistic regression
2. Last model → aggregation of features' importance → linear mixed model
 - estimates
 - submodel fit
$$y = 1 + Z\mathbf{u} + \mathbf{e}$$
3. Selecting important features
 - **1D-K-means clustering**
4. DL based classification (CNN + DNN)



MD - SRA

Supervised Rank Aggregation

1. Fitting multiple submodels to subsets of features → multinomial logistic regression
2. ~~Last model → aggregation of features' importance → linear mixed model estimates~~
~~— submodel fit~~
3. Selecting important features
 - **MD-K-means clustering weighted by submodel's fit**
4. DL based classification (CNN + DNN)

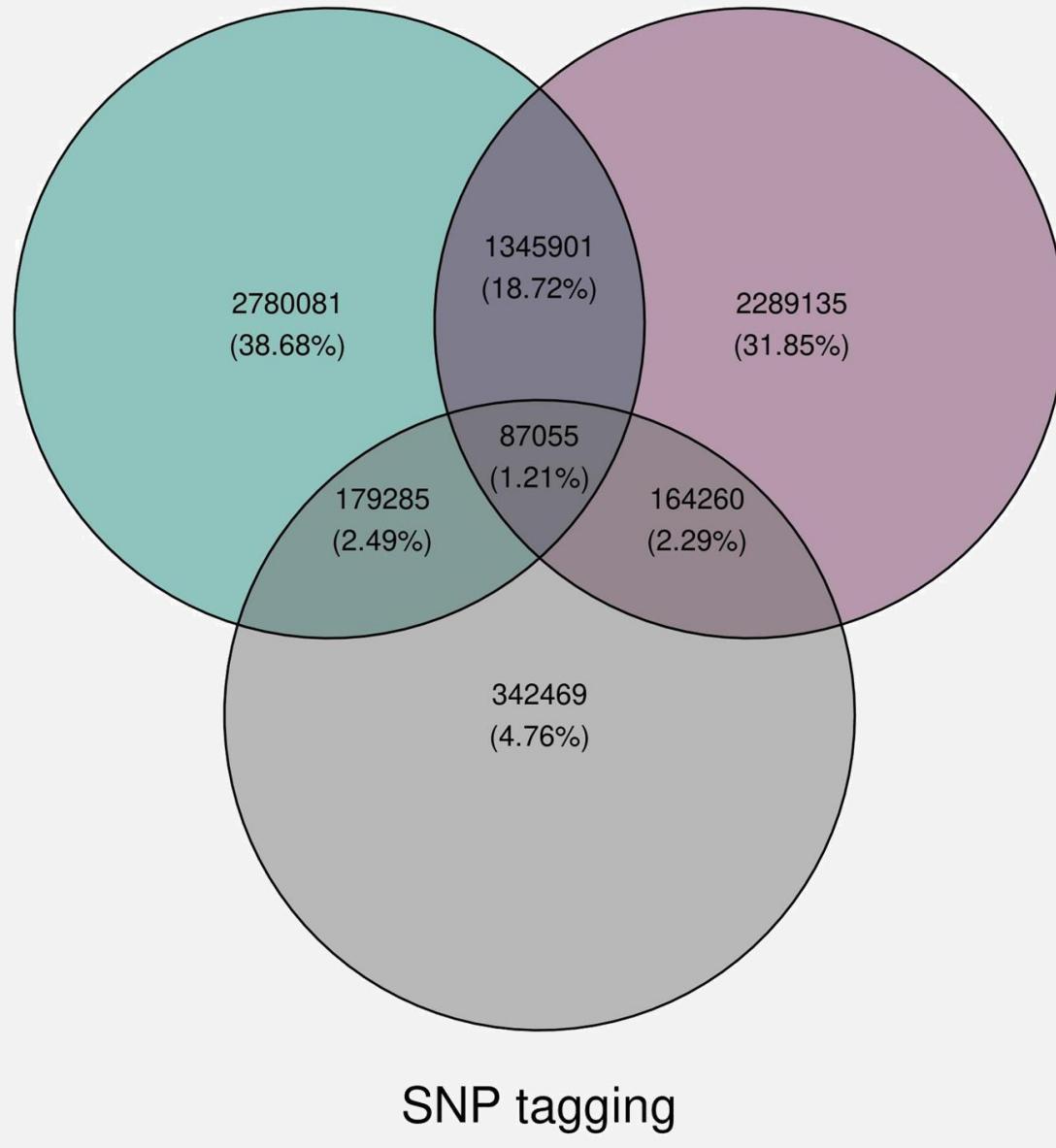


Comparison

Low overlap

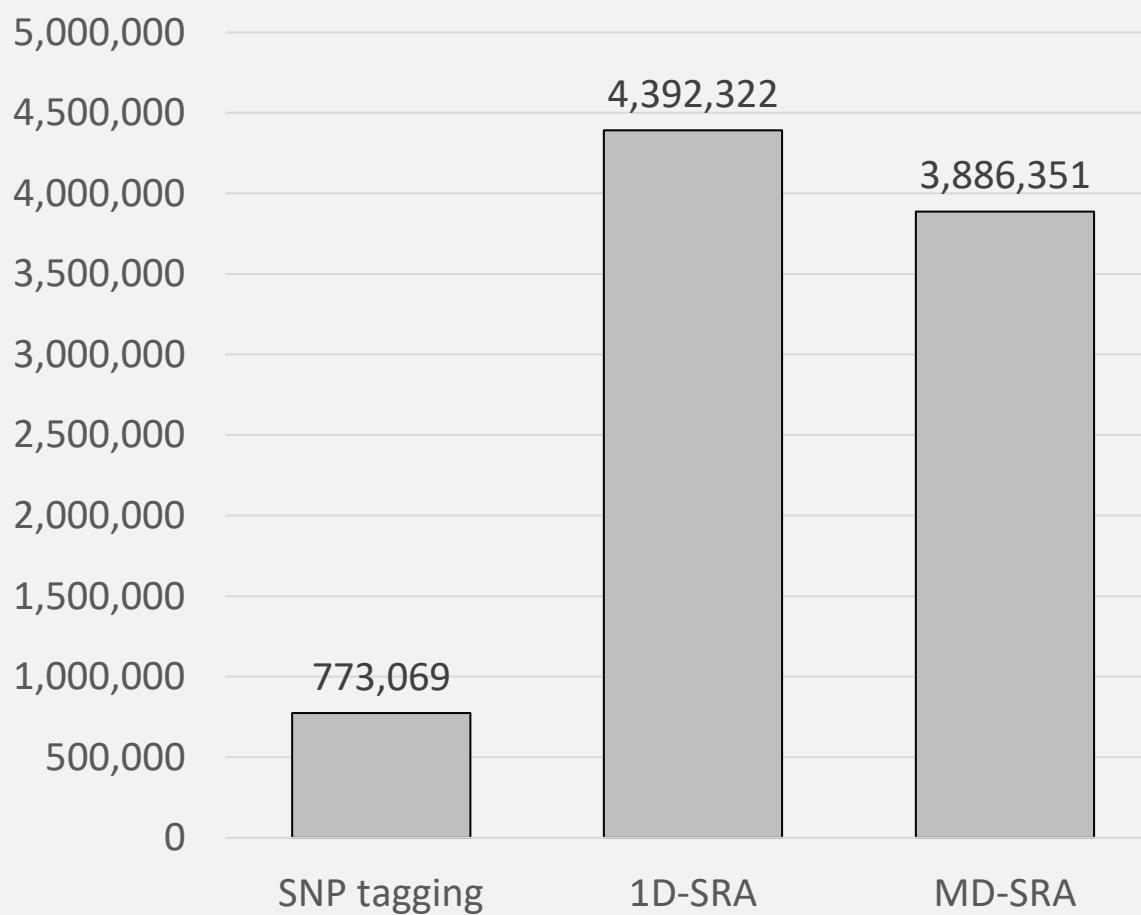
1D-SRA

MD-SRA

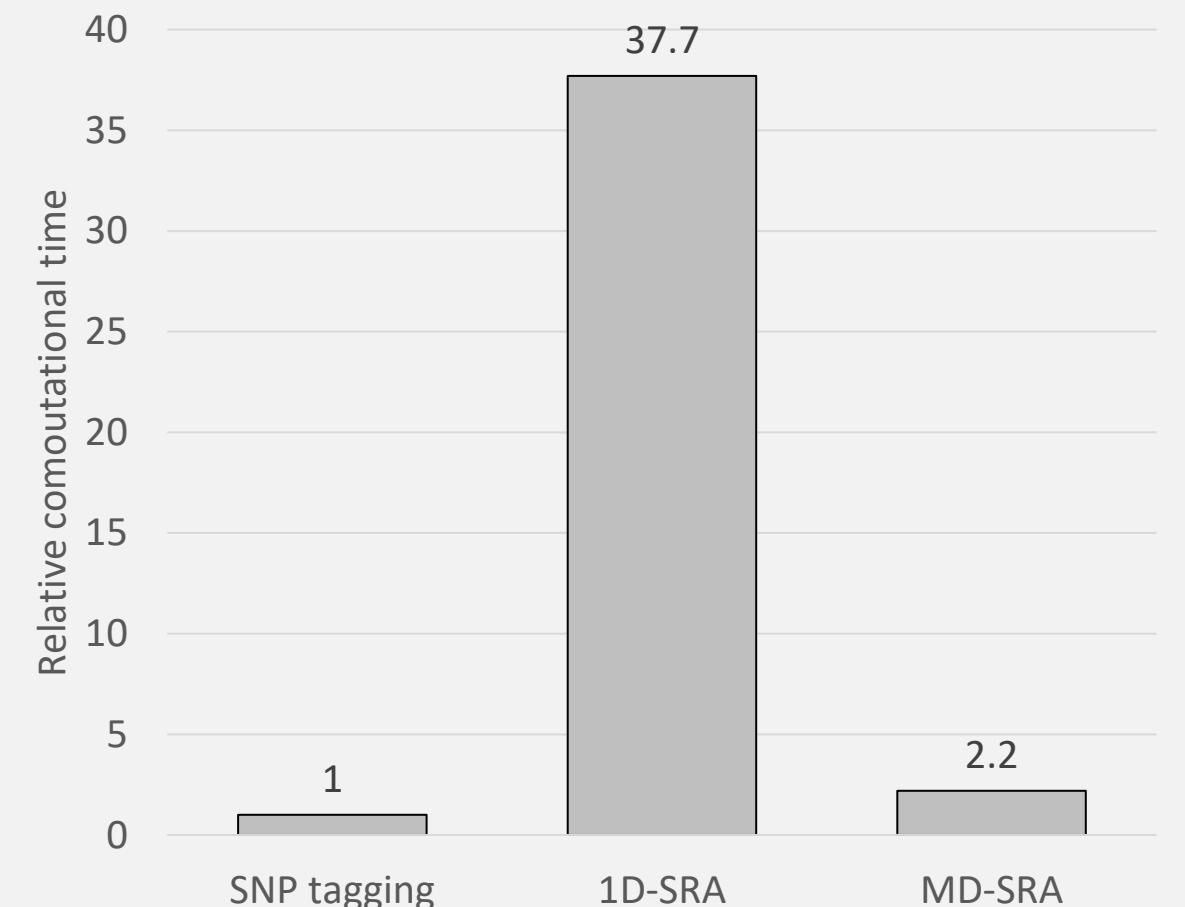


Comparison

Very different numbers of SNPs selected



Very different computational times



Comparison

SNP tagging

		True breed			
		Norwegian Red	Jersey	Brown Swiss	Holstein
Predicted breed	Norwegian Red	30	0	0	0
	Jersey	1	33	11	1
	Brown Swiss	0	0	171	0
	Holstein	2	1	0	46
	Angus	15	0	3	0

1D-SRA

		True breed			
		Norwegian Red	Jersey	Brown Swiss	Holstein
Predicted breed	Norwegian Red	47	0	0	1
	Jersey	1	34	3	0
	Brown Swiss	0	0	182	0
	Holstein	0	0	0	45
	Angus	0	0	0	49

MD-SRA

		True breed			
		Norwegian Red	Jersey	Brown Swiss	Holstein
Predicted breed	Norwegian Red	43	0	0	3
	Jersey	2	34	3	1
	Brown Swiss	0	0	182	0
	Holstein	0	0	0	46
	Angus	3	0	0	48

