

## Genomic Prediction



## Genomic Prediction

1. Fit vs. Prediction in modern genomics
2. The large p, small n paradigm
3. Variable selection vs. Ridge regression methods
4. Merging data: Single Step
5. GS using sequence data

Copyright © 2001 by the Genetics Society of America

## Prediction of Total Genetic Value Using Genome-Wide Dense Marker Maps

T. H. E. Meuwissen,<sup>\*</sup> B. J. Hayes<sup>†</sup> and M. E. Goddard<sup>†,‡</sup>

<sup>\*</sup>Research Institute of Animal Science and Health, 8200 AB Lelystad, The Netherlands, <sup>†</sup>Victorian Institute of Animal Science, Attwood 3049, Victoria, Australia and <sup>‡</sup>Institute of Land and Food Resources, University of Melbourne, Parkville 3052, Victoria, Australia

Manuscript received August 17, 2000  
Accepted for publication January 17, 2001

### ABSTRACT

Recent advances in molecular genetic techniques will make dense marker maps available and genotyping many individuals for these markers feasible. Here we attempted to estimate the effects of ~50,000 marker haplotypes simultaneously from a limited number of phenotypic records. A genome of 1000 cM was simulated with a marker spacing of 1 cM. The markers surrounding every 1-cM region were combined into marker haplotypes. Due to finite population size ( $N = 100$ ), the marker haplotypes were in linkage disequilibrium with the QTL located between the markers. Using least squares, all haplotype effects could not be estimated simultaneously. When only the biggest effects were included, they were overestimated and the accuracy of predicting genetic values of the offspring of the recorded animals was only 0.32. Best linear unbiased prediction of haplotype effects assumed equal variances associated to each 1-cM chromosomal segment, which yielded an accuracy of 0.73, although this assumption was far from true. Bayesian methods that assumed a prior distribution of the variance associated with each chromosome segment increased this accuracy to 0.85, even when the prior was not correct. It was concluded that selection on genetic values predicted from markers could substantially increase the rate of genetic gain in animals and plants, especially if combined with reproductive techniques to shorten the generation interval.

## Principles of Genomic Prediction

- ❖ It consists of using molecular data to predict phenotypic performance.
- ❖ It was formally proposed by T Meuwissen, B Hayes and M Goddard in 2001.
- ❖ It has revolutionized animal breeding, dairy cattle in particular, resulting in increases in accuracy >30%

## Advantages of Genomic Prediction

- ❖ It increases accuracy for individuals without phenotypes (e.g., newborn).
- ❖ It decreases generation interval.
- ❖ Computes relationship (similarity) without the need for a pedigree, e.g., distinguishes between full sibs.

## Limitations

- ❖ Expensive, also in terms of logistics.
- ❖ It requires large datasets.
- ❖ Difficult to predict across breeds.
- ❖ Not clearly beneficial in all scenarios.

*But let's go back a little bit . . .*



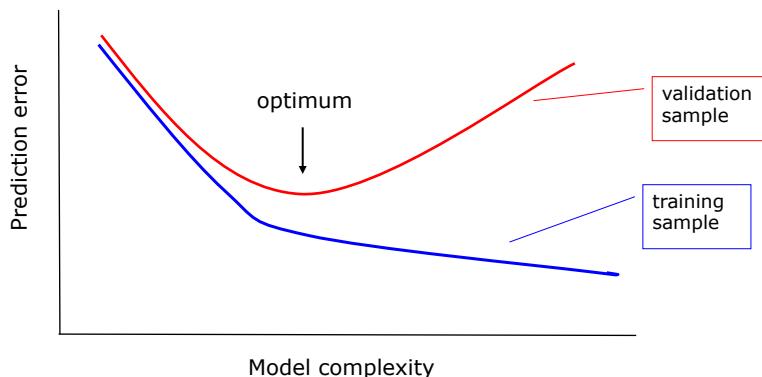
## Recall Statistics roles

1. For the largest part of 20th century, Statistics predominant role has been inference, i.e., inferring unknown parameters, usually **model** based.
2. More recently, **prediction** has become central in Statistics, all the more with the enormous amounts of easily obtained data and powerful computer based methods.
3. Because of that, machine learning has evolved independently – occasionally converging with statistics.

## Two desirable characteristics: parsimony and goodness of fit

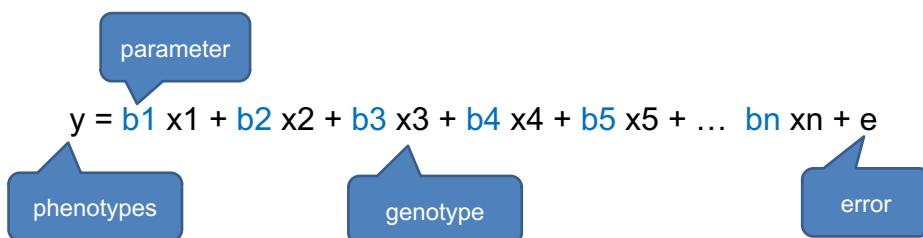
- ✓ For low  $p$  (# variables), goodness of fit is the most important issue.
- ✓ For high  $p$ , parsimony becomes critical (collinearity problems).
- ✓ For  $p > n$  (# observations), additional restrictions on the solutions must be posed. This is called **prior information** (Bayesian Statistics), **shrinkage** (Frequentist Statistics) or **regularization** (Machine Learning jargon)
- ✓ For  $n \gg p$ , we are usually no longer interested in inference but only in prediction.

## A modeling tradeoff



Hastie et al. Elements of Statistical Learning

## Large p small n paradigm



Under the normal setting of breeding in which more variables than data are available, only two solutions exist (or a combination).

