Introduction
○○○○○○○○○○○○○○

Genomic selection
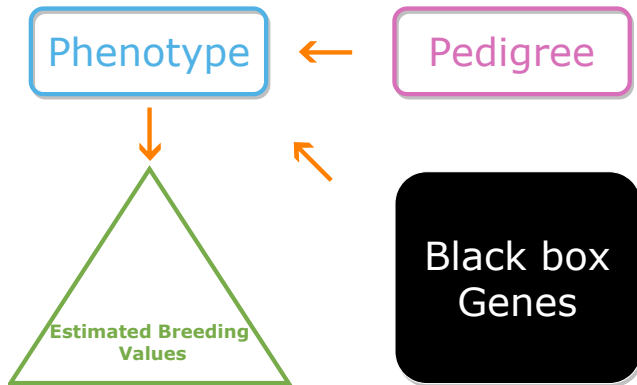○○○○○○○○○○○○

SNP-BLUP
○○○○○○○

GBLUP
○

# Genomic Selection in Animal Breeding

Hongding Gao

Natural Resources Institute Finland (LUKE)

Aug. 23$^{rd}$, 2024

@ University of Helsinki

# Traditional genetic prediction

▶ Successful but Limited

## Recap: Animal model

▶ Pedigree BLUP (PBLUP)

$$\mathbf{y} = \mathbf{Xb} + \mathbf{Za} + \mathbf{e}$$

$$\left[ \begin{array}{cc} \mathbf{X'X} & \mathbf{X'Z} \\ \mathbf{Z'X} & \mathbf{Z'Z} + \mathbf{A}^{-1}\frac{\sigma_e^2}{\sigma_a^2} \end{array} \right] \left[ \begin{array}{c} \hat{\boldsymbol{b}} \\ \hat{\boldsymbol{a}} \end{array} \right] = \left[ \begin{array}{c} \mathbf{X'y} \\ \mathbf{Z'y} \end{array} \right]$$

# Playing Lego

$$
\left[ \begin{array}{cc} \mathbf{X'X} & \mathbf{X'Z} \\ \mathbf{Z'X} & \mathbf{Z'Z} + \mathbf{A}^{-1}\frac{\sigma_e^2}{\sigma_a^2} \end{array} \right] \left[ \begin{array}{c} \hat{\boldsymbol{b}} \\ \hat{\boldsymbol{a}} \end{array} \right] = \left[ \begin{array}{c} \mathbf{X'y} \\ \mathbf{Z'y} \end{array} \right]
$$

## Limitations

▶ Slow genetic progress for low heritability traits (reproduction, health traits)

▶ Young animals (no own records, no offspring)

▶ Long generation intervals (e.g. dairy cattle 5 yr)
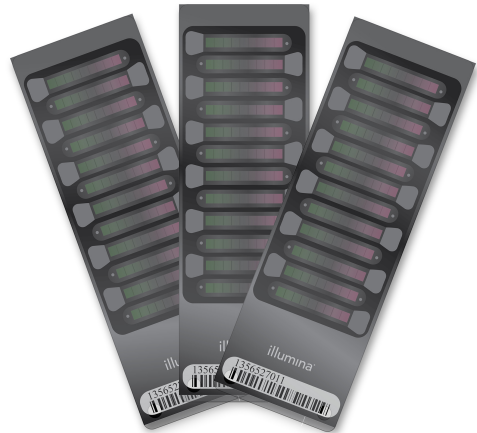
▶ High cost for progeny testing

## QTL

▶ QTL: Quantitative Trait Locus

▶ QTL: A locus/region of DNA which is correlated with variation of a quantitative trait in the phenotypes

▶ QTL mapping: identify the position of genes or markers that influence the trait

## Genetic markers

► A fragment of DNA that is associated with a certain location within the genome

► Many types: Microsatellite, RFLP, AFLP, Single nucleotide polymorphism (SNP), ...