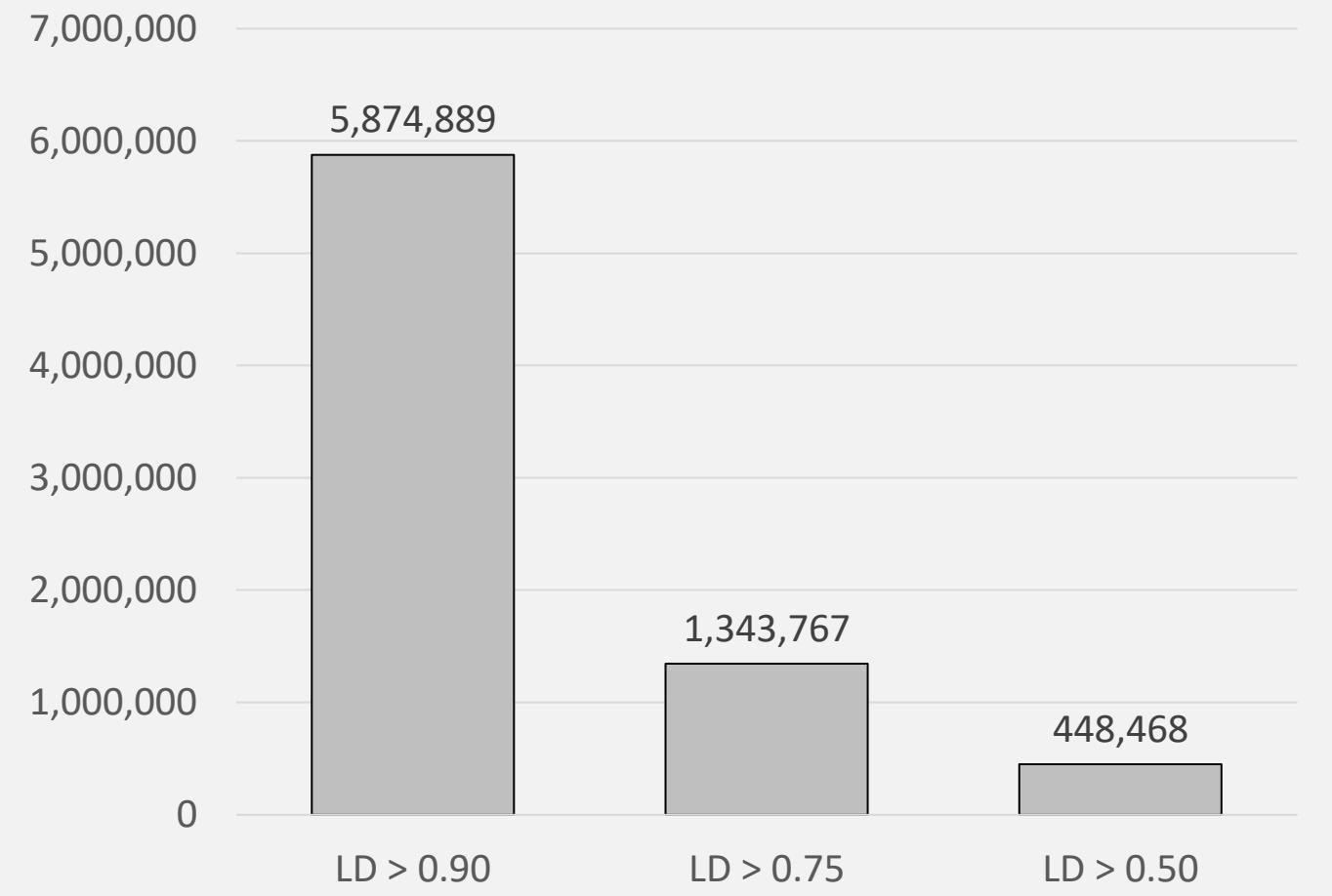


Other applications



material

- 1000 Bull Genomes Project, Run 9
- 528 HF bulls
 - EBVs for stature
 - 10,771,020 SNPs
 - SNP subsets selected by LD