## Large p small n paradigm

$$y = b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 x_4 + b_5 x_5 + \dots \dots \dots b_n x_n + e$$

Under the normal setting of breeding in which more variables than data are available, only two solutions exist (or a combination):

$$y = b_1 x_1 + b_2 \cancel{x_2} + b_3 x_3 + b_4 \cancel{x_4} + b_5 \cancel{x_5} + \cancel{\dots} \dots b_n x_n + e$$

Variable selection

$$y = b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 x_4 + b_5 x_5 + \dots \dots \dots b_n x_n + e$$

Shrinkage

## Large p small n paradigm

Either variable selection or shrinkage are two kinds of regularization (in Machine Learning terminology).

Type 1 (L1) imposes a penalty on the sum of absolute values of coefficients

$$\min_w \|Xw - y\|_2^2 + \alpha \|w\|_1$$

Lasso: results  
in variable  
selection

Type 2 (L2) imposes a penalty on the sum of squared coefficients

$$\min_w \|Xw - y\|_2^2 + \alpha \|w\|_2^2$$

Ridge  
regression:  
results in  
shrinkage

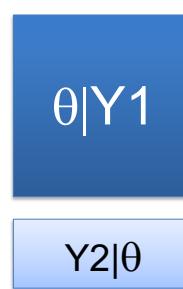
Both regularization and combined methods have been used in GS

## Prediction is what matters in genomic selection

Prediction, assessed by crossvalidation, is what matters. The population of reference is a critical choice here.



Whole dataset



Training  
population

Target  
population

$$y = \mu + \sum_{k=1}^m \mathbf{x}_k \beta_k + \epsilon$$

**Bayes A:** Each marker effect  $\beta$  follows, conditionally, a N distribution, variance is different for each marker, with prior inverted chi-2 distributed. Marginal of  $\beta$  is a t.

**Bayes B:** Same as Bayes A except that only a fraction of markers are assumed to have an effect.

**GBLUP:** Marker effects follow a normal distribution with equal variance for all markers.

$$y = \mu + \sum_{k=1}^m \mathbf{x}_k \beta_k + \epsilon$$

Bayes A

$$p(\beta_j, \sigma_{\beta_j}^2, S_\beta) = \left\{ \prod_k N(\beta_{jk} | 0, \sigma_{\beta_{jk}}^2) \chi^{-2}(\sigma_{\beta_{jk}}^2 | df_\beta, S_\beta) \right\} G(S_\beta | r, s)$$

Bayes B

$$p(\beta_j, \sigma_{\beta_j}^2, \pi) = \left\{ \prod_k \left[ \pi N(\beta_{jk} | 0, \sigma_{\beta_j}^2) + (1 - \pi) 1(\beta_{jk} = 0) \right] \right\} \times \chi^{-2}(\sigma_{\beta_j}^2 | df_\beta, S_\beta) B(\pi | p_0, \pi_0)$$

GBLUP

$$p(\beta_j, \sigma_{\beta_j}^2) = \left\{ \prod_k N(\beta_{jk} | 0, \sigma_{\beta_j}^2) \right\} \chi^{-2}(\sigma_{\beta_j}^2 | df_\beta, S_\beta)$$

$$y = \mu + \sum_{k=1}^m \mathbf{x}_k \beta_k + \sum_{l=1}^L u_l + \varepsilon$$

Likelihood

$$p(\mathbf{y}|\boldsymbol{\theta}) = \prod_{i=1}^n N\left(y_i | \mu + \sum_{j=1}^J \sum_{k=1}^{K_j} x_{ijk} \beta_{jk} + \sum_{l=1}^L u_{li}, \sigma_\varepsilon^2 w_i^2\right)$$

Prior

$$p(\boldsymbol{\theta}) = p(\mu)p(\sigma_\varepsilon^2) \prod_{j=1}^J p(\beta_j) \prod_{l=1}^L p(u_l).$$

Genome-Wide Regression and Prediction with the BGLR Statistical Package  
Paulino Pérez and Gustavo de los Campos

---

model= Join distribution of effects and hyper-parameters

---

FIXED  $p(\beta_j) \propto 1$

---

BRR  $p(\beta_j, \sigma_\beta^2) = \left\{ \prod_k N(\beta_{jk}|0, \sigma_\beta^2) \right\} \chi^{-2}(\sigma_\beta^2 | df_\beta, S_\beta)$

---

BayesA  $p(\beta_j, \sigma_{\beta_j}^2, S_\beta) = \left\{ \prod_k N(\beta_{jk}|0, \sigma_{\beta_{jk}}^2) \chi^{-2}(\sigma_{\beta_{jk}}^2 | df_\beta, S_\beta) \right\} G(S_\beta | r, s)$

---

$p(\beta_j, \tau_j^2, \lambda^2 | \sigma_\varepsilon^2) = \left\{ \prod_k N(\beta_{jk}|0, \tau_{jk}^2 \times \sigma_\varepsilon^2) \text{Exp} \left\{ \tau_{jk}^2 | \frac{\lambda^2}{2} \right\} \right\} \times G(\lambda^2 | r, s) , \text{ or}$

BL  $p(\beta_j, \tau_j^2, \lambda | \sigma_\varepsilon^2, \max) = \left\{ \prod_k N(\beta_{jk}|0, \tau_{jk}^2 \times \sigma_\varepsilon^2) \text{Exp} \left\{ \tau_{jk}^2 | \frac{\lambda^2}{2} \right\} \right\} \times B(\lambda / \max | p_0, \pi_0), \text{ or}$

$p(\beta_j, \tau_j^2 | \sigma_\varepsilon^2, \lambda) = \left\{ \prod_k N(\beta_{jk}|0, \tau_{jk}^2 \times \sigma_\varepsilon^2) \text{Exp} \left\{ \tau_{jk}^2 | \frac{\lambda^2}{2} \right\} \right\}$

---

BayesC  $p(\beta_j, \sigma_\beta^2, \pi) = \left\{ \prod_k \left[ \pi N(\beta_{jk}|0, \sigma_\beta^2) + (1-\pi)1(\beta_{jk}=0) \right] \right\} \times \chi^{-2}(\sigma_\beta^2 | df_\beta, S_\beta) B(\pi | p_0, \pi_0)$

