



Effect of the size of the reference population on the validation reliability of national genomic evaluations



Esa Mäntysaari¹

Mohammad Nilforooshan²

¹MTT, Biotechnology and Food Research

²Interbull Center

Introduction

- Previous talk by Mohammad:

A review of the validation of national genomic evaluations

- Interbull GEBV validation test since 2010
Tests if the national GEBV are unbiased – useful for GMACE and
- GEBV validation test includes also requirement $R^2_{\text{GEBV}} > R^2_{\text{EBV-PA}}$
- In 2013/14: 74 breed/country/trait tests:
5 **FAILED** because of $R^2_{\text{GEBV}} > R^2_{\text{EBV-PA}}$



Introduction

- In talk before Mohammad:

GMACE pilot #4: Adjusting the national reliability input data. (Sullivan and Jakobsen 2014)

- What is the effect of size of reference pop to model based R^2_{GEBV}

$$\exp(Grel_n) = trait + \sum rel_loc + \sum rel_for$$

- Differences on R^2_{GEBV} values submitted to ITB and predicted by ref pop size: -5.7 - +7.25 (protein)

Results from
Interim
Reports 2013

- Suggestion: For the stability of GMACE the country submitted should be scaled towards the predicted



Introduction

Goal of this presentation:

- Relate the **validation** R^2_{GEBV} with the size of reference population !
- Interests:
 1. Value of domestic and foreign MACE information
 2. Behavior of R^2_{GEBV} different traits
 3. Behavior of R^2_{GEBV} different breeds
 4. Behavior of R^2_{GEBV} different evaluation models



Accuracy of Genomic evaluation

- Several equations exist for predicting the accuracy of DGV
 - Daetwyler et al, 2008; Goddard, 2009; Hayes et al. 2009; Goddard et al. 2011; Meuwissen et al. 2013)
 - Generally reliability of prediction for the animals that have no phenotypes **themselves:**

$$R_{DGV}^2 = w \frac{N_{ref} * h^2}{N_{ref} * h^2 + Me}$$

where

- w is the proportion of genetic variance that can be predicted by genomic model
- N_{ref} is the number of animals with genotypes and phenotypes
- h^2 is the prediction accuracy of the phenotypes
- Me is the number of haplotypes segregating in the population



Accuracy of Genomic evaluation

- Several equations exist for predicting the accuracy of DGV
 - Daetwyler et al, 2008; Goddard, 2009; Hayes et al. 2009; Goddard et al. 2011; Meuwissen et al. 2013)
 - Generally reliability of prediction for the animals that have no phenotypes themselves:

$$R_{DGV}^2 = w \frac{N_{ref} * h^2}{N_{ref} * h^2 + Me}$$

where

- w is the proportion of genetic variance that can be predicted by genomic model
- N_{ref} is the number of animals with genotypes and phenotypes
- h^2 is *Accuracy of observation: heritability or reliability*
- Me is the number of haplotypes segregating in the population



Accuracy of Genomic evaluation (2)

- The prediction generally fits poorly to our data
- Erbe et al. (2013)
A Function Accounting for Training Set Size and Marker Density to Model the Average Accuracy of Genomic Prediction
Used ML estimation to obtain R^2 prediction model parameters
- We reparametrized the base model to the simplest form:

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Me/h^2 were estimated with non-linear model
(using function nls in R)



Accuracy of Genomic evaluation (2)

Base Model

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model I

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model II

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model III

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model IV

$$(R_{DGV}^2 - R_{EBV-PA}^2) = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model IV

$$R_{DGV}^2 = w_1 \frac{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}}}{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}} + Me/h^2}$$



Accuracy of Genomic evaluation (2)

Base Model

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model I

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model II

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model III

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model IV

$$(R_{DGV}^2 - R_{EBV-PA}^2) = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model IV

$$R_{DGV}^2 = w_1 \frac{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}}}{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}} + Me/h^2}$$



Accuracy of Genomic evaluation (2)

Base Model

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model I

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model II

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model III

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model IV

$$(R_{DGV}^2 - R_{EBV-PA}^2) = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model IV

$$R_{DGV}^2 = w_1 \frac{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}}}{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}} + Me/h^2}$$



