

Genomic evaluation

Vincent Ducrocq



GABI, Jouy-en-Josas



Whole Genome Sequence

- After Human (2001) and mice, WGS of the chicken (2004), the dog (2005), bovine (2006), horse (2007), pig (2009), ...
- Entirely in the public domain
- Sequencing of different individuals => **polymorphisms**
 - 3.2 million bovine polymorphisms in dbSNP
 - >30 millions known today
- New technologies for genotyping and sequencing

SNP : Single Nucleotide Polymorphism

DNA Variation of one base

..GAATCTTATGCTATACTACATAATTATATACTAAT**C**GGGTATTGTTCTTAT..

..GAATCTTATGCTATACTACATAATTATATACTAAT**A**GGGTATTGTTCTTAT..

↑
SNP

Genotyping chips

- Simultaneous genotyping of many SNP
- From few dozens up to several million SNP
- Two main technology providers, Illumina and Affymetrix
- Illumina products in cattle
 - 3000 (□7000□ 1000□ 20000=« LD »)
 - **54 000=« 50k »**
 - 777 000=« HD »

FIGURE 1: BOVINESNP50 BEADCHIP

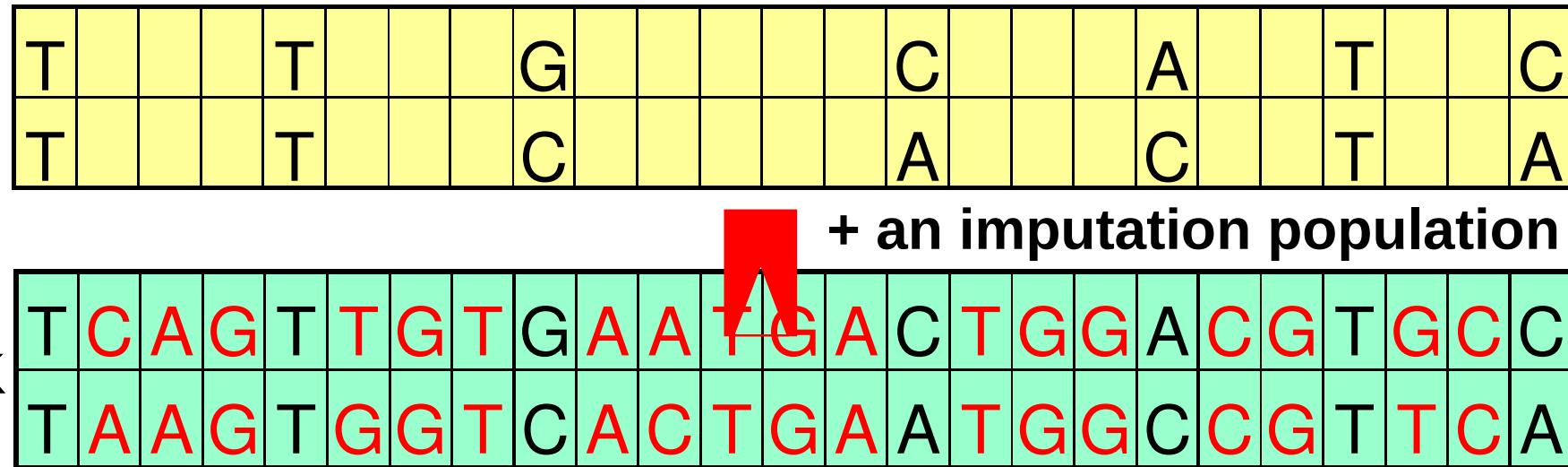


The BovineSNP50 BeadChip features more than 54,000 evenly-spaced SNPs across the entire bovine genome.

Chips of different densities:

IMPUTATION

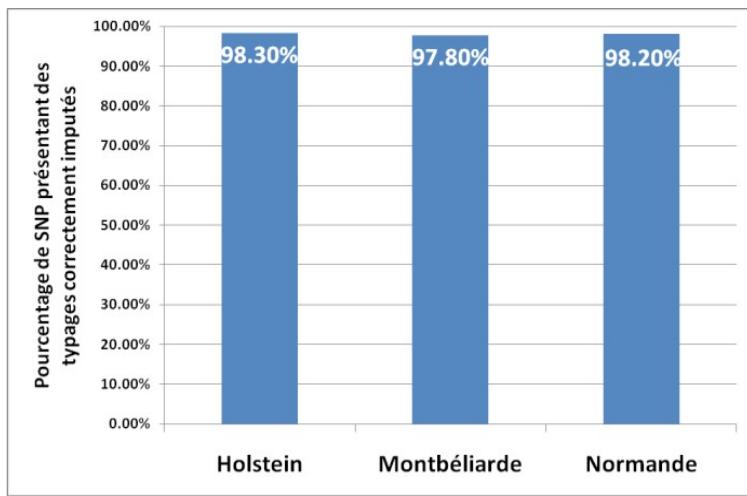
- i_p_ta_i_o_c_nsi_t_i_pr_di_t_n_t_e_m_s_g_l_t_e_r_i_h_w_d_o_a_s_t_c_ + a dictionary
 - imputation consists in predicting the missing letters within a word or a sentence



Works well ! (~1% error in Holstein if sire genotyped with 54k)

Imputation

Beagle/DagPhase



Fimpute



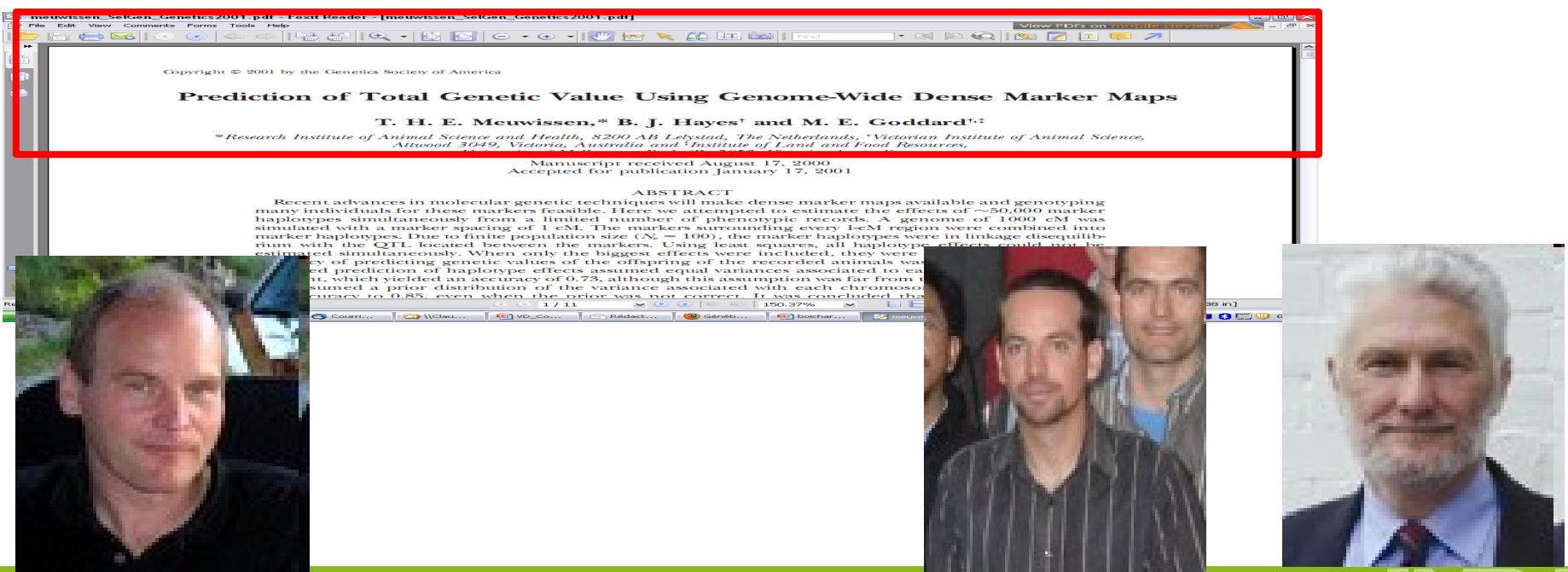
(Dassonneville et al.,
2011)

(Saintilan et al., 2014)

Computation time divided by 3 – 200-300 SNP/animal corrected

Genomic selection

selection based on the estimation of the genetic value of candidates using information on dense markers covering the whole genome
= selection based on results from a genomic evaluation



What is genomic selection?