GWAS

$$y = 1 + Zu + e$$

y EBVs

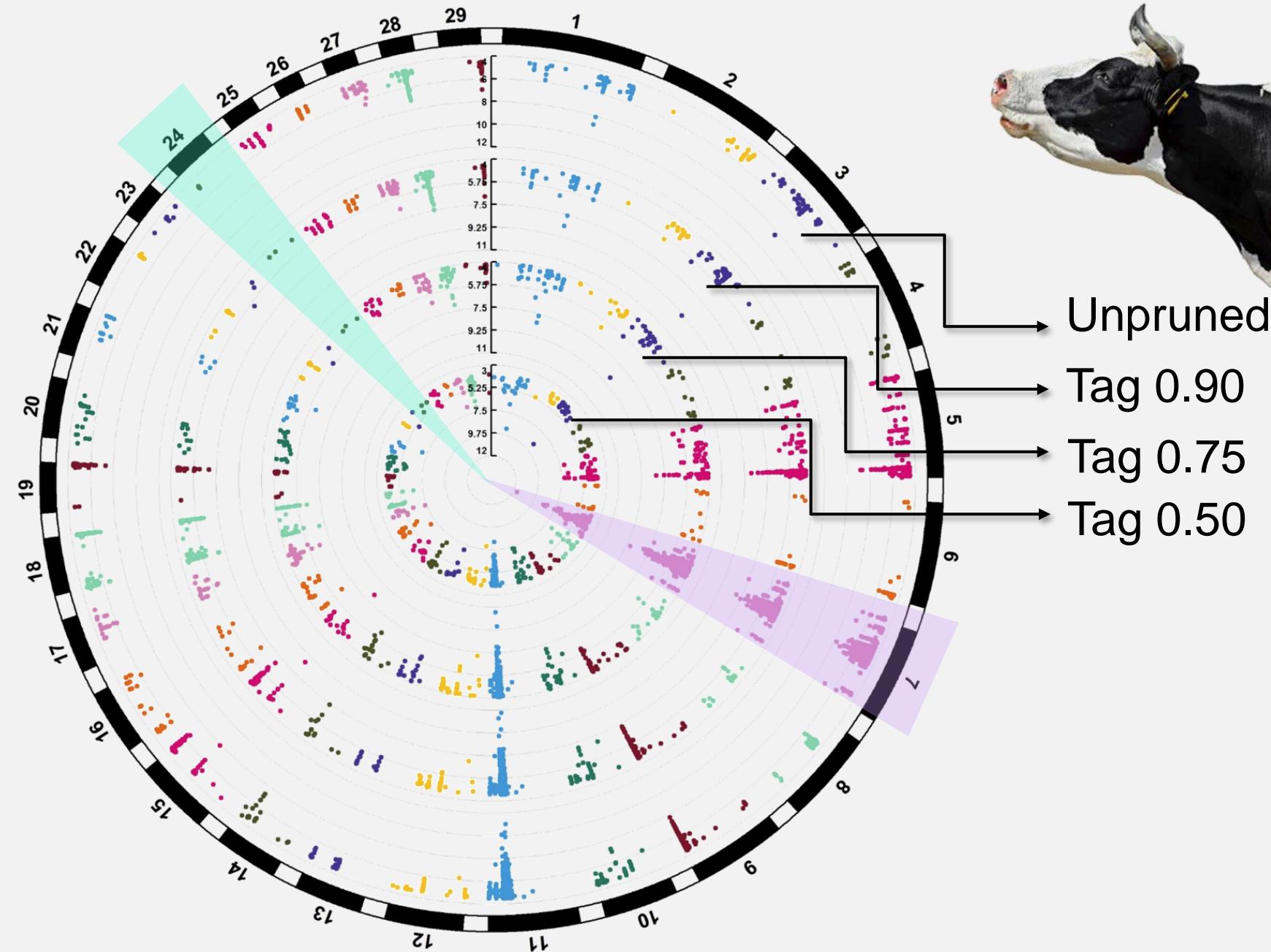
Z SNP genotypes

u SNP effect

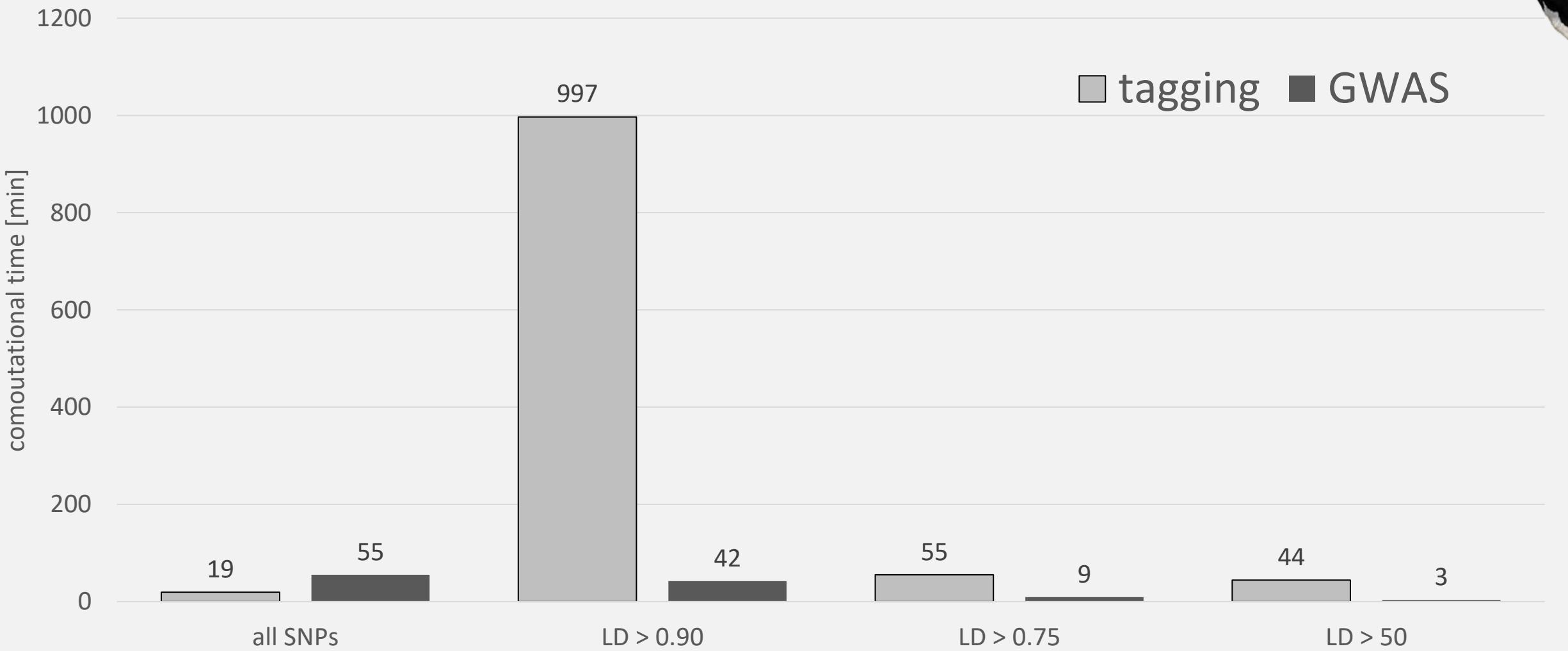
e residual

$$u \sim N(0, I\sigma_u^2)$$

$$e \sim N(0, I\sigma_e^2)$$



Differences in computational times





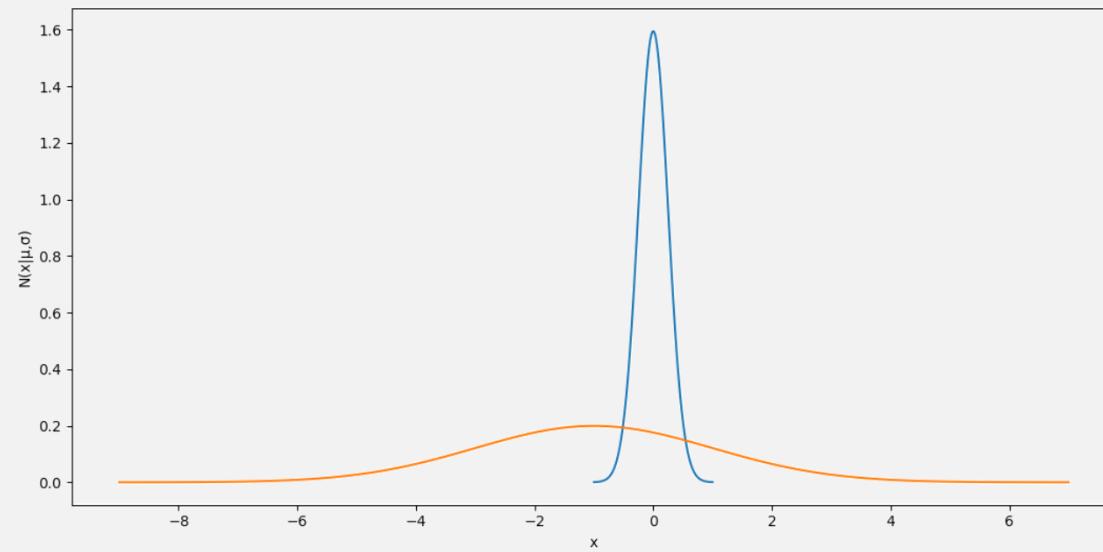
1D-SRA humans

- 1000 Polish Genomes
 - 1,222 persons with WGS → 41,836,187 SNPs → ~7,000,000 SNPs
 - Severe / non severe COVID-19 infection → binary classification
1. Submodels → logistic regression