---

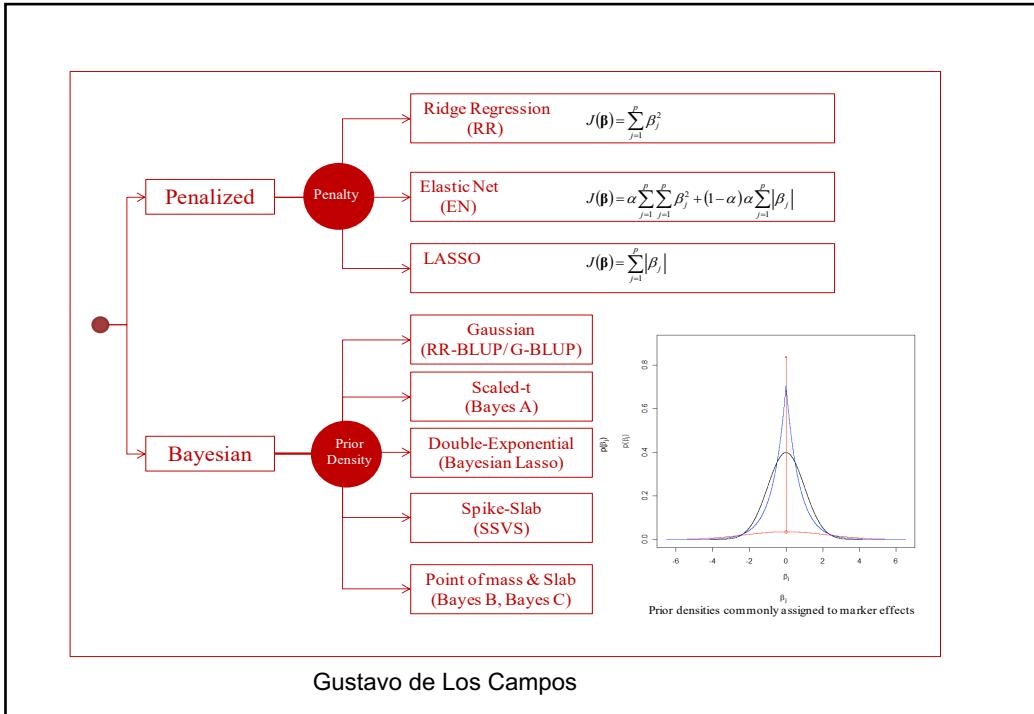
BayesB  $p(\beta_j, \sigma_\beta^2, \pi) = \left\{ \prod_k \left[ \pi N(\beta_{jk}|0, \sigma_\beta^2) + (1-\pi)1(\beta_{jk}=0) \right] \chi^{-2}(\sigma_{\beta_{jk}}^2 | df_\beta, S_\beta) \right\} B(\pi | p_0, \pi_0) \times G(S_\beta | r, s)$

---

RKHS  $p(\mathbf{u}_l, \sigma_{u_l}^2) = N(\mathbf{u}_l | \mathbf{0}, \mathbf{K}_l \times \sigma_{u_l}^2) \chi^{-2}(\sigma_{u_l}^2 | df_l, S_l)$

---

Genome-Wide Regression and Prediction with the BGLR Statistical Package  
Paulino Pérez and Gustavo de los Campos



## GBLUP (VanRaden)

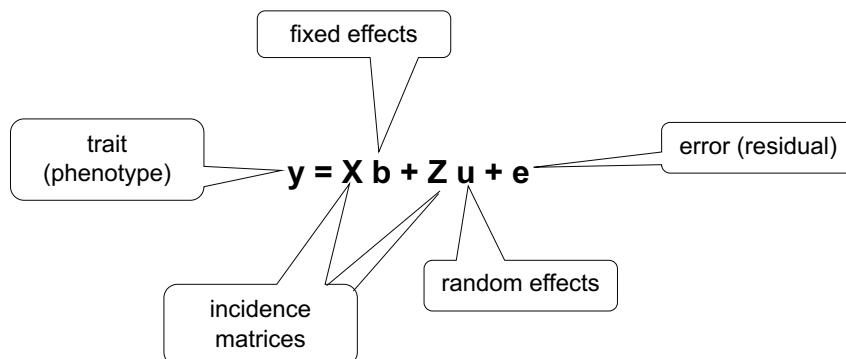
G-BLUP is a widely used method for genomic selection. It consists in replacing the pedigree based matrix by a molecular marker based one (GRM).

$$\mathbf{y} = \mathbf{X} \mathbf{b} + \mathbf{Z} \mathbf{u} + \mathbf{e}$$

$$\begin{pmatrix} \mathbf{X}' \mathbf{R}^{-1} \mathbf{X} & \mathbf{X}' \mathbf{R}^{-1} \mathbf{Z} \\ \mathbf{Z}' \mathbf{R}^{-1} \mathbf{X} & \mathbf{Z}' \mathbf{R}^{-1} \mathbf{Z} + \mathbf{G}^{-1} \end{pmatrix} \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}' \mathbf{R}^{-1} \mathbf{y} \\ \mathbf{Z}' \mathbf{R}^{-1} \mathbf{y} \end{pmatrix}$$

- There are several potential ways of calculating **G**.
- What is the relationship between **G** and **A** (pedigree NRM).
- **G** must be inverted brute force (good approx algorithms by Misztal).