Pheno= 0



+8



+15



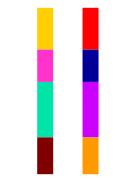
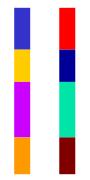
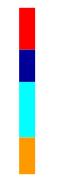
+6



-6



+12



Genotyped and
phenotyped animals



Reference
population

What is genomic selection?

Pheno= 0



+8



+15



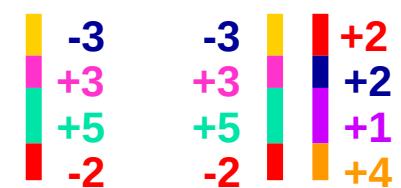
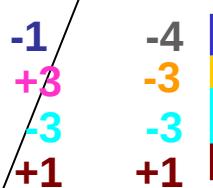
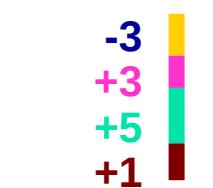
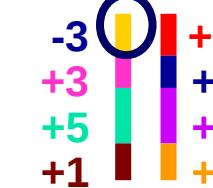
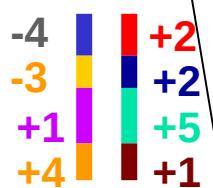
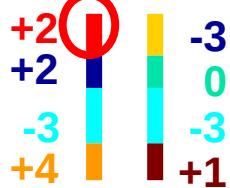
+6



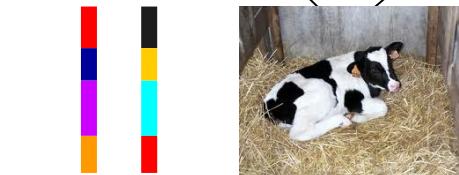
-6



+12

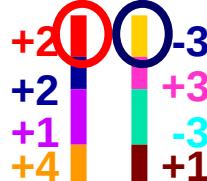


2/ Genotype
at birth



3/ Apply!