Introduction
○○○○○○○●○○○○○

Genomic selection
○○○○○○○○○○○

SNP-BLUP
○○○○○○○○

GBLUP
○
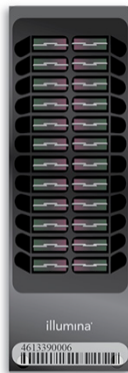
# Revolution of genotyping technology

▶ Highly dense molecular markers
covering the whole genome

  ▶ Single nucleotide polymorphisms
  (SNPs)

  ▶ High throughput genotyping
  technology (SNP chips in 2007)

  ▶ Available for most livestock species

  ▶ Genotyping costs continue to
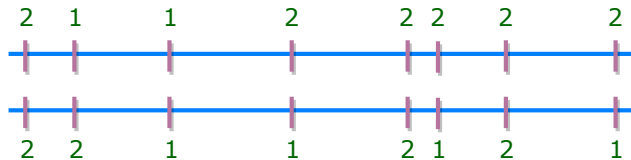  decrease

# Various SNP chips



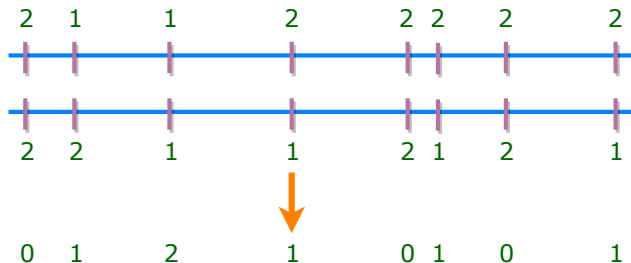BovineLD v2.0 (3k)      BovineSNP50 v3 (50k)      BovineHD (777k)

# SNP information

▶ SNPs have 2 alleles (diallelic)

# Recode SNP data

▶ Minor allele frequency (MAF): the frequency of less frequent allele at each locus in the population

▶ SNPs are commonly coded based on the counts of minor allele

▶ 0, 2 homozygote; 1 heterozygote

▶ Assume 1 is the minor allele for the example below

# One more example

▶ Assume lower case letter is the minor allele in that locus

|      | Locus 1 | Locus 2 | Locus 3 | Locus 4 |
|------|---------|---------|---------|---------|
| Ani1 | Ag      | AA      | GG      | Ct      |
| Ani2 | AA      | AA      | Ga      | Ct      |
| Ani3 | AA      | tt      | GG      | CC      |

## One more example

|      | Locus 1 | Locus 2 | Locus 3 | Locus 4 |
|------|---------|---------|---------|---------|
| Ani1 | 1       | 0       | 0       | 1       |
| Ani2 | 0       | 0       | 1       | 1       |
| Ani3 | 0       | 2       | 0       | 0       |

## SNP data

10001112200200121110111121111011110011211000201220022220111
12021012002111221100211120011110010110110102200110022011101
12002011010202221211221020100111000112022212221120211201201
20100202202000021100011202011221112111022011110000212202002
02210120200022112011101210011121110211211002010210002200022
22010000201100002202211022112101121110122220012112122200200
02002020201222110022222220022212111121002111120011011101120
02022200011120110102111212111020221002112012110011111102111
21102111220001011011102022002111010201112111101120210210212
12110110221220012110112110120220110022200210021100011100212
10211011100022200202212121100022201020022221212211211112002
01102020012222221122120212112101100121101102002200022001002
00011110110012110212121112010101212022210101011111021102112
21111112121112101101200111110211110111112201210121211101022
20202121122212022200212121012121020110011122212111011

Introduction
0000000000000

Genomic selection
●00000000000

SNP-BLUP
00000000

GBLUP
○

# Genomic selection

▶ First proposed by Meuwissen et al., 2001

▶ All SNPs fitted simultaneously

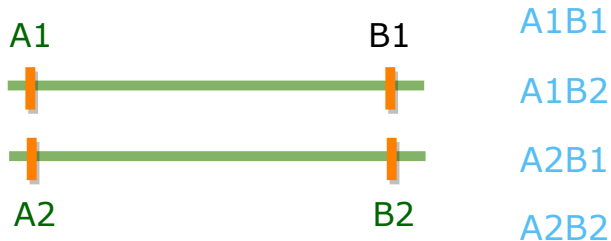▶ Innovating & Boosting the breeding cycle in animal and plant