## Usual models: Mixed Linear Models



## Usual models: Mixed Linear Models

$$y = Xb + Zu + e$$

$$\begin{pmatrix} y \\ u \\ e \end{pmatrix} \sim N \left( \begin{pmatrix} Xb \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} ZGZ' \sigma_u^2 + I\sigma_e^2 & ZG\sigma_u^2 & I\sigma_e^2 \\ GZ' \sigma_u^2 & G\sigma_u^2 & 0 \\ I\sigma_e^2 & 0 & I\sigma_e^2 \end{pmatrix} \right)$$

- In classical breeding,  $\text{Var}(u) = G$  is computed from the pedigree.
- In GS,  $G$  is computed using marker information.

## Genomic relationship matrix

$$G = \left\{ \frac{\sum (z - 2p)(z' - 2p)}{\sum 2pq} \right\}$$

Covariance  
matrix (averaged)

## In general, it seems that

- Most GP methods work similarly as size of data increases.
- Genetic architecture does influence the performance of the methods. In general, variable selection is better when the number of causal loci is small (Daetwyler et al., 2010).

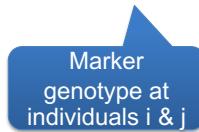
## Important problema in practice:

Many individuals are available but only a few are genotyped. All are related through a common pedigree

### Merging Information: Single Step Approach (Legarra, Aguilar, Misztal, JDS 2009; Mogens, Lund, GSE, 2010)

- ❖ It addresses the issue of a large connected pedigree where only a subset of animals are genotyped.
- ❖ It models genotypes as multivariate normal variables, using **A** and **G** as covariance matrices.
- ❖ It derives an expression for **A** of ungenotyped individuals given **G**.

$$\text{Cov}(z_i, z_j) = A_{ij} 2 p q$$



## Single Step Approach

Ungenotyped  
indivs

Genotyped  
indivs

$$\begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} + \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} + \boldsymbol{\varepsilon}$$

$$\text{Var}\left(\begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} \mid \text{pedigree}\right) = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix} \sigma_u^2$$

$$\text{Var}\left(\begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} \mid \text{marker}\right) = \begin{pmatrix} ? & ? \\ ? & \mathbf{G}_{22} \end{pmatrix} \sigma_g^2$$

## Single Step Approach

Ungenotyped  
indivs

Genotyped  
indivs

$$\begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} + \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} + \boldsymbol{\varepsilon}$$

$$p\left(\begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix}\right) = N\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}\right)$$

prior

Assumes prior  $\mathbf{A}$   
has no influence  
on  $\mathbf{u}_2$  (ie,  
molecular  
information fully  
dominates over  
pedigree)

$$p(\mathbf{u}_1, \mathbf{u}_2 \mid \text{markers}) = p(\mathbf{u}_1 \mid \mathbf{u}_2)p(\mathbf{u}_2 \mid \text{markers}) = N(\mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{u}_2, \sigma_u^2(\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})). N(0, \sigma_u^2\mathbf{G}).$$

'posterior'

$$= [\mathbf{A}_{21}\mathbf{A}_{22}^{-1}\mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{12} + \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}\mathbf{A}_{21}] \sigma_u^2$$

## Single Step Approach

Putting all together

$$\text{Var} \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} | \text{markers} = \mathbf{H} \sigma_u^2 = \sigma_u^2 \begin{pmatrix} \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} + \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{G} \\ \mathbf{G}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{G} \end{pmatrix}$$

$\mathbf{H}$  can be interpreted as a projection (ie, multivariate regression) of  $\mathbf{A}$  on  $\mathbf{G}$

Amazingly,  $\mathbf{H}$  has a very simple inverse

$$\mathbf{H}^{-1} = \mathbf{A}^{-1} + \begin{pmatrix} 0 & 0 \\ 0 & \mathbf{G}^{-1} - \mathbf{A}_{22}^{-1} \end{pmatrix}$$

## Single Step Approach

1. Automatic accounting of all relatives of genotyped individuals and their performances.
2. Simultaneous fit of genomic information and estimates of other effects (e.g., contemporary groups). Therefore not loss of information.
3. Feedback: the extra accuracy in genotyped individuals is transmitted to all their relatives (e.g. [Christensen et al., 2012](#)).
4. Simple extensions. Because this is a linear BLUP-like estimator, the extension to more complicated models (multiple trait, threshold traits, test day records) is immediate. Any model fit using relationship matrices can be fit using combined relationship matrices.
5. Analytical framework. The Single Step provides an analytical framework for further developments. This is notoriously difficult with pseudo-data.

[Legarra et al., 2014](#)

## Single Step Approach: Limitations

1. Programming complexity to fit complicated models for marker effects (Bayesian Regressions, machine learning algorithms, etc.).
2. Lack of experience on very large data sets.
3. Long computing times with current Single Step algorithms methods, for very large data sets.
4. Lack of an easy and elegant way of considering major genes in a multiple trait setting, this is a drawback of multiple step methods as well.

## Genomic Selection in Dairy Cattle: The USDA Experience\*

**Annual Review of Animal Biosciences**

Vol. 5:309-327 (Volume publication date February 2017)  
 First published online as a Review in Advance on November 16, 2016  
<https://doi.org/10.1146/annurev-animal-021815-111422>

**George R. Wiggans,<sup>1</sup> John B. Cole,<sup>1</sup> Suzanne M. Hubbard,<sup>1</sup> and Tad S. Sonstegard<sup>2</sup>**

<sup>1</sup>Animal Genomics and Improvement Laboratory, Agricultural Research Service, US Department of Agriculture, Beltsville, Maryland 20705-2350; email: george.wiggans@ars.usda.gov, john.cole@ars.usda.gov, suzanne.hubbard@ars.usda.gov

<sup>2</sup>Acceligen of Recombinetics Inc., St. Paul, Minnesota 55104; email: tad@recombinetics.com

[Full Text HTML](#) | [Download PDF](#) | [Article Metrics](#) | [Permissions](#) | [Reprints](#) | [Download Citation](#) | [Citation Alerts](#)

\*This is a work of the U.S. Government and is not subject to copyright protection in the United States.