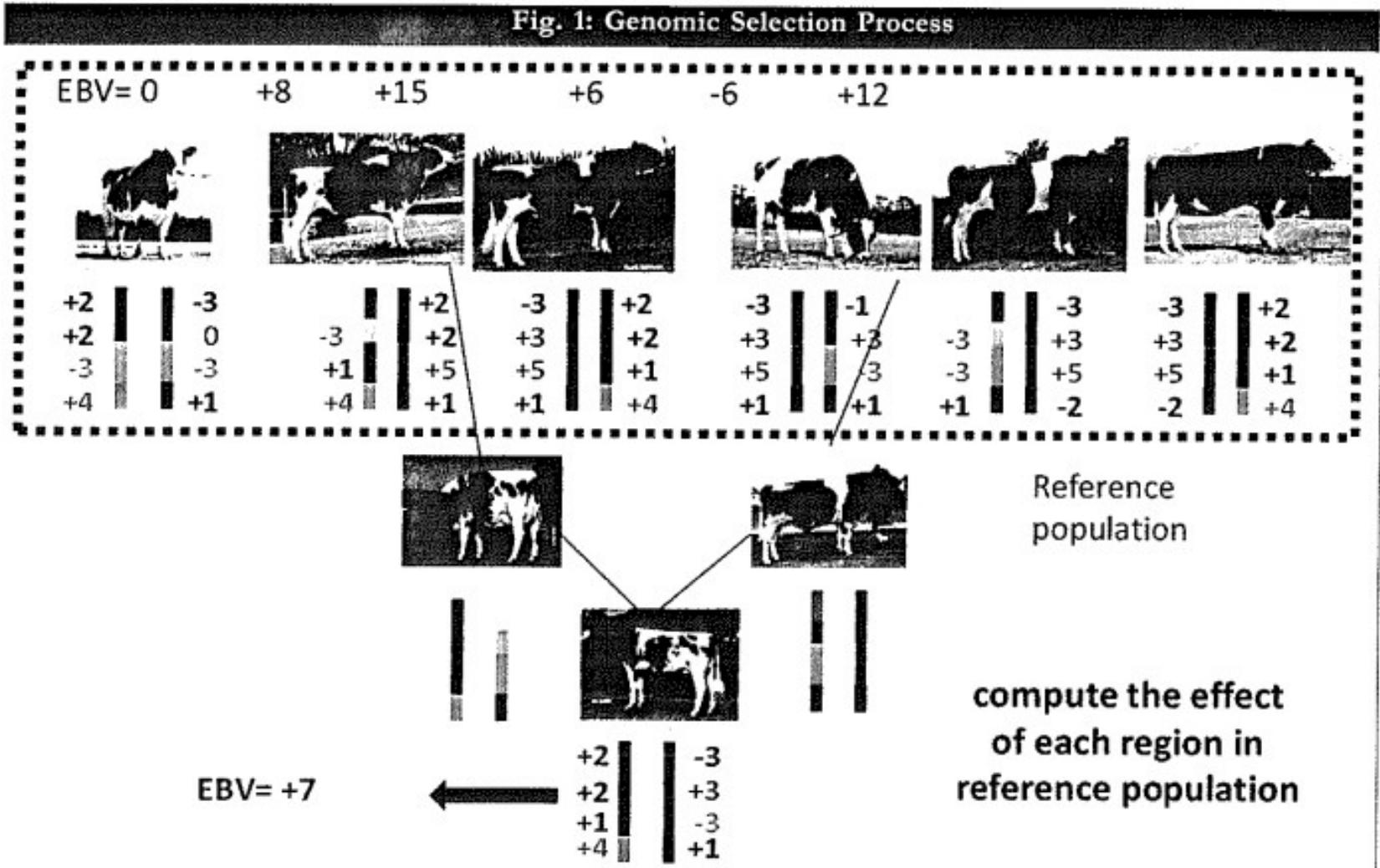
(G)EBV= +7



Reference
population

1/ Compute the effect
of each chromosome
region in the
reference population

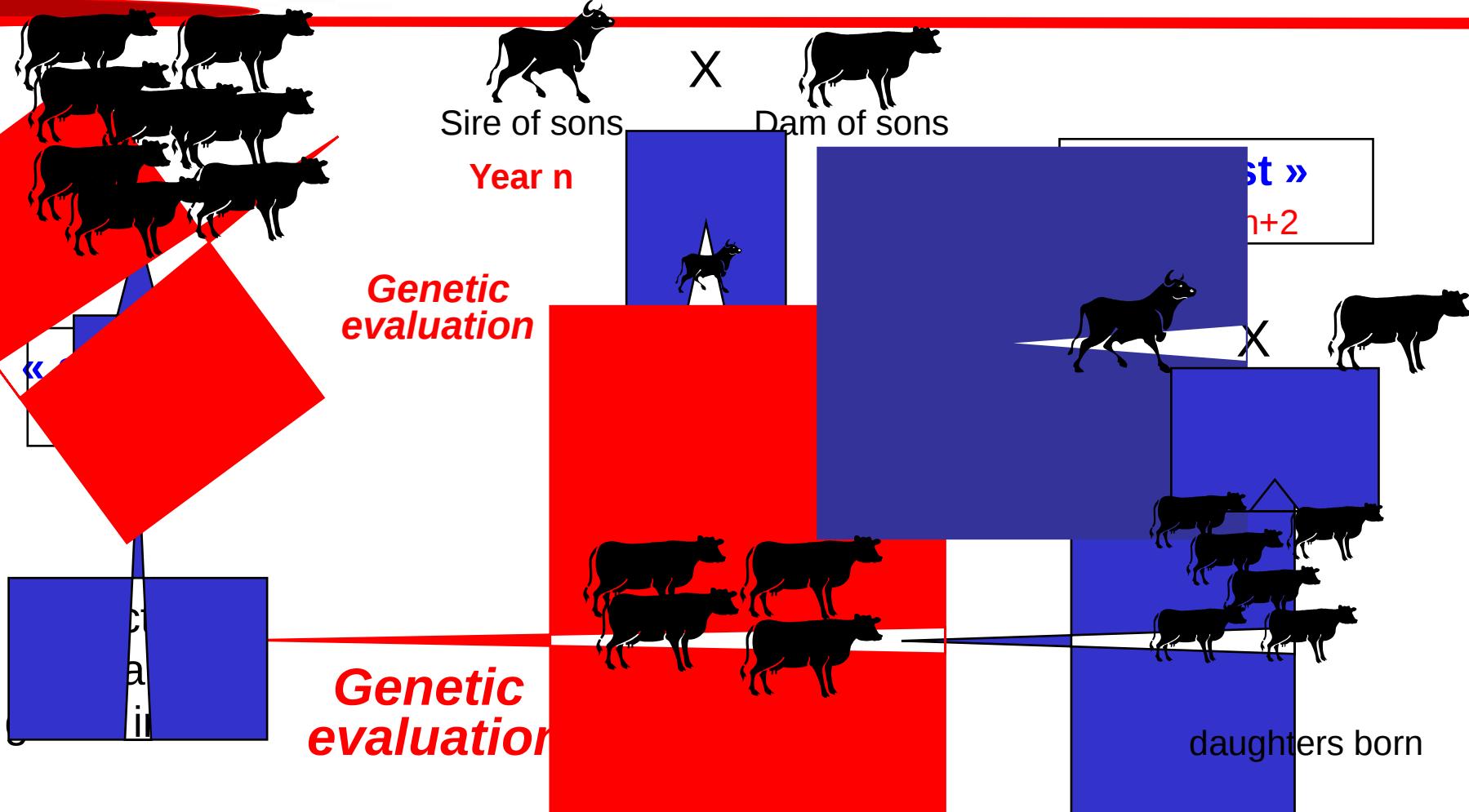
Digression: don't do this !



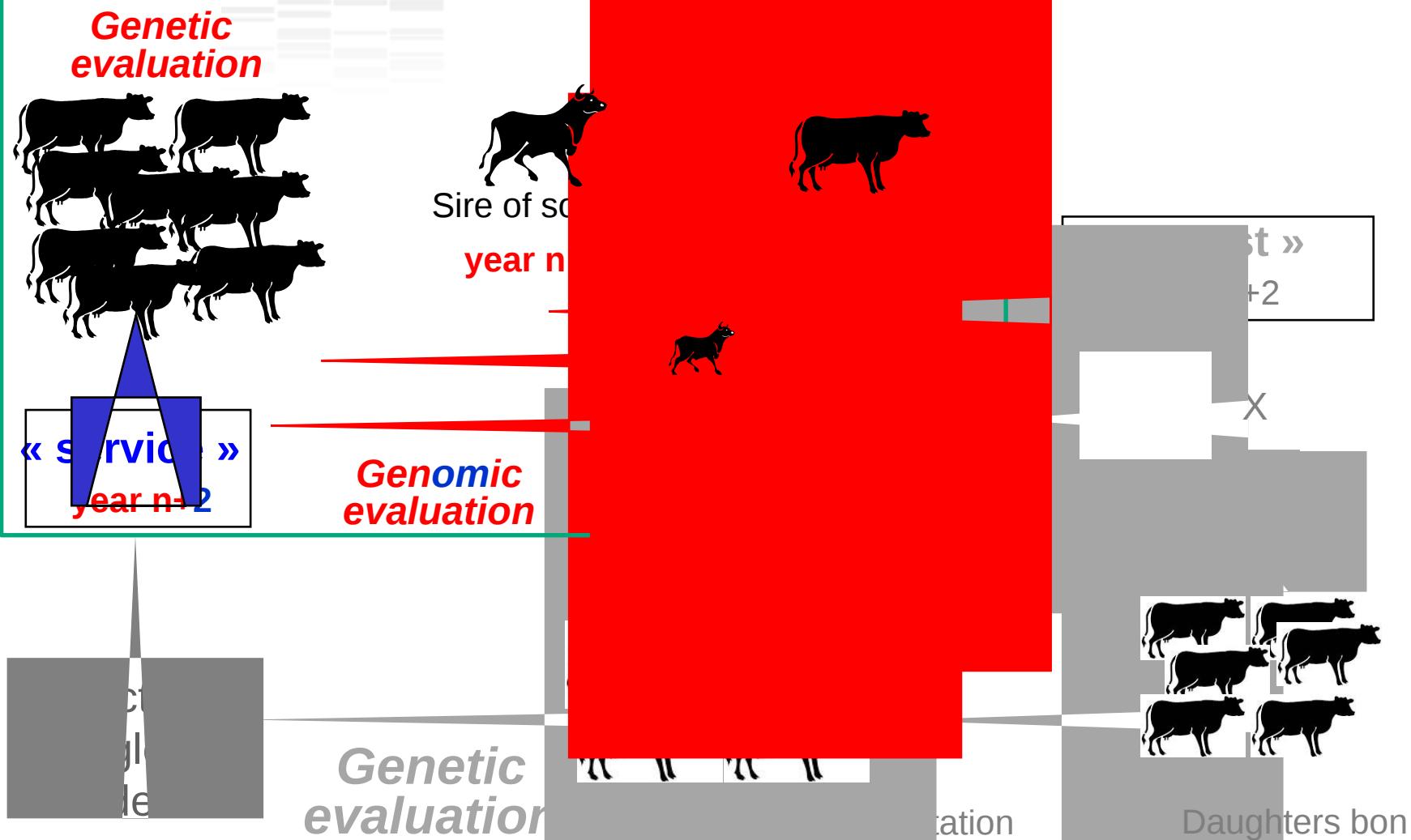
and another set of validation bulls; using stored semen | the Netherlands, Nordic countries, Poland and Spain

Remember: “Classical” progeny test

Genetic evaluation



A typical genomic breeding scheme



In practice ...