Introduction
0000000000000

Genomic selection
0●0000000000

SNP-BLUP
00000000

GBLUP
0

# Linkage Disequilibrium (LD)

▶ Alleles present at the two loci are not independent

# Linkage Disequilibrium (LD)

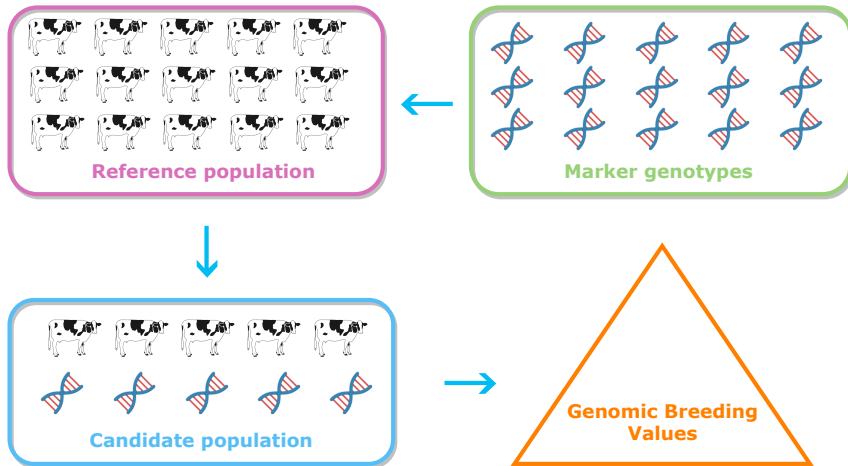- A1 is more often associated with B1 than A2 with B1

- $freq(A1B1) \neq freq(A1)freq(B1)$

Introduction
Genomic selection
SNP-BLUP
GBLUP

000000000000
000●00000000
00000000
0

## Linkage Disequilibrium (LD)

▶ Genomic selection requires LD

▶ SNPs are in LD with the QTLs across the whole population

▶ Assume all QTLs in the genome can be traced by SNPs

▶ SNP density must be sufficiently high to ensure that all QTLs are in LD with at least a marker

# Genomic selection



Reference population

Marker genotypes

Candidate population

Genomic Breeding Values

Introduction
00000000000000

Genomic selection
00000●000000

SNP-BLUP
00000000

GBLUP
0

## Genomic breeding values

▶ Sum of single marker effects

$$\mathbf{GEBV} = \sum_{i=1}^{m} \mathbf{M}_i \hat{\mathbf{g}}_i$$
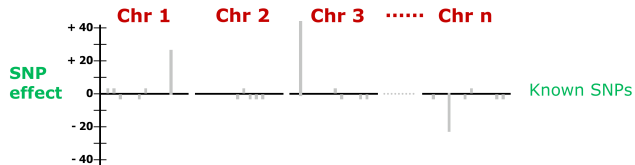
where
$m$ is the number of markers
$\mathbf{M}_i$ is the SNP code at locus $i$
$\hat{\mathbf{g}}_i$ is SNP effect at locus $i$

Introduction
○○○○○○○○○○○○○○○

Genomic selection
○○○○○○●○○○○○

SNP-BLUP
○○○○○○○○

GBLUP
○

# Genomic breeding values



**1 + 1 - 1 - 1 + 1 + 25 - 1 + 1 - 1 - 1 - 1 + 42 + 1 - 1 - 1 - 1 - 1 - 22 - 1 + 1 - 1 - 1 = +38**
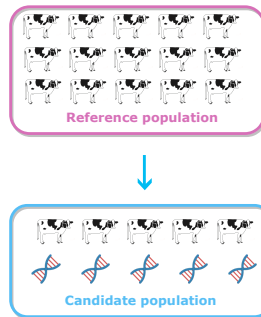
# Benefits of genomic selection

- ▶ Double the genetic progress (Schaeffer, 2006)
  - ▶ Increase accuracy of selection
  - ▶ Reduce generation interval

- ▶ Reduce costs in the breeding plan by 92% (Schaeffer, 2006)
  - ▶ Reduce the scale for progeny testing (dairy cattle)