### Sections

ABSTRACT

KEYWORDS

INTRODUCTION

HISTORY

CURRENT US GENOMIC EVALUATION SYSTEM

EFFECT OF GENOMIC SELECTION ON THE DAIRY INDUSTRY

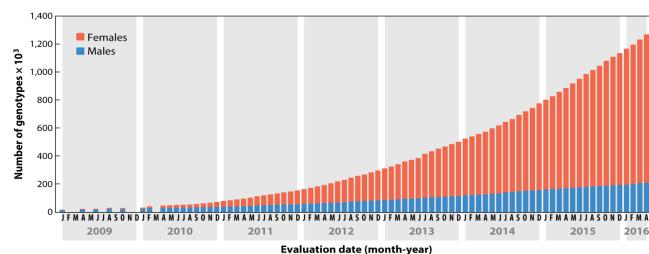
FUTURE OF GENOMIC SELECTION IN DAIRY CATTLE

CONCLUSION

### Abstract

Genomic selection has revolutionized dairy cattle breeding. Since 2000, assays have been developed to genotype large numbers of single-nucleotide polymorphisms (SNPs) at relatively low cost. The first commercial SNP genotyping chip was released with a set of 54,001 SNPs in December 2007. Over 15,000 genotypes were used to determine which SNPs should be used in genomic evaluation of US dairy cattle. Official USDA genomic evaluations were first released in January 2009 for Holsteins and Jerseys, in August 2009 for Brown Swiss, in April 2013 for Ayrshires, and in April 2016 for Guernseys. Producers have accepted genomic evaluations as accurate indications of a bull's eventual daughter-based evaluation. The integration of DNA marker technology and genomics into the traditional evaluation system has doubled the rate of genetic progress for traits of economic importance, decreased generation interval, increased selection accuracy, reduced previous costs of progeny testing, and allowed identification of recessive lethals.

## Number of genotyped bulls and cows in USA



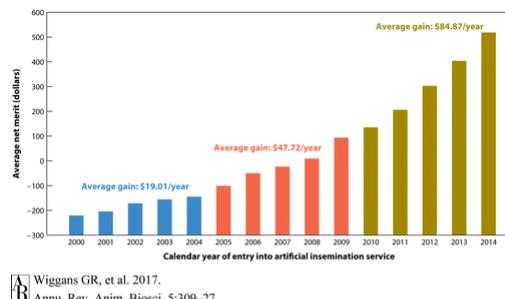
Wiggans GR, et al. 2017.  
Annu. Rev. Anim. Biosci. 5:309–27

Figure 1 Number of genotyped animals included in US genomic evaluations for dairy cattle since January 2009. Official US genomic evaluations were first released to the dairy industry in January 2009 for Holsteins and Jerseys, in August 2009 for Brown Swiss, in April 2013 for Ayrshires, and in April 2016 for Guernseys. Data for figure generation were reported by the Council on Dairy Cattle Breeding (27). Months without data represent months in which official evaluations were not released.

**Table 4** Increases in reliability of US genetic evaluations of Holsteins from including genomic information<sup>a</sup>

Trait	August 2011 reliability (%)		
	Parent average	Genomic evaluation	Gain <sup>b</sup>
Milk (kg)	38.5	72.5	34.0
Fat (kg)	38.5	72.2	33.8
Protein (kg)	38.4	63.3	24.9
Fat (%)	38.5	96.9	58.5
Protein (%)	38.4	87.4	49.0

## Evolution of genetic merit



Wiggans GR, et al. 2017.  
Annu. Rev. Anim. Biosci. 5:309–27

Figure 5 Net merit in April 2015 of marketed US Holstein bulls that entered artificial-insemination service in 2000 and later. Net merit is a genetic-economic index that was developed as a lifetime profit function that uses actual incomes and expenses for traits of economic importance (37). The economic values and traits included in the net merit index are updated as needed to reflect changes in the dairy industry, and the latest revisions were made in 2014 (38). The data for figure generation were provided by the Council on Dairy Cattle Breeding (<https://www.cdcb.us>).

## Genomic selection has profoundly affected genetic improvement of dairy cattle

- ✓ The AI organizations no longer rely on progeny-test herds to determine which bulls to market, and they purchase young bulls based on genomic evaluations.
- ✓ Another benefit of genomics is the detection of carriers of undesirable recessive characteristics.
- ✓ The age of parents of marketed bulls has steadily decreased since the start of genomic selection, essentially halving the generation interval.
- ✓ The rate of improvement in average net merit has nearly doubled for Holstein bulls since the implementation of genomic evaluation in 2010.

## How about sequence?

### Sequence (NGS) vs. Chip data

- About ~ 1000 € + 24h computing per sample for initial processing.
  - Sophisticated, tricky bioinformatic procedures.
  - ‘Unbiased’ estimates of SNP variability.
  - Missing data is unavoidable
- About ~ 100 € per sample (20x less than NGS).
  - Standard software already available.
  - Missing data can be controlled.
  - Bias is unavoidable, even in the highest density array.

## Claimed advantages of sequence

- Causal variants are in the data.
- No decrease in LD over generations.