Accuracy of Genomic evaluation (2)

Base Model

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model I

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model II

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model III

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model IV

$$(R_{DGV}^2 - R_{EBV-PA}^2) = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model IV

$$R_{DGV}^2 = w_1 \frac{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}}}{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}} + Me/h^2}$$



Accuracy of Genomic evaluation (2)

Base Model

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model I

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model II

$$R_{DGV}^2 = \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model III

$$R_{DGV}^2 = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2} + w_3 R_{EBV-PA}^2$$

Model IV

$$(R_{DGV}^2 - R_{EBV-PA}^2) = w_1 \frac{N_{ref}}{N_{ref} + Me/h^2}$$

Model IV

$$R_{DGV}^2 = w_1 \frac{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}}}{N_{ref_{Domestic}} + w_2 N_{ref_{foreign}} + Me/h^2}$$





Unfortunate Realism

Maanyy GEBV tests....

```
> table(data$brd,data$trt)/2
```

	ang	bcs	bde	cc1	cc2	crc	cwi	dce	dlo	dsb	fan	fat	ftl	ftp	fua	hco	int	loc	mas	mce	mil	msb	msp	ocs	ofl	ous	pro	ran	r1r	r1s	rtp	ruh	rwi	scs	sta	tem	ude	usu
BSW	1	0	0	1	1	1	1	2	1	1	1	4	1	1	1	1	2	0	0	2	4	1	1	1	0	0	4	1	1	1	0	1	1	3	1	1	1	1
HOL	8	5	8	10	11	11	8	8	11	6	8	14	8	8	8	6	8	6	3	7	14	6	6	8	7	6	17	8	8	8	5	8	7	14	14	5	9	14
JER	2	0	1	1	1	1	2	1	2	1	2	3	2	2	2	1	2	1	0	1	3	1	1	2	1	1	5	2	1	2	1	2	2	2	2	1	2	2
NOR	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
RDC	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
SIM	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	2	0	0	0	0	0	3	0	0	0	0	0	0	0	0	0	0	0
MON	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	
RHO	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	

But

- Only limited number by trait (at maximum 17 per breed)
- Only few on breeds other than Holstein
- And the domestic vs. foreign information was too weak to be used

Therefore

prediction models were fitted:

- 8 traits: milk, protein, fat, fertility (cc1), SCS, longevity, direct calving ease, stature
- Holstein only



Summary of model Fits:

Residual Mean Squares of Models

Trait	0	I	II	III	IV
Milk	15,96	9,81	18,23	8,85	8,00
Fat	18,16	12,60	6,19	5,95	5,50
Protein	17,71	11,78	6,49	7,23	6,59
SCS	11,28	7,17	14,09	7,34	6,91
Fertility	18,24	17,46	4,96	5,49	***
Direct Calving Ease	25,08	24,29	16,24	10,74	10,13
Direct Longevity	10,91	9,58	7,34	6,46	6,07
Stature	17,91	8,13	20,37	4,49	4,06