Introduction
○○○○○○○○○○○○○○

Genomic selection
○○○○○○○○○●○○○

SNP-BLUP
○○○○○○○

GBLUP
○

# Benefits of genomic prediction

▶ Accuracy exceeds 0.8 and gains in reliability up to 48% in US Holsteins (Wiggans et al., 2011)

▶ Accuracy increases by 50% in pig breeding (Knol et al., 2016)

▶ Accuracy increases by 20%-50% in poultry breeding (Wolc et al., 2015; Wang et al., 2013)

▶ Promising in plant breeding, e.g. genetic gain 1.4 to 2.7 times higher than phenotypic selection in wheat (Battenfield et al., 2016)

Introduction
0000000000000

Genomic selection
0000000000●00

SNP–BLUP
00000000

GBLUP
0

# Factors affecting accuracy of GP

▶ Prediction models

▶ Size of reference population

▶ Heritability of trait

▶ Marker density

▶ Population structure

▶ Variance components

▶ . . .

Introduction
000000000000000

Genomic selection
0000000000●0

SNP-BLUP
00000000

GBLUP
o

## Two equivalent models

▶ SNP-BLUP: marker based model
  ▶ Estimated marker effects
  ▶ Better when no. of animals $>$ no. of markers

▶ GBLUP: breeding value based model
  ▶ Estimated breeding values directly
  ▶ Better when no. of animals $<$ no. of markers

Introduction
000000000000

Genomic selection
00000000000●

SNP-BLUP
00000000

GBLUP
○

# Response variables

- ▶ Deregressed proofs (DRP)
  - ▶ Derived from EBV which is a regressed variable

- ▶ Daughter yield deviations (DYD)
  - ▶ Average of the daughter's actual performance adjusted for fixed and non-genetic random effects and genetic effects of the daughters' dams

- ▶ Estimated breeding values (EBV)

Introduction
00000000000000

Genomic selection
00000000000

SNP-BLUP
●0000000

GBLUP
○

## SNP-BLUP

▶ We assume markers can explain all the genetic variance

▶ We assume marker effects are identically and independently distributed

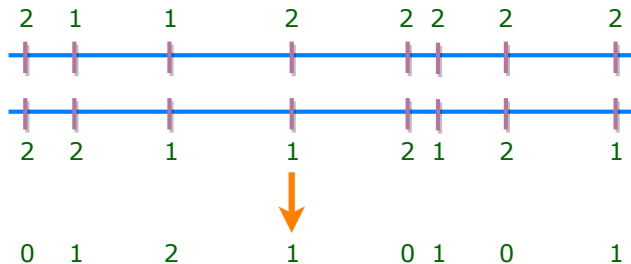$$\mathbf{y} = \mathbf{Xb} + \sum_{i}^{m} \mathbf{M}_i \mathbf{g}_i + \mathbf{e}$$

where
**M** is a n (no. of animals) by m (no. of markers) matrix
**g** is a vector of marker effects
**e** is a vector of random residuals

Introduction
○○○○○○○○○○○○○○

Genomic selection
○○○○○○○○○○○○

SNP-BLUP
○●○○○○○○

GBLUP
○

# SNP coding

▶ **M** is a marker covariates matrix containing SNP codes
▶ SNPs are commonly coded based on counts of minor allele
▶ 0, 2 homozygote; 1 heterozygote

Introduction
0000000000000

Genomic selection
000000000000

SNP-BLUP
00●00000

GBLUP
0

# Centering

▶ To set the mean value of allele effects equal to 0

$$\mathbf{Z} = \mathbf{M} - \mathbf{P}$$

where
$\mathbf{P}$ is a matrix with column $j$ equal to $2p_j$
$p_j$ is the allele frequency of the second allele at locus $j$

## SNP-BLUP

$$\mathbf{y} = \mathbf{Xb} + \mathbf{Zg} + \mathbf{e}$$

$$\mathbf{g} \sim N(\mathbf{0}, \mathbf{I}\sigma_g^2) \quad \mathbf{e} \sim N(\mathbf{0}, \mathbf{I}\sigma_e^2)$$

where
$\sigma_g^2$ is the SNP variance for each SNP
$\sigma_e^2$ is the residual variance