Pérez-Enciso et al. *Genetics Selection Evolution* (2015) 47:43  
DOI 10.1186/s12711-015-0117-5



RESEARCH

Open Access

### Sequence- vs. chip-assisted genomic selection: accurate biological information is advised

Miguel Pérez-Enciso<sup>1,2,3\*</sup>, Juan C Rincón<sup>1,4</sup> and Andrés Legarra<sup>5</sup>

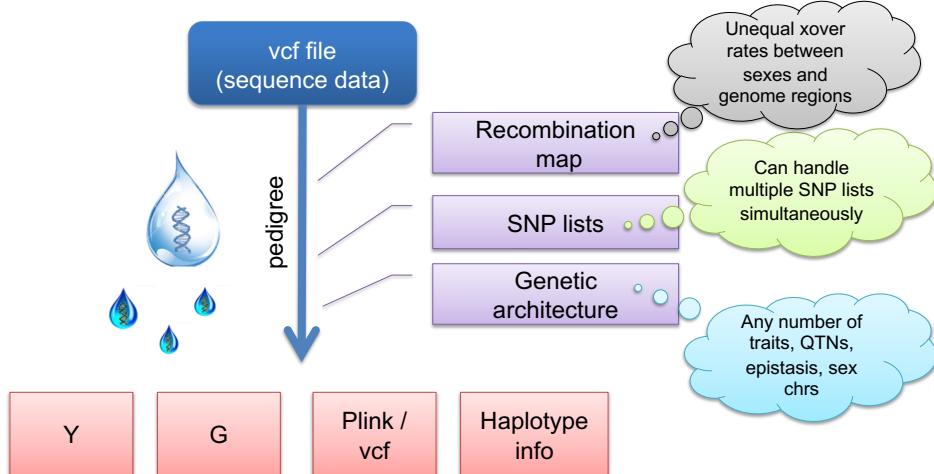
### Evaluating Sequence-Based Genomic Prediction with an Efficient New Simulator

Miguel Pérez-Enciso,<sup>\*,†,‡,§</sup> Natalia Forneris,<sup>\*</sup> Gustavo de los Campos,<sup>§,\*</sup> and Andrés Legarra<sup>†</sup>

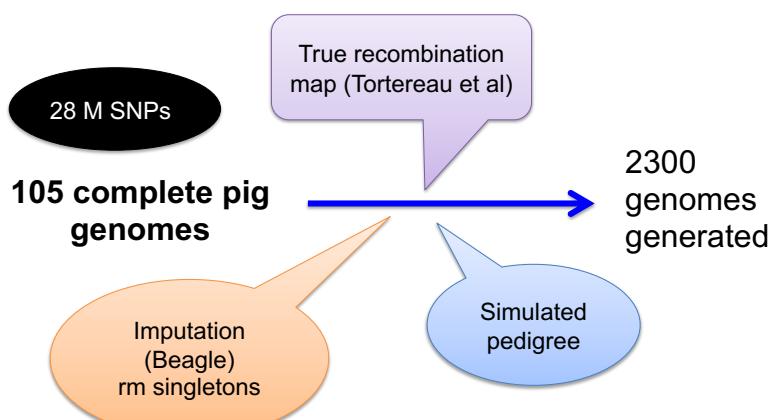
\*Centre for Research in Agricultural Genomics (CRAG), Consejo Superior de Investigaciones Científicas - Institut de Recerca i Tecnologia Agroalimentàries - Universitat Autònoma de Barcelona - Universitat de Barcelona (CSIC-IRTA-UAB-UB) Consortium and <sup>†</sup>Departament de Ciència Animal i dels Aliments, Universitat Autònoma de Barcelona, 08193 Bellaterra, Barcelona, Spain, <sup>‡</sup>Institut Català de Recerca i Estudis Avançats (ICREA), 08010 Barcelona, Spain, <sup>§</sup>Department of Epidemiology and Biostatistics and <sup>\*\*</sup>Department of Statistics, Michigan State University, East Lansing, Michigan 48824, and <sup>††</sup>Institut National de la Recherche Agronomique (INRA), Unité Mixte de Recherche 1388 GÉNPHYSE, Castanet-Tolosan 31326, France

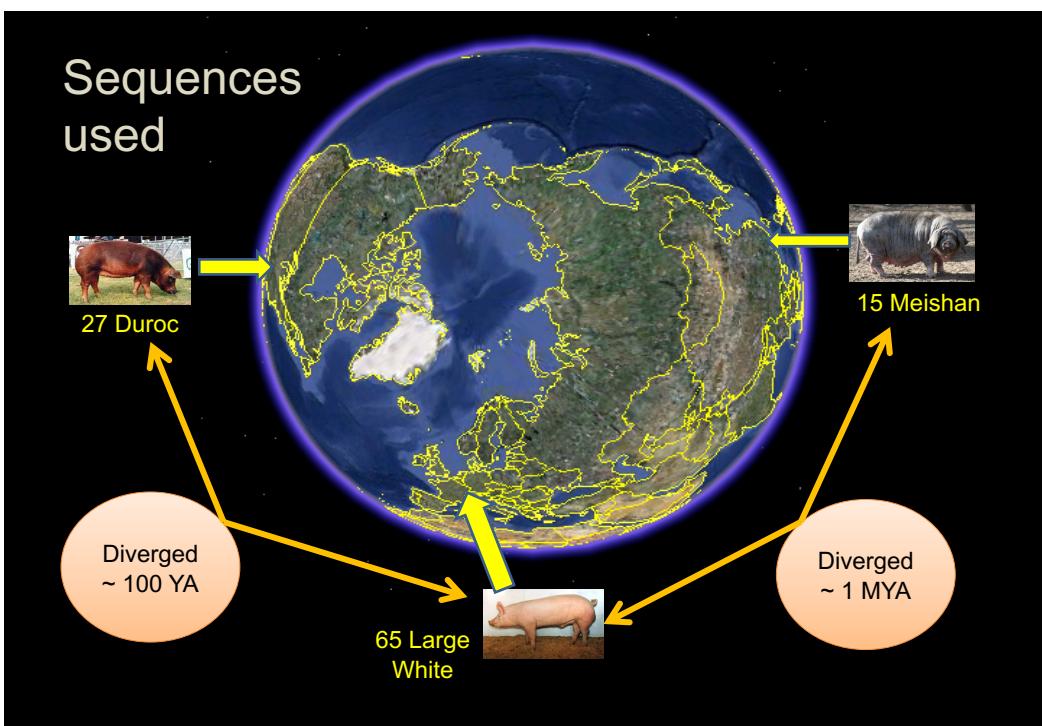
ORCID IDs: 0000-0003-3524-995X (M.P.-E.), 0000-0001-8893-7620 (A.L.)