- Genomic selection
 - Requires **phenotypes and genotypes**
 - To minimize costs, genotype bulls in countries with large existing progeny test programs
 - But bulls don't have phenotypes (of interest)
 - ➔ summarize daughters' phenotypes at bull level
 - ✓ *Advantage* : more precise phenotypes

Basic principle for implementation

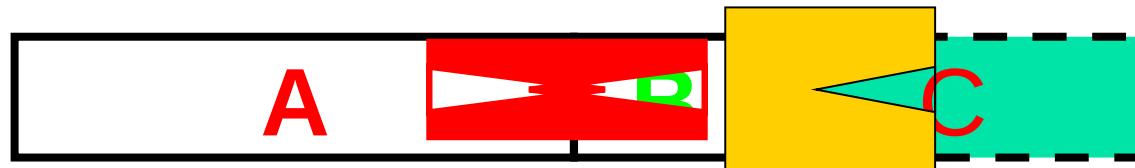
A **reference population A+B** : several thousands of genotyped animals with phenotype / performances

- From population A (**training population**), find the relationship between the phenotypes and the markers (=a prediction equation)
- **purely statistical approach** (no/very few genetic assumptions)



Basic principle for implementation

- then **validate** this prediction equation in **population B** (**validation population**)
- If satisfactory, redo the analysis on A+B and **apply** the resulting prediction equation to obtain a « genomic breeding value » for selection candidates C, for example at birth

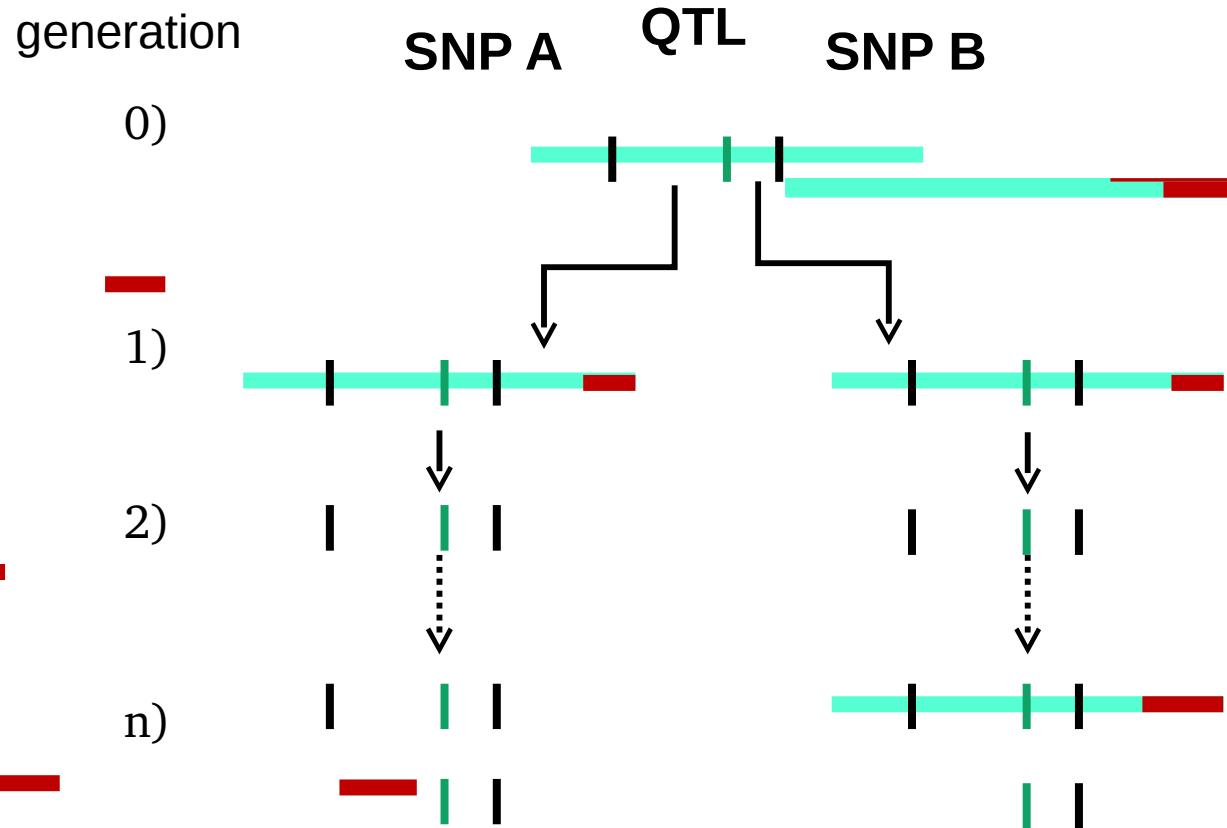


Genomic Selection: efficiency

- Two main factors :
 - Accuracy of SNP effect estimation
 - size of reference population
 - heritability of the trait
 - statistical methodology used
 - Linkage Disequilibrium (LD) between markers and QTL
 - marker density
 - effective size of the population => number of « independent » segments
 - Relationship between candidates and reference population

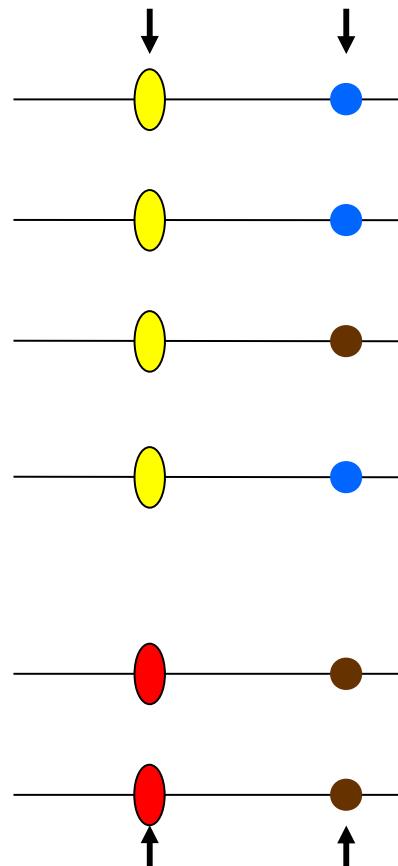
Linkage Disequilibrium

QTL-SNP association transmitted over generations



Linkage Disequilibrium

QTL frequent allele SNP



QTL rare allele SNP

SNP ● more often associated with QTL allele

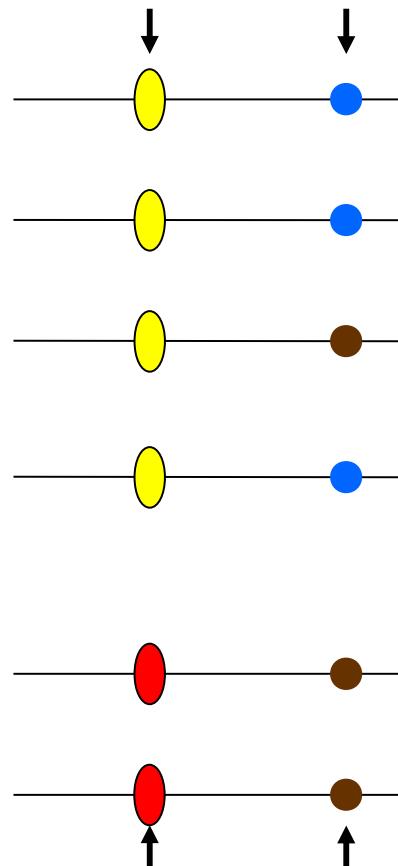
SNP ● more often associated with QTL allele