 2. Important features → 1D-K-means / ~~Gaussian mixture~~
 3. Last model → linear mixed model $y = 1 + \mathbf{Z}\boldsymbol{u} + \boldsymbol{e}$
 4. DL based classification (DNN)

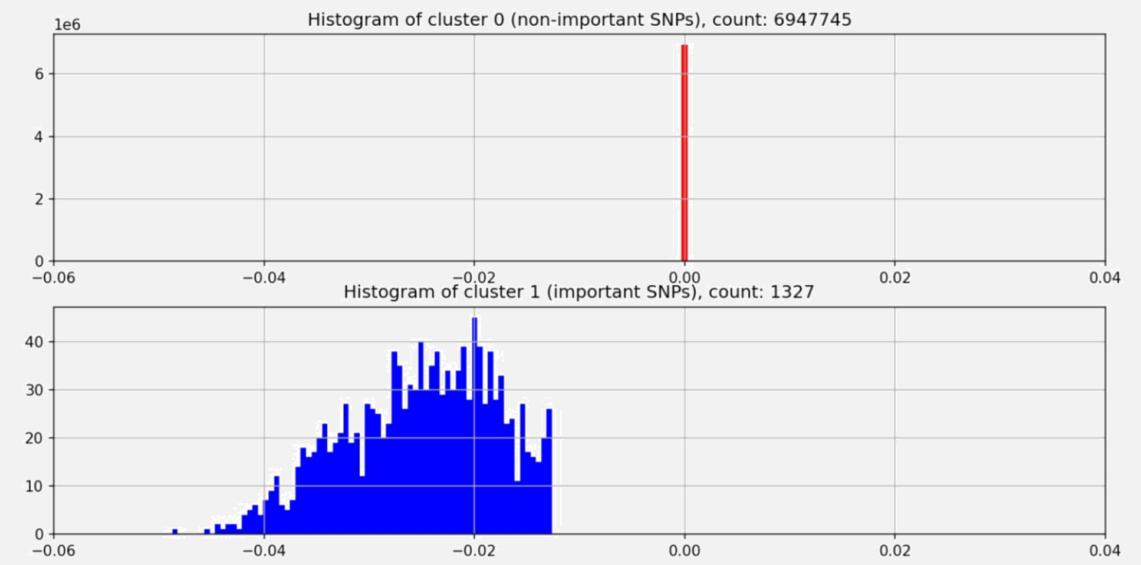


Approaches to 1-D clustering

Gaussian mixture



K-means



Conclusions

1. Feature selection necessary for ultra-high-dimensional data
 - Fitting pattern to noise → redundant information
 - Fitting models statistically impossible $p \gg n$
 - Fitting models computationally impossible → numerical or time constraints
2. Feature selection requires design
 - Selection process computationally intensive as well
 - Retain informative features
3. MD-SRA → promising algorithm



People



THETA

Statistical Genetics Group

Institute of Animal Genetics

► Leading Research Group **THETA**

THE BIOSTATISTIC GROUP

LEADER

PROFESSOR JOANNA SZYDA



Respect the data ...



Krzysztof Kotlarz

Magda Mielczarek

Dawid Słomian

Jakub Liu