## Simulation Strategy



## Simulation Strategy

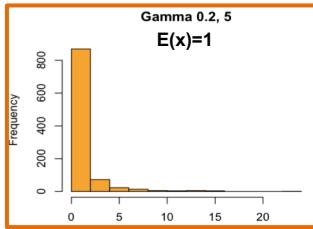




## Genetic architectures compared

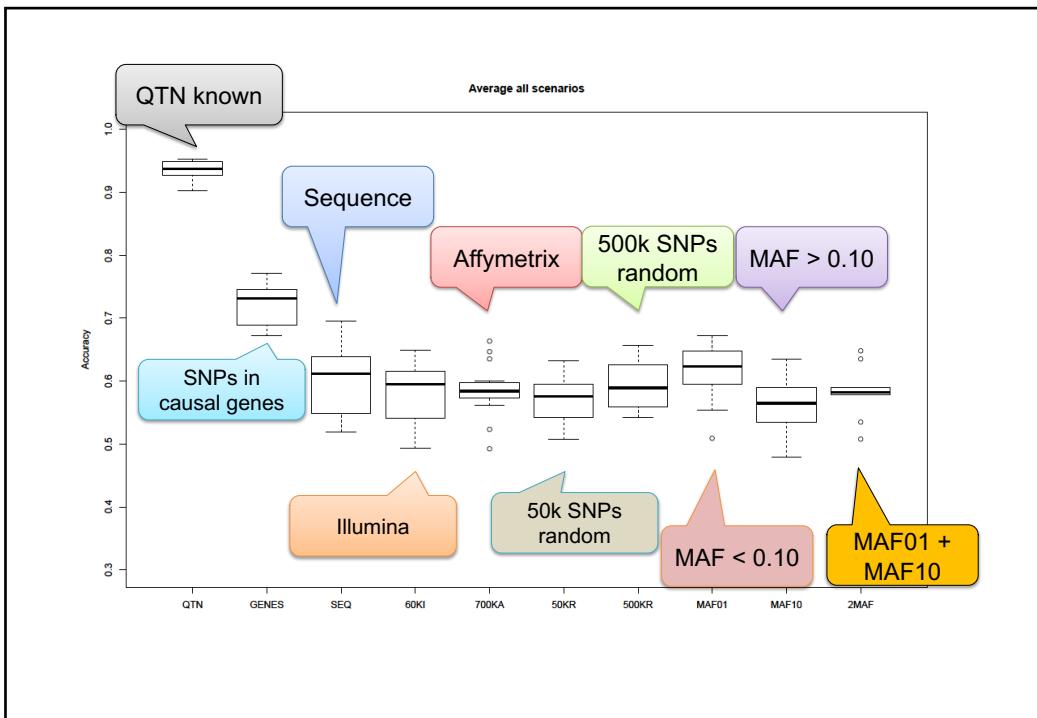
### 'NEUTRAL'

- ✓ 1000 SNPs randomly chosen as QTNs
- ✓ QTN effects  $\sim \Gamma(0.2, 5)$
- ✓ No correlation btw effect and frequency



### 'SELECTIVE'

- ✓ 200 genes with highest Fst (wild – domestic) and lowest Tajima's D
- ✓ 1000 SNPs with 'moderate' or 'high' impact + UTRs
- ✓ QTN effects  $\sim \Gamma(0.2, 5)$
- ✓ Negative correlation btw effect and frequency



## What do real data say?

Genomic prediction from whole genome sequence in livestock: the 1000 Bull Genomes Project

Conference Paper · August 2014 with 218 Reads

[Cite this publication](#)

Conference: 10th World Congress on Genetics Applied to Livestock Production

In a dairy data set, predictions using BayesRC and imputed sequence data from 1000 Bull Genomes were 2% more accurate than with 800k data.

1st Ben J Hayes

2nd Iona Macleod  
26.73 · University of Melbourne

3rd Hans D. Daetwyler  
33.74 · Department of Economic...

21

Last Michael E. Goddard

## What do real data say?

### Whole-genome sequence-based genomic prediction in laying chickens with different genomic relationship matrices to account for genetic architecture

Guilan Ni<sup>1\*</sup>, David Caverio<sup>2</sup>, Anna Fangmann<sup>1</sup>, Malena Erbe<sup>1,3</sup> and Henner Simianer<sup>1</sup>

Efficient genomic prediction based on whole-genome sequence data using split-and-merge Bayesian variable selection

Mario P. L. Calus , Aniek C. Bouwman, Chris Schrooten and Roel F. Veerkamp

*Genetics Selection Evolution* 2016 48:49 | <https://doi.org/10.1186/s12711-016-0225-x> | © The Author(s) 2016

### Accuracy of genomic prediction using imputed whole-genome sequence data in white layers

M. Heidaribar , M.P.L. Calus, H.-J. Megens, A. Vereijken,  
M.A.M. Groenen, J.W.M. Bastiaansen



Our results show that little or no benefit was gained when using all imputed WGS data to perform genomic prediction compared to using HD array data regardless of the weighting factors tested. However, using only generic SNPs from WGS data had a positive effect on prediction ability.

Predictions computed as the average of the predictions computed for each subset achieved the highest accuracies, i.e. 0.5 to 1.1 % higher than the accuracies obtained with the 50k-SNP chip, and yielded the least biased predictions.



With sequence data, there was a small increase (~1%) in prediction accuracy over the 60 K genotypes.

## So far

- ✓ Our simulations show that increasing SNP number only does not necessarily improve extant tools.
- ✓ **Conjecture:** It is by judicious use of biology that we are to make the most of sequence for prediction of genetic merit.

## ARE WE HEADING TOWARDS A BIOLOGICALLY INFORMED BREEDING?

- ❖ An advantage of biology is that (some) information can be transferred across species.
- ❖ But recall highly heterogeneous sources: reactome, omim,...
- ❖ Multiple options: machine learning or parametric methods that combine multiple sources of information.
- ❖ So far results are mixed.