use effect of ● as proxy of effect of ○
and effect of ● as proxy of effect of ○

Linkage Disequilibrium

QTL frequent allele SNP



QTL rare allele SNP

SNP ● more often associated with QTL allele

SNP ● more often associated with QTL allele

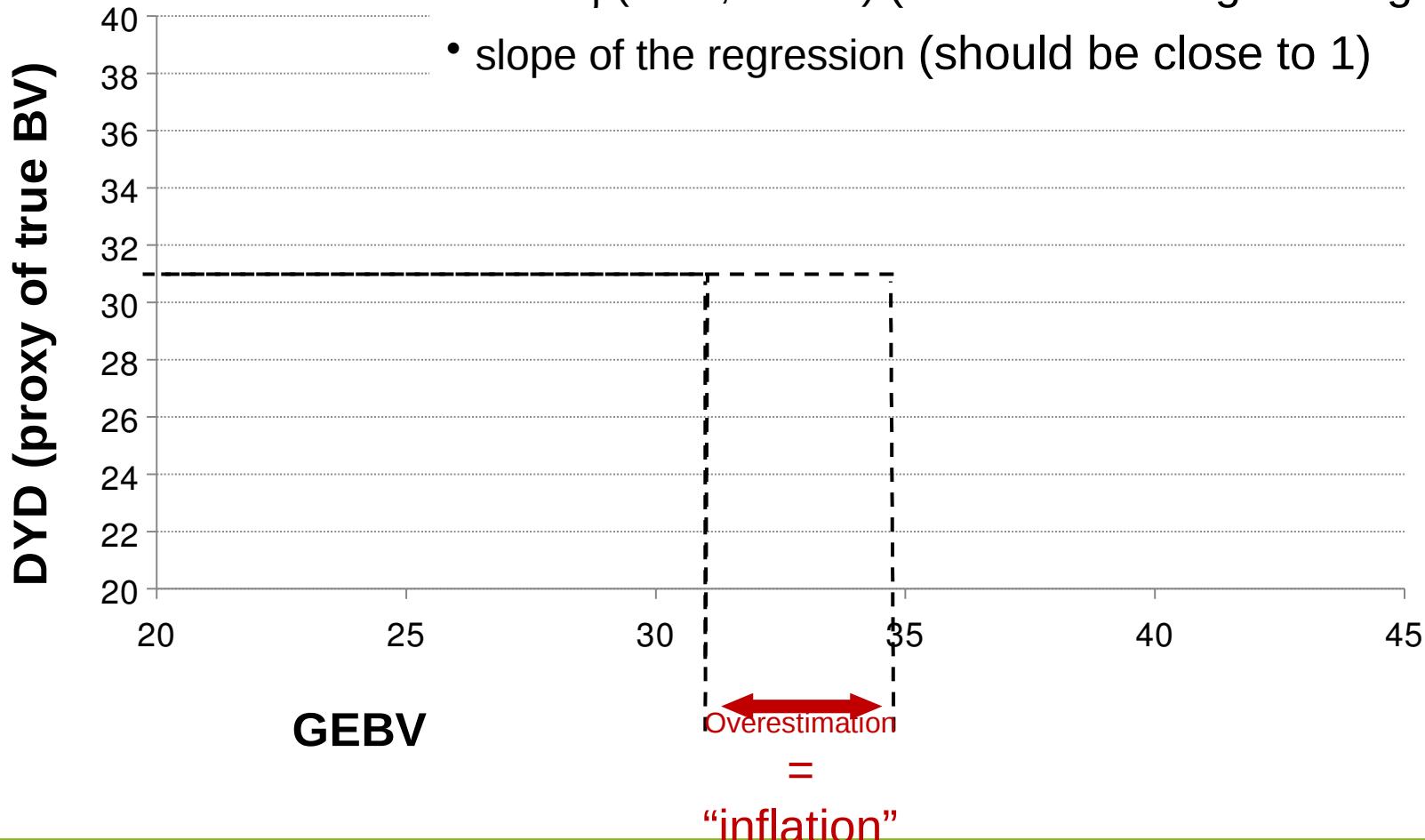


use effect of ● as proxy of effect of ○
and effect of ● as proxy of effect of ○

Validation test

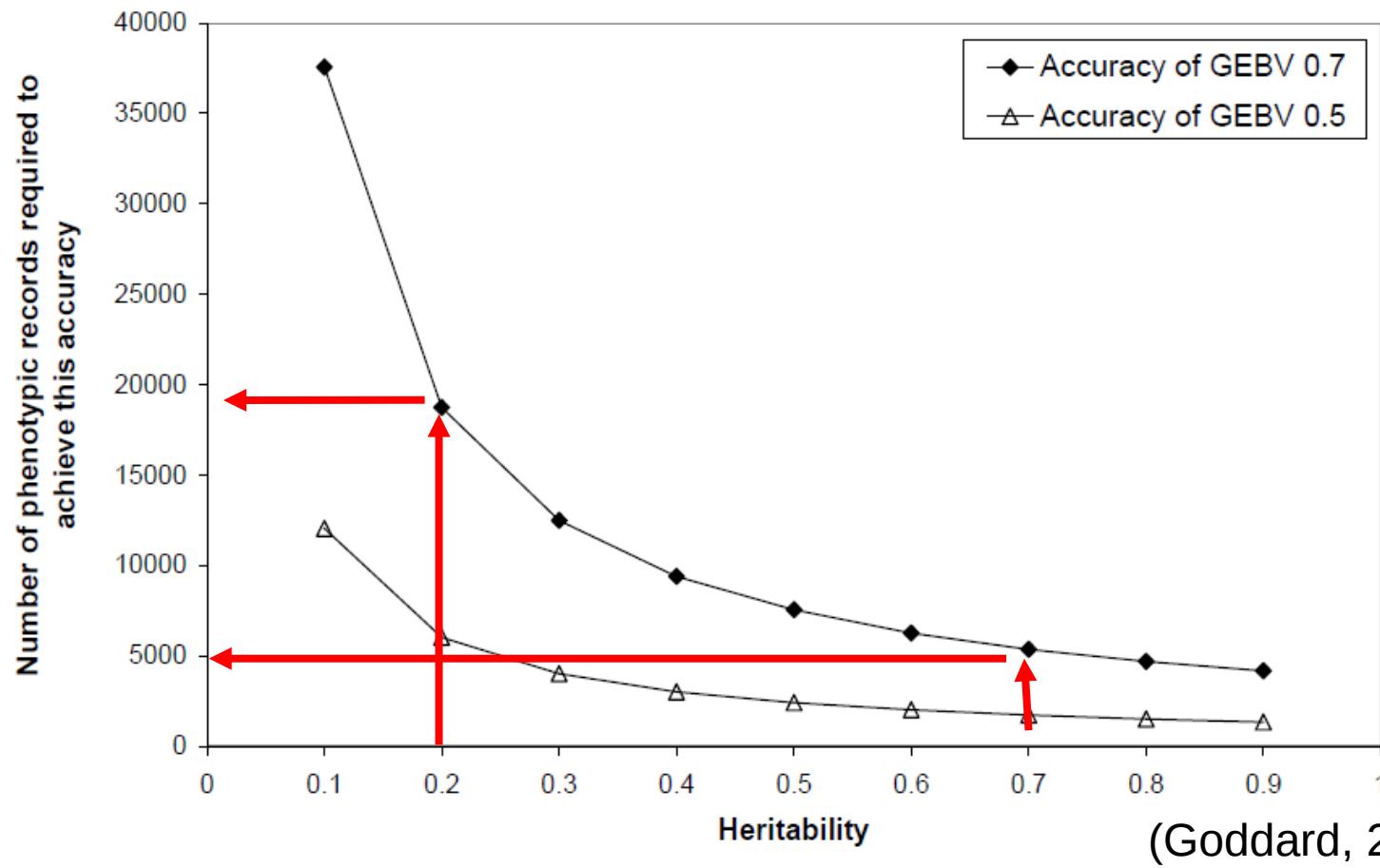
Two important parameters :