Model w. R^2_{EBV-pa} as a covariable
and w as maximum reliability

General view of R^2_{GEBV} and $R^2_{\text{EBV-PA}}$ Holstein



Protein

Mean

R^2_{GEBV} 41%

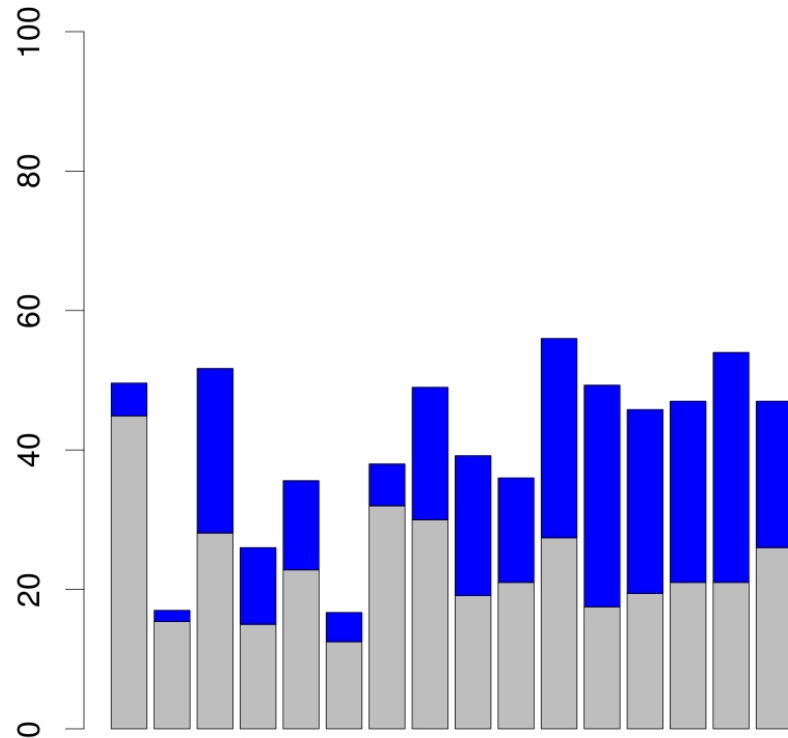
Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 18%

Ordered by size of Nref

Clear indication of increasing

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$



Protein GEBV validation reliability
Grey is EBV-PA validation reliability

General view of R^2_{GEBV} and $R^2_{\text{EBV-PA}}$ Holstein



Milk

Mean

R^2_{GEBV} 46%

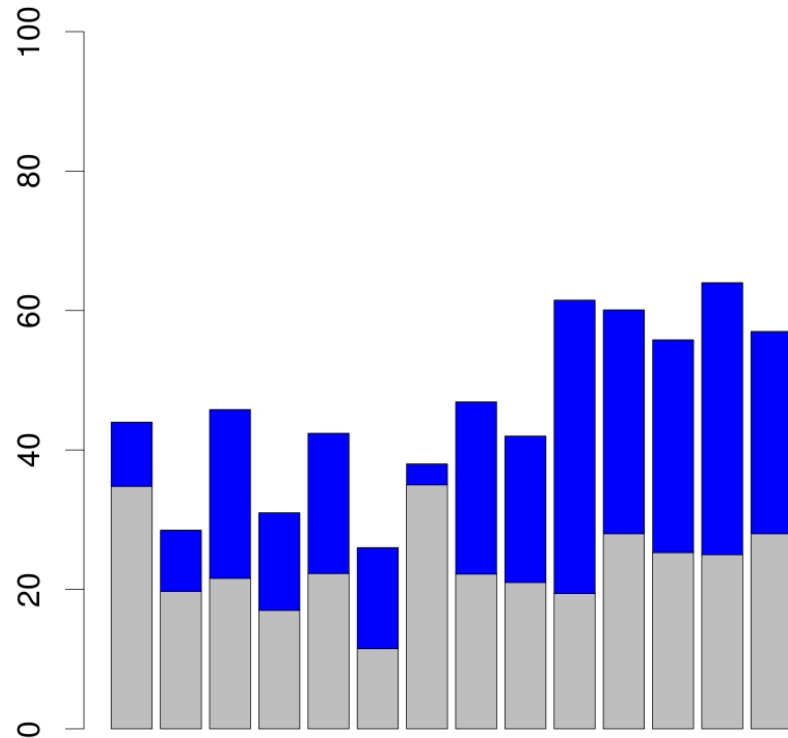
Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 22%