## HOW MANY FUNCTIONAL KINDS?



National Center for Biotechnology Information

- [NCBI Home](#)
- [Resource List \(A-Z\)](#)
- [All Resources](#)
- [Chemicals & Bioassays](#)
- [Data & Software](#)
- [DNA & RNA](#)
- [Domains & Structures](#)
- [Genes & Expression](#)
- [Genetics & Medicine](#)
- [Genomes & Maps](#)
- [Homology](#)
- [Literature](#)
- [Proteins](#)
- [Sequence Analysis](#)
- [Taxonomy](#)
- [Training & Tutorials](#)
- [Variation](#)

**OMIA - ONLINE MENDELIAN INHERITANCE IN ANIMALS**



Gene Ontology Consortium



UniProt



OMIM  
STRUCTURAL CLASSIFICATION OF PROTEINS

**Epigenomics**



Animal QTLdb

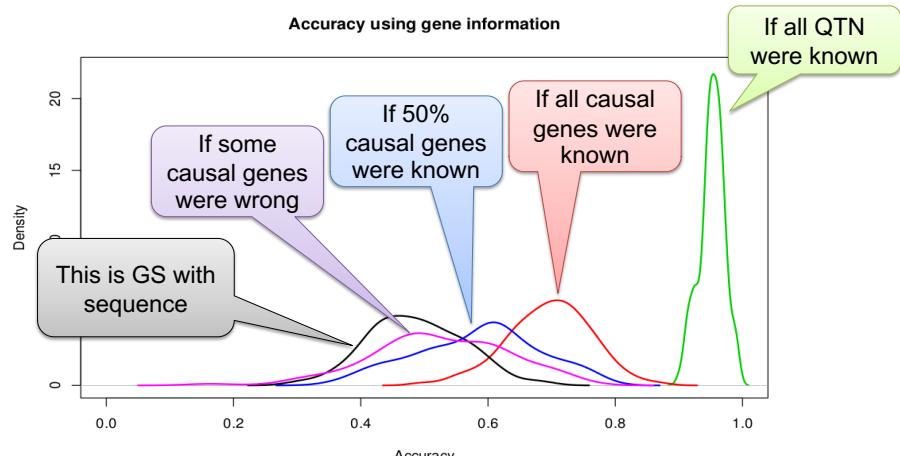


REACTOME  
A CURATED PATHWAY DATABASE

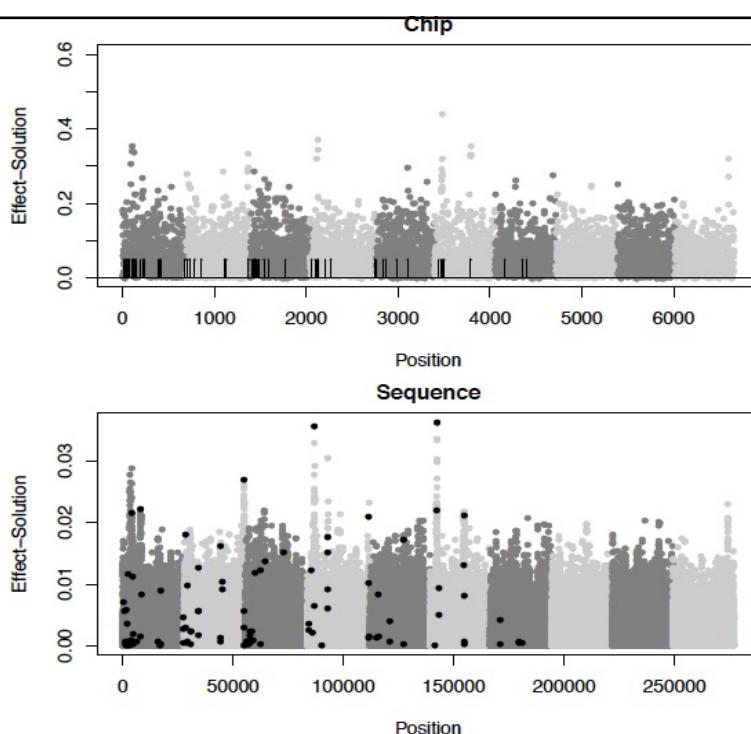


AnimalTFDB  
Animal Transcription Factor Database

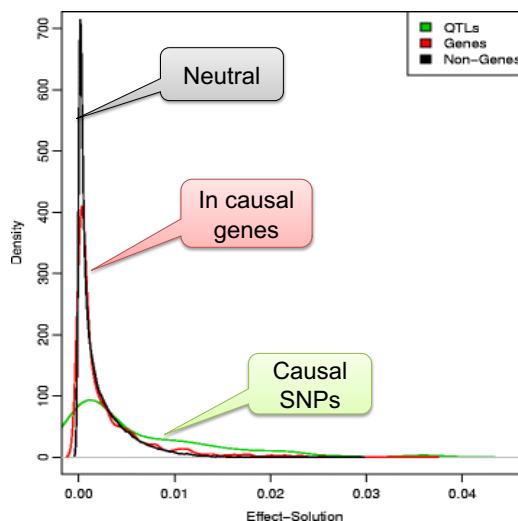
## Biology can make a difference... if true



Pérez-Enciso et al (2015), Genet Sel Evol



## Distribution of SNP estimated effects according to annotation



## In summary

- ✓ There are numerous GP algorithms, but performance tend to be similar across them.
- ✓ In general, GP methods are robust to extreme genetic architectures of complex traits.
- ✓ There is a quick law of diminishing returns with increasing SNP density. Sequence and array data are operationally equivalent for genomic selection purposes.
- ✓ Unless accurate biological information is provided, it is unlikely an increase in accuracy over 4% with sequence (in agreement with experimental results).