Introduction
○○○○○○○○○○○○○

Genomic selection
○○○○○○○○○○○

SNP-BLUP
○○○○●○○○

GBLUP
○

## MME for SNP-BLUP

$$\left[\begin{array}{cc} \mathbf{X'R^{-1}X} & \mathbf{X'R^{-1}Z} \\ \mathbf{Z'R^{-1}X} & \mathbf{Z'R^{-1}Z} + (\mathbf{I}\sigma_g^2)^{-1} \end{array}\right] \left[\begin{array}{c} \hat{\mathbf{b}} \\ \hat{\mathbf{g}} \end{array}\right] = \left[\begin{array}{c} \mathbf{X'R^{-1}y} \\ \mathbf{Z'R^{-1}y} \end{array}\right]$$

$$\left[\begin{array}{cc} \mathbf{X'X} & \mathbf{X'Z} \\ \mathbf{Z'X} & \mathbf{Z'Z} + \mathbf{I}\alpha \end{array}\right] \left[\begin{array}{c} \hat{\mathbf{b}} \\ \hat{\mathbf{g}} \end{array}\right] = \left[\begin{array}{c} \mathbf{X'y} \\ \mathbf{Z'y} \end{array}\right]$$

where

$$\alpha = \sigma_e^2/\sigma_g^2$$

$$\sigma_g^2 = \frac{\sigma_a^2}{\sum_{i=1}^{m} 2p_i q_i}$$

Introduction
0000000000000

Genomic selection
00000000000

SNP-BLUP
00000●00

GBLUP
0

# Direct genomic value (DGV)

▶ An "old" terminology for GEBV
▶ Means GEBV calculated only by marker information

$$\textbf{DGV} = \textbf{a} = \textbf{Z}\hat{\textbf{g}}$$

Introduction
0000000000000

Genomic selection
00000000000

SNP-BLUP
0000000●0

GBLUP
0

## SNP-BLUP with polygenic effect

▶ We assume markers cannot explain all the genetic variance

▶ We assume marker effects are identically and independently distributed

$$\mathbf{y} = \mathbf{Xb} + \sum_{i}^{m} \mathbf{M}_i \mathbf{g}_i + \mathbf{Wu} + \mathbf{e}$$

where
**M** is a n (no. of animals) by m (no. of markers) matrix
**g** is a vector of marker effects
**u** is a vector of polygenic effects
**W** is a design matrix linking records to animals

## SNP-BLUP with polygenic effect

$$\mathbf{y} = \mathbf{Xb} + \sum_i^m \mathbf{M}_i \mathbf{g}_i + \mathbf{Wu} + \mathbf{e}$$

where

$$\mathbf{g} \sim N(\mathbf{0}, \mathbf{I}\sigma_g^2) \quad \mathbf{u} \sim N(\mathbf{0}, \mathbf{A}\sigma_u^2)$$

Let $\omega$ equal to the proportion of total genetic variance cannot be explained by markers, then

$$\sigma_u^2 = \omega \sigma_a^2 \quad \sigma_g^2 = \frac{(1-\omega)\sigma_a^2}{\sum_{i=1}^m 2p_i q_i}$$

# PBLUP ⇒ GBLUP

▶ You can understand **GBLUP** is a "improved" version of traditional **BLUP**

$$\begin{bmatrix} \mathbf{X'X} & \mathbf{X'Z} \\ \mathbf{Z'X} & \mathbf{Z'Z} + \mathbf{A}^{-1}\frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{b}} \\ \hat{\boldsymbol{a}} \end{bmatrix} = \begin{bmatrix} \mathbf{X'y} \\ \mathbf{Z'y} \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{X'X} & \mathbf{X'Z} \\ \mathbf{Z'X} & \mathbf{Z'Z} + \mathbf{G}^{-1}\frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{b}} \\ \hat{\boldsymbol{a}} \end{bmatrix} = \begin{bmatrix} \mathbf{X'y} \\ \mathbf{Z'y} \end{bmatrix}$$