Ordered by size of Nref

Clear indication of increasing

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$



Milk GEBV validation reliability
Grey is EBV-PA validation reliability

General view of R^2_{GEBV} and $R^2_{\text{EBV-PA}}$ Holstein



Fat

Mean

R^2_{GEBV} 44%

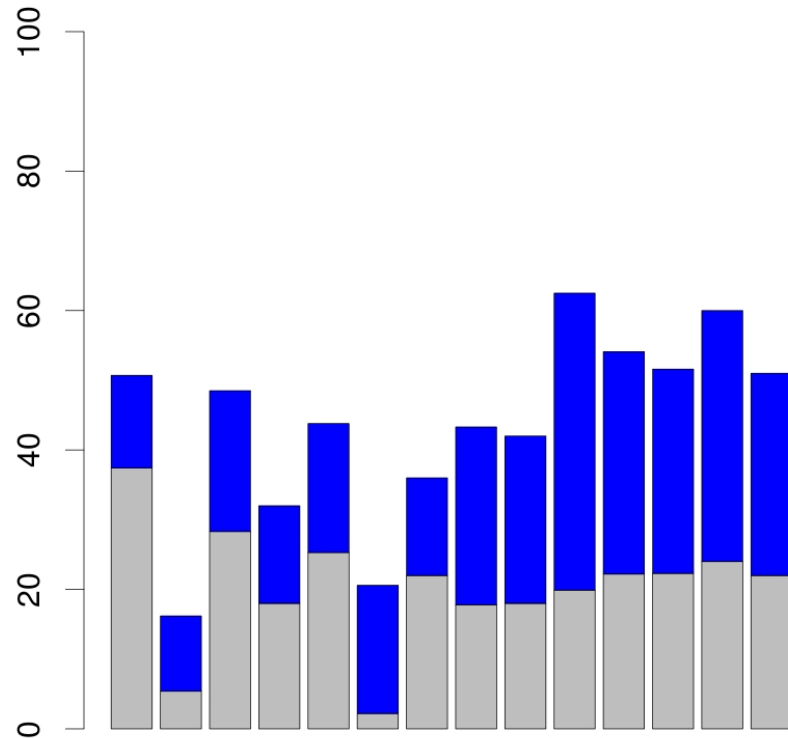
Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 23%

Ordered by size of Nref

Clear indication of increasing

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$



Fat GEBV validation reliability
Grey is EBV-PA validation reliability

General view of R^2_{GEBV} and $R^2_{\text{EBV-PA}}$ Holstein



SCS

Mean

R^2_{GEBV} 42%

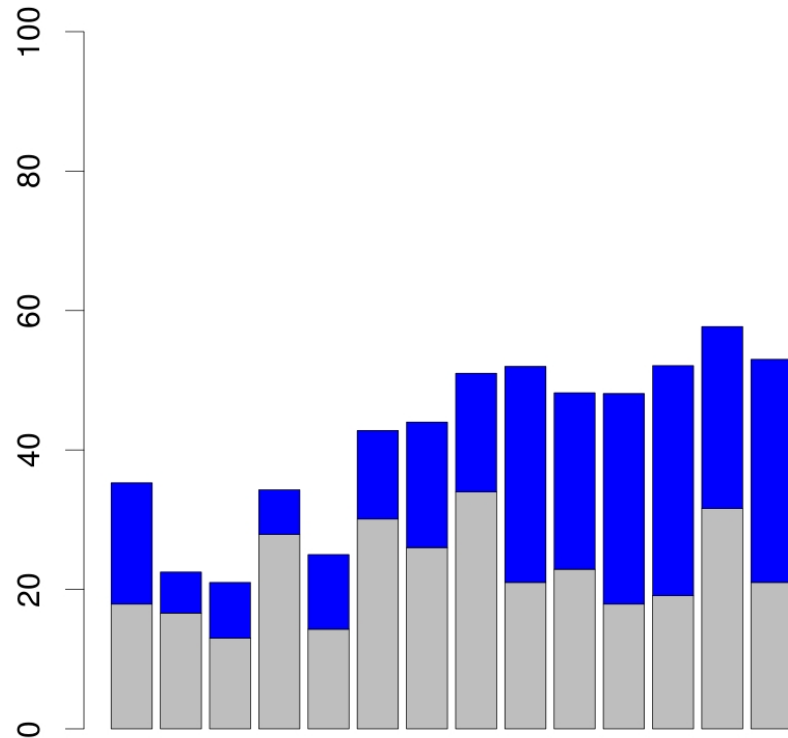
Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 20%

Ordered by size of Nref

Clear indication of increasing

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$



SCS GEBV validation reliability

Grey is EBV-PA validation reliability

General view of R^2_{GEBV} and $R^2_{\text{EBV-PA}}$ Holstein



Stature

Mean

R^2_{GEBV} 52%

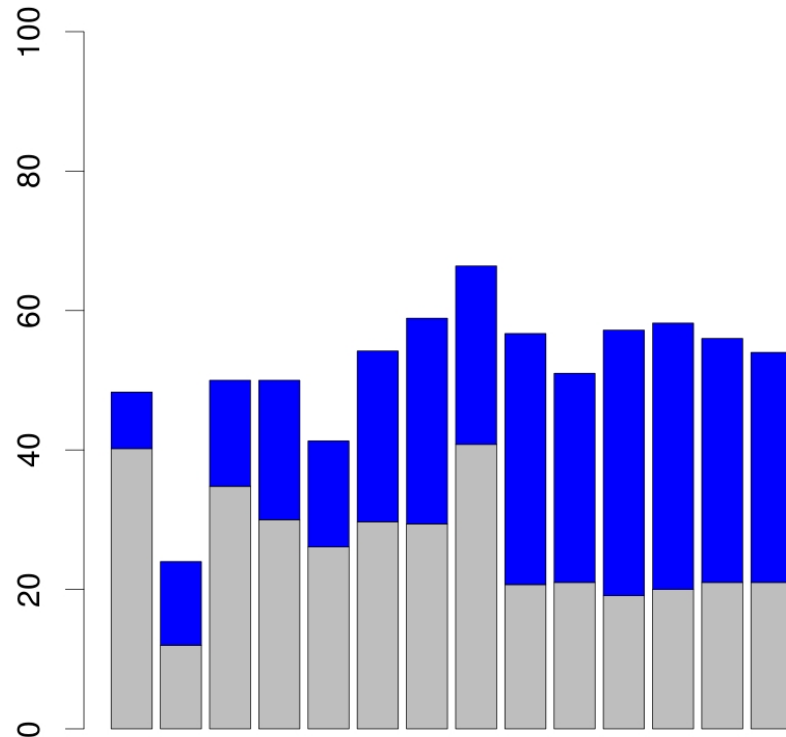
Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 26%

Ordered by size of Nref

Clear indication of increasing

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$



Stature GEBV validation reliability

Grey is EBV-PA validation reliability

General view of R^2_{GEBV} and $R^2_{\text{EBV-PA}}$ Holstein



Fertility

Mean

R^2_{GEBV} 19%

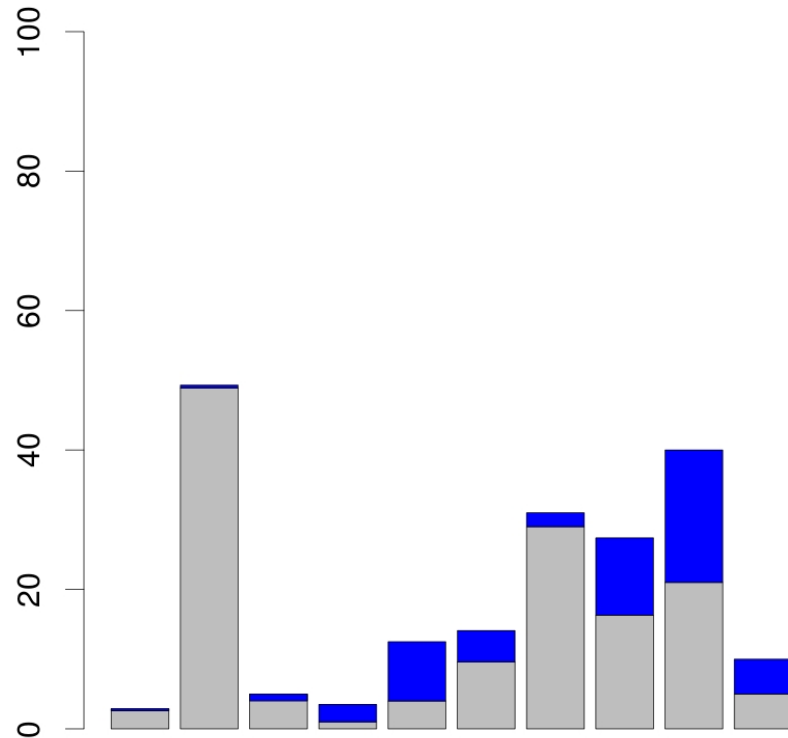
Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 5%