- R^2 or $p(DYD, GEBV)$ (should be « large enough »)
- slope of the regression (should be close to 1)



Size of the reference population

- Greatly influences the accuracy of genomic evaluations



Increase size of reference populations

- Genotype as many progeny tested bulls as possible
- International collaborations
 - **Holstein**: 2 big consortia
 - **USA + Canada** ~**>25000** bulls (?) + Italy + GB
 - **Eurogenomics** France, The Netherlands, (Germany), Nordic countries, Spain, Poland
~**34000** bulls



- **For small breeds**: not enough AI bulls
 - combine with many genotyped cows

Information from bulls and cows

Number of cows with 1 own record
equivalent to a bull with performances on daughters

CD*	h ²				
	0.1	0.2	0.3	0.4	0.5
0.40	6.0	2.7	1.6	1.0	0.7
0.50	9.0	4.0	2.3	1.5	1.0
0.60	13.5	6.0	3.5	2.3	1.5
0.70	21.0	9.3	5.4	3.5	2.3
0.80	36.0	16.0	9.3	6.0	4.0
0.90	81.0	36.0	21.0	13.5	9.0

* CD : coefficient of determination

= reliability of EBV based on progeny performance only (Boichard, 2015)

Basic data for genomic evaluation

- $y = YD$ (Yield deviation) = Individual record corrected for all fixed and non genetic random effects
 - or*
- $y = DYD$ (Daughter yield deviation) = (2 x) average for a bull of all YD of their daughters corrected for $\frac{1}{2}$ genetic effect of their dam (with associated weight = EDC (Equivalent Daughter Contribution))
 - or*
- $y = \text{deregressed proofs}$ (= DYD reconstructed from EBV, EDC and relationships: a genetic evaluation based on deregressed proofs will give back the bulls' EBV))
 - or*
- $y = \text{EBV}$ (~~← proxy → of DYD, valid only if reliabilities are high~~)

Different genomic evaluation methods

Can be grouped into 4 families:

- 1. Derived from BLUP: GBLUP**
- 2. General multidimension methods when the number of unknowns (p) is much larger than the number of observations (n) : the « p >>n problem » (PLS, LASSO, Elastic Net...)**
- 3. Bayesian methods A, B, C, C π , D, R ...**
- 4. « Machine learning » methods (big black box)**

To be applied on individual **SNP** or **haplotypes**

GBLUP: Basic « genomic » model

Hypotheses :

- 1) QTLs (quantitative Trait loci) have a purely additive effect
- 2) All the variability at QTLs is explained by the markers (!)

N biallelic markers

2 alleles A/B but only one effect defined = substitution effect mi

$$w_i = \frac{x_i m_i - \bar{x} m_i}{\sqrt{N}}$$

example:	AA	-1mi	or	0	or	-2pi mi
	AB	0		+mi		(1 - 2pi) mi
	BB	+1mi		+2mi		(2 - 2pi) mi

where pi is the allelic frequency of allele B of marker i



November 20-24 2017

vincent.ducrocq@inra.fr

Genomic evaluation

Distribution of marker effects

$$g = BV = \sum_{i=1}^N x_i m_i$$

$$\text{? GEBV} = \sum_{i=1}^N x_i \hat{m}_i$$

With the third form, and assuming that each marker contributes the same way to the variance of the trait, we have (Van Raden, 2009):

$$\sigma_m^2 \quad \sigma_m^2 = \frac{\sigma_g^2}{2 \sum_j p_j(1-p_j)}$$

$m_i \sim N(0, \sigma_m^2)$ with

Why this formula?: for consistency with the additive genetic variance. But what allelic frequency should be used? Quite a

few were proposed (observed ones/ in the base population..)



Evaluation model = BLUP on markers

If all animals with records are genotyped, and if we take allelic frequencies into account by centering x_i :

$w_i = x_i - 2\pi$ (grouped in a matrix W) , we can write:

The corresponding Mixed Model Equations are:

$$\begin{matrix} X'X & X'W \\ W'X & W'W + \frac{\sigma_e^2}{\sigma_m^2} \end{matrix} \begin{matrix} ? \\ ? \\ ? \\ ? \end{matrix} \begin{matrix} \hat{m} \\ ? \\ ? \\ ? \end{matrix} = \begin{matrix} X'y \\ W'y \end{matrix} \begin{matrix} ? \\ ? \end{matrix}$$

For any new animal * without performance: $GEBV = \sum_{i=1}^N w_i^* \hat{m}_i$

An equivalent model

- We can also write: $\mathbf{g} = \mathbf{Wm}$

But, since we have $m_i \sim N(\mathbf{0}_m^T, \sigma_m^2 I_m)$, we also have:

$$\mathbf{g} \sim N(\mathbf{0}_g^T, \mathbf{G})$$

$$\begin{aligned} \text{with } \mathbf{G} &= \text{var}(\mathbf{Wm}) = \mathbf{W} \text{ var}(\mathbf{m}) \mathbf{W}' \\ &= \mathbf{WW}'/k \end{aligned}$$

where k is a standardisation factor connecting the marker effects variance to the additive genetic variance

$$k = 2 \sum_j p_j (1 - p_j)$$

An equivalent model

- **G** plays the same role as the classical relationship matrix **A**
G is called « genomic relationship matrix »
- So, reparameterizing the model
 - if all recorded animals are genotyped:
 - or more generally (adding columns of 0 to **Z** if no phenotype):

GBLUP

The Mixed Model Equations are the same as in the classical case except that the relationship matrix \mathbf{A} (=« expectation » of relationship coefficients) is replaced by the genomic relationship matrix \mathbf{G} :

$$\begin{matrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \frac{\sigma_e^2}{\sigma_g^2} \mathbf{G}^{-1} \end{matrix} \begin{matrix} \mathbf{?} & \mathbf{?} & \mathbf{?} & \mathbf{?} \\ \mathbf{?} & \mathbf{?} & \mathbf{?} & \mathbf{?} \\ \mathbf{?} & \mathbf{?} & \mathbf{?} & \mathbf{?} \\ \mathbf{?} & \mathbf{?} & \mathbf{?} & \mathbf{?} \end{matrix} = \begin{matrix} \mathbf{?} & \mathbf{?} & \mathbf{?} & \mathbf{?} \\ \mathbf{?} & \mathbf{?} & \mathbf{?} & \mathbf{?} \\ \mathbf{?} & \mathbf{?} & \mathbf{?} & \mathbf{?} \\ \mathbf{?} & \mathbf{?} & \mathbf{?} & \mathbf{?} \end{matrix} \begin{matrix} \mathbf{X}'\mathbf{y} \\ \mathbf{Z}'\mathbf{y} \end{matrix}$$

$\hat{\mathbf{g}}$

(some) limits of GBLUP

- Assuming that **G** is correct, how to compute it efficiently?
 $\mathbf{G} = \mathbf{W}\mathbf{W}'/k$ **W** is dense!
Very long to compute, can be too big to store and to invert
 - iterative solutions
- How to compute reliabilities □ approximations
- A marker model seems more attractive.... (at least, its complexity does not increase too much when the number of genotyped animals increases!)