Ordered by size of Nref

Moderate indication of increasing

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$



Fertility GEBV validation reliability

Grey is EBV-PA validation reliability

General view of R^2_{GEBV} and $R^2_{\text{EBV-PA}}$ Holstein



Longevity

Mean

R^2_{GEBV} 21%

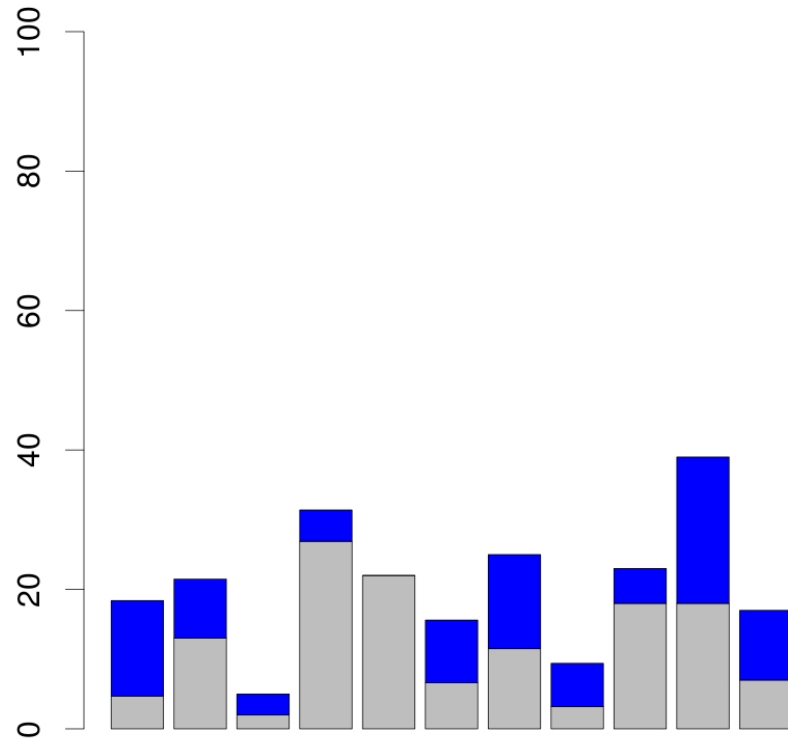
Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 9%

Ordered by size of Nref

Weak indication of increasing

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$



Longevity GEBV validation reliability

Grey is EBV-PA validation reliability

General view of R^2_{GEBV} and $R^2_{\text{EBV-PA}}$ Holstein



Calving ease

Mean

R^2_{GEBV} 41%

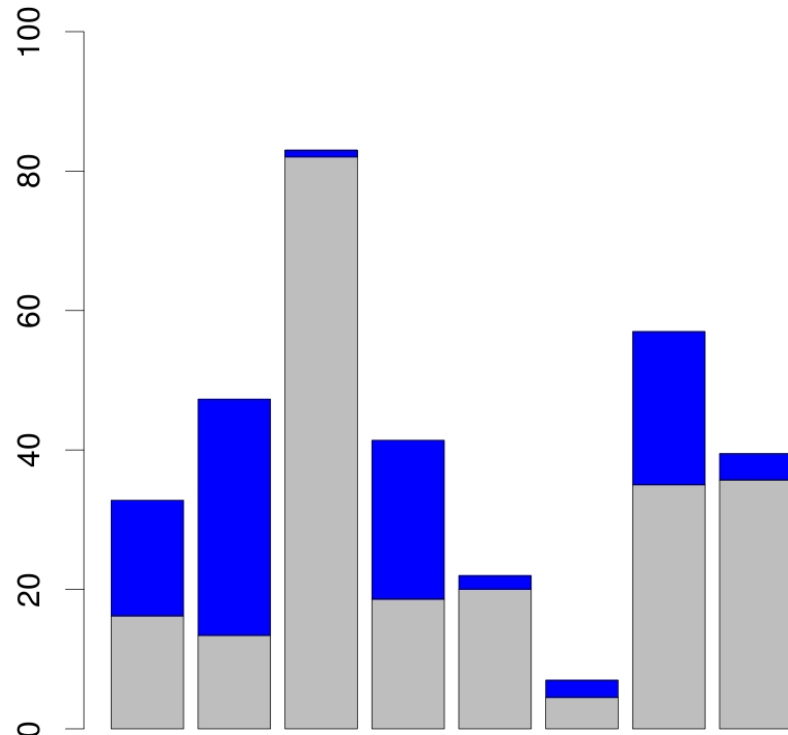
Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 13%

Ordered by size of Nref

No indication of increasing R^2_{GEBV} -

$R^2_{\text{EBV-PA}}$



dlo GEBV validation reliability
Grey is EBV-PA validation reliability

JERSEY

Fat

Mean

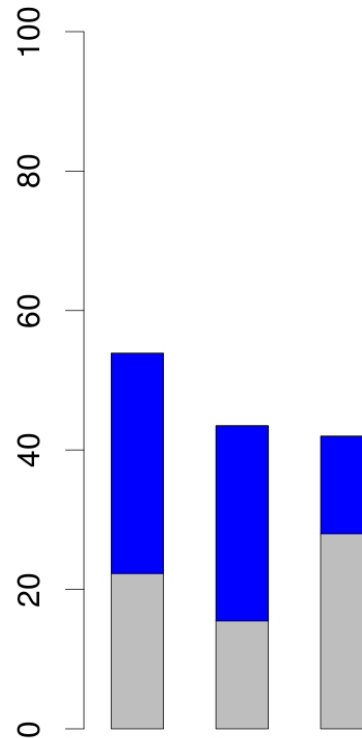
R^2_{GEBV}

46%

Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$

25%



Fat GEBV validation reliability
Grey is EBV-PA validation reliability

JERSEY

Protein

Mean

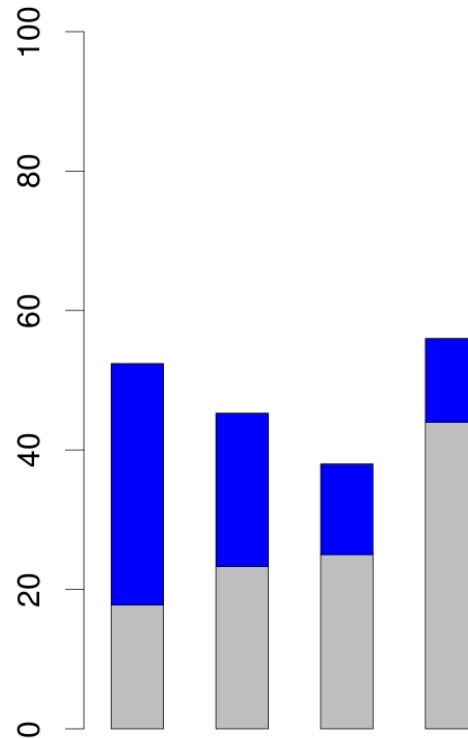
R^2_{GEBV}

48%

Average increase

$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$

20%



Protein GEBV validation reliability

Grey is EBV-PA validation reliability

JERSEY

Milk

Mean

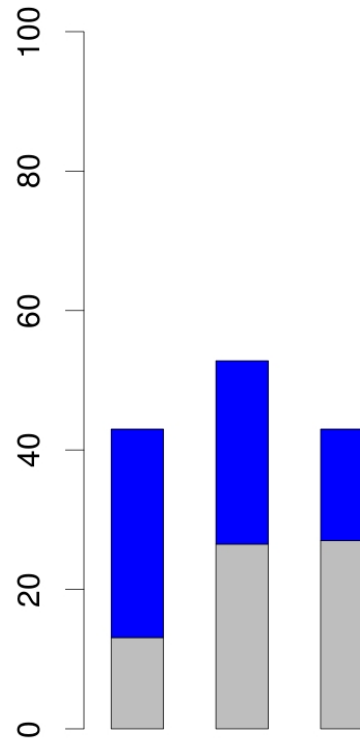
R^2_{GEBV}

46%

Average increase

$R^2_{GEBV} - R^2_{EBV-PA}$

24%



Milk GEBV validation reliability
 Grey is EBV-PA validation reliability

JERSEY