(some) limits of GBLUP

- It is not realistic to think that all the genetic variation at QTLs is captured by markers. This leads to genomic values which show «inflation » □ requires an additional « residual » polygenic effect:

Depending on the trait, $\omega = 5$ to 30%

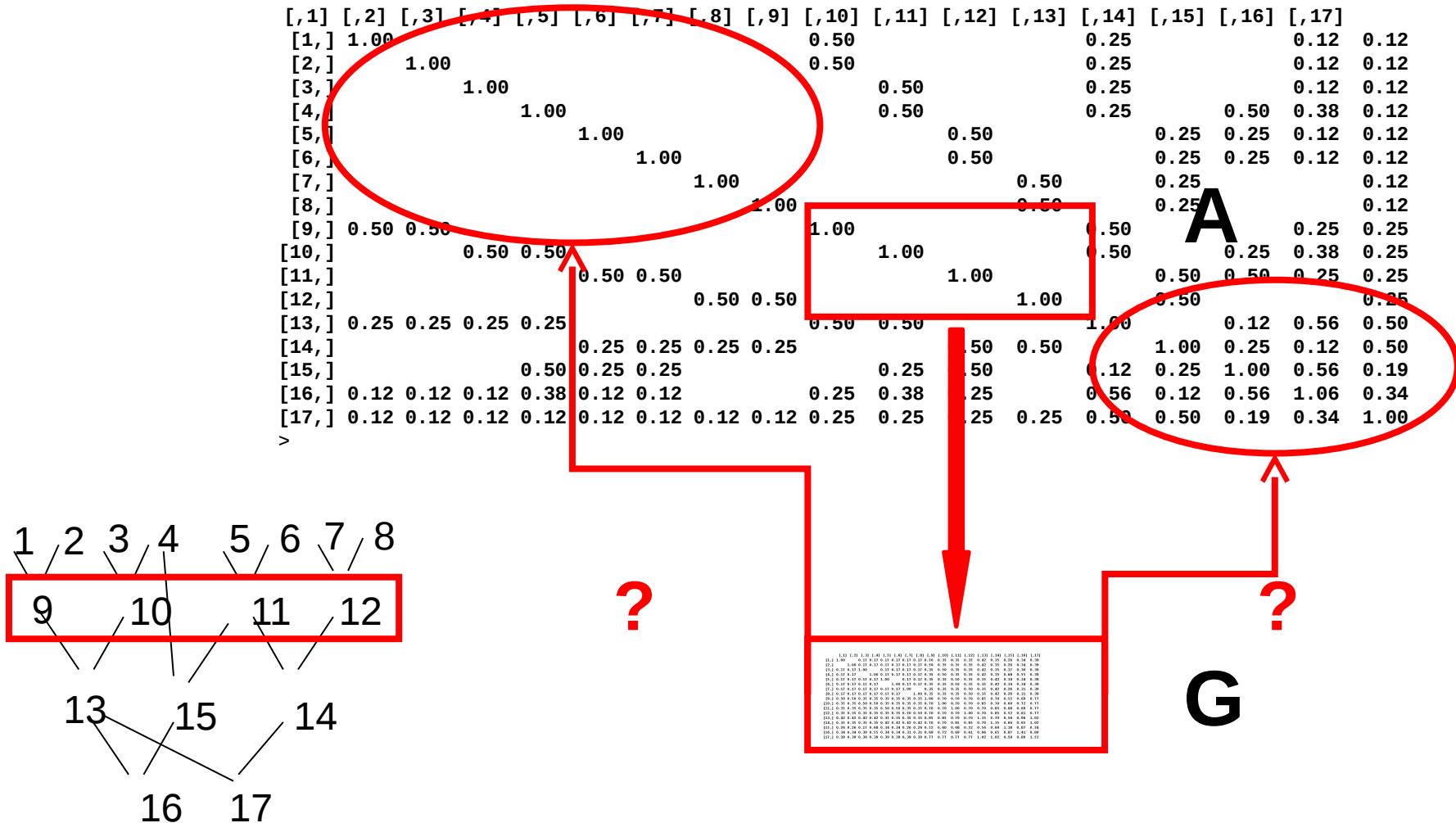
(some) limits of GBLUP

- How to combine genomic and classical (e.g., phenotypes from progeny) information?

Van Raden's (2009) approach:

- 1/ run a classical genetic evaluation with all animals
 - 2/ run a genomic evaluation with all genotyped animals
 - 3/ rerun the same (classical) evaluation with genotyped animals
...(note: BLUP hypotheses are not fulfilled!)
- 4/ combine the 3 sources of information (more or less: 1 + 2 -3)
with selection index theory, for each animal
- How to make non genotyped animals benefit from genomic information of relatives?

An example



Transmission of information (to non genotyped animals)

Let \mathbf{g}_2 be the genetic (genomic) values of **genotyped** animals and \mathbf{g}_1 the genetic values of **non genotyped** animals

Let $\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}$ and $p(\mathbf{g}_2) \sim N(0, \mathbf{G}\sigma_g^2)$

By regression of \mathbf{g}_1 on \mathbf{g}_2 :

$$p(\mathbf{g}_1 | \mathbf{g}_2) \sim N(\mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{g}_2, \mathbf{V}\sigma_g^2)$$

$$\text{with } \mathbf{V} = (\mathbf{A}^{11})^{-1} = \mathbf{A}_{11} + \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$$

Transmission of information (to non genotyped animals)

Let **g2** be the genetic (genomic) values of **genotyped** animals and
g1 the genetic values of **non genotyped** animals

An estimate of **g1** based on genomic information is obtained
by regression of **g1** on **g2**

This information is added to regular BLUP equations

«Single-Step» approach

Legarra, Misztal, Aguilar; Christensen & Lund; D. Johnson

- $H = \text{Variance of }$



«Single-step» approach

$$H = \begin{matrix} H_{11} & H_{12} & H_{21} \\ H_{12} & H_{22} & \end{matrix} = \begin{matrix} ? & ? & ? \\ ? & ? & ? \end{matrix} = A_{11} + A_{12}A_{22}^{-1}(G - A_{22})A_{22}^{-1}A_{21}$$

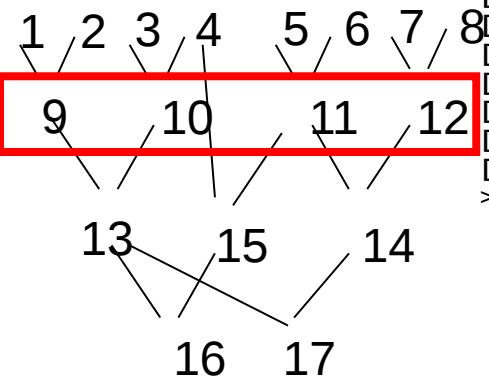
non genotyped *genotyped*

$$GA_{22}^{-1}A_{21} \quad \quad \quad A_{12}A_{22}^{-1}G \quad \begin{matrix} ? & ? \\ ? & ? \end{matrix}$$
$$G \quad \quad \quad G \quad \begin{matrix} ? & ? \\ ? & ? \end{matrix}$$

- Incredible but true ! H^{-1} has a relatively simple form:

$$H^{-1} = A^{-1} + \begin{bmatrix} 0 & 0 \\ 0 & G^{-1} - A_{22}^{-1} \end{bmatrix}$$

Back to the example



```
[,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13] [,14] [,15] [,16] [,17]
[1,] 1.00 0.50 0.25 0.12 0.12
[2,] 1.00 0.50 0.25 0.12 0.12
[3,] 1.00 0.50 0.25 0.12 0.12
[4,] 1.00 0.50 0.25 0.50 0.38 0.12
[5,] 1.00 0.50 0.25 0.25 0.12 0.12
[6,] 1.00 0.50 0.25 0.25 0.12 0.12
[7,] 1.00 0.50 0.25 0.12 0.12
[8,] 1.00 0.50 0.25 0.12 0.12
[9,] 0.50 0.50 1.00 0.50 0.25 0.25 0.25
[10,] 0.50 0.50 1.00 0.50 0.25 0.38 0.25
[11,] 0.50 0.50 1.00 0.50 0.25 0.25 0.25
[12,] 0.50 0.50 1.00 0.50 0.25 0.25 0.25
[13,] 0.25 0.25 0.25 0.25 0.50 0.50 1.00 0.12 0.56 0.50
[14,] 0.25 0.25 0.25 0.25 0.50 0.50 1.00 0.25 0.12 0.25 0.12 0.50
[15,] 0.50 0.25 0.25 0.25 0.25 0.50 0.50 0.12 0.25 1.00 0.56 0.19
[16,] 0.12 0.12 0.12 0.38 0.12 0.12 0.25 0.38 0.25 0.56 0.12 0.56 1.06 0.34
[17,] 0.12 0.12 0.12 0.12 0.12 0.12 0.25 0.25 0.25 0.25 0.50 0.50 0.19 0.34 1.00
>
```

```
[,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13] [,14] [,15] [,16] [,17]
[1,] 1.00 0.17 0.17 0.17 0.17 0.17 0.17 0.17 0.50 0.35 0.35 0.35 0.42 0.35 0.26 0.34 0.39
[2,] 1.00 0.17 0.17 0.17 0.17 0.17 0.17 0.17 0.50 0.35 0.35 0.35 0.42 0.35 0.26 0.34 0.39
[3,] 0.17 0.17 1.00 0.17 0.17 0.17 0.17 0.35 0.50 0.35 0.35 0.42 0.35 0.17 0.30 0.39
[4,] 0.17 0.17 1.00 0.17 0.17 0.17 0.17 0.35 0.50 0.35 0.35 0.42 0.35 0.68 0.55 0.39
[5,] 0.17 0.17 0.17 1.00 0.17 0.17 0.35 0.35 0.50 0.35 0.35 0.42 0.34 0.34 0.34 0.39
[6,] 0.17 0.17 0.17 0.17 1.00 0.17 0.17 0.35 0.50 0.35 0.35 0.42 0.34 0.34 0.34 0.39
[7,] 0.17 0.17 0.17 0.17 0.17 1.00 0.35 0.35 0.50 0.35 0.35 0.42 0.26 0.31 0.39
[8,] 0.17 0.17 0.17 0.17 0.17 0.17 1.00 0.35 0.35 0.50 0.35 0.42 0.26 0.31 0.39
[9,] 0.50 0.50 0.35 0.35 0.35 0.35 0.35 1.00 0.70 0.70 0.70 0.85 0.70 0.52 0.69 0.77
[10,] 0.35 0.35 0.50 0.50 0.35 0.35 0.35 0.70 1.00 0.70 0.70 0.85 0.70 0.60 0.72 0.77
[11,] 0.35 0.35 0.35 0.50 0.50 0.35 0.35 0.70 0.70 1.00 0.70 0.70 0.85 0.68 0.69 0.77
[12,] 0.35 0.35 0.35 0.35 0.35 0.50 0.50 0.70 0.70 0.70 1.00 0.70 0.85 0.52 0.61 0.77
[13,] 0.42 0.42 0.42 0.42 0.35 0.35 0.35 0.85 0.85 0.70 0.70 1.35 0.70 0.56 0.96 1.02
[14,] 0.35 0.35 0.35 0.42 0.42 0.42 0.42 0.70 0.70 0.85 0.85 0.70 1.35 0.60 0.65 1.02
[15,] 0.26 0.26 0.17 0.68 0.34 0.34 0.26 0.52 0.60 0.68 0.52 0.56 0.60 1.18 0.87 0.58
[16,] 0.34 0.34 0.30 0.55 0.34 0.34 0.31 0.31 0.69 0.72 0.69 0.61 0.96 0.65 0.87 1.41 0.80
[17,] 0.39 0.39 0.39 0.39 0.39 0.39 0.39 0.77 0.77 0.77 0.77 1.02 1.02 0.58 0.80 1.52
```



«Single-step» approach

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{Z}\mathbf{a} + \mathbf{e} \quad \square \text{ (G)BLUP on all data}$$

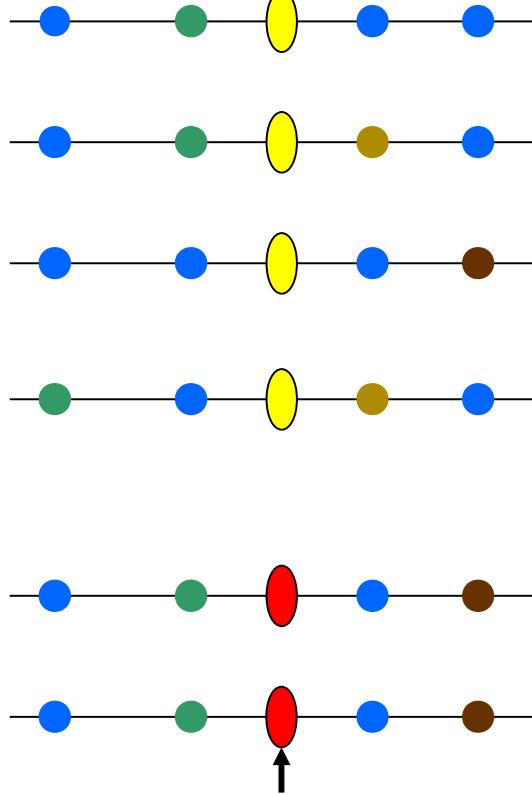
$$\begin{array}{ccc} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} & \begin{matrix} ? & ? & ? \\ ? & ? & ? \\ ? & ? & ? \end{matrix} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \alpha \mathbf{H}^{-1} & \begin{matrix} ?' & \mathbf{g} & ? \\ ? & ? & ? \\ ?' & ?' & ? \end{matrix} \end{array}$$

- Huge (and relatively dense) system but with some tricks, can be applied to large data sets (Legarra)
- Avantages:
 - Conceptually very simple (=BLUP)
 - quasi-optimal use of information from genotyped animals
 - automatically accounts for (genomic) selection

«Single-step» approach: current issues

Haplotypes vs SNP

QTL frequent allele



Increase linkage disequilibrium marker(s)-QTL

- Requires marker phasing (but imputation)
- more effects to estimate

QTL rare allele

The French genomic model (since April 2015)

QTL size	Genomic evaluation
Large	
Moderate	traced with markers <input type="checkbox"/> haplotype effects
Small	
Tiny	Consider their sum only: $\hat{u} = \sum_{j'} \hat{m}_{j'} \sim N(0, \cancel{\text{pedigree relationship matrix}}) \sim N(0, \text{genomic relationship matrix})$

SNP & Haplotypes

$$g_i = u_i + \sum_{j=1}^n (h_{ij1} + h_{ij2}) + \sum_{j=1}^k (SNP_{ij1} + SNP_{ij2})$$

The equation is divided into two main parts by a plus sign. The first part, enclosed in a red box, is labeled "Trait dependent". The second part, enclosed in a red oval, is labeled "Trait independent". A magnifying glass icon is positioned above the trait-independent term.

Trait dependent

Trait independent

Comparison of genomic evaluation methods

- correlation between GEBV and DYD in validation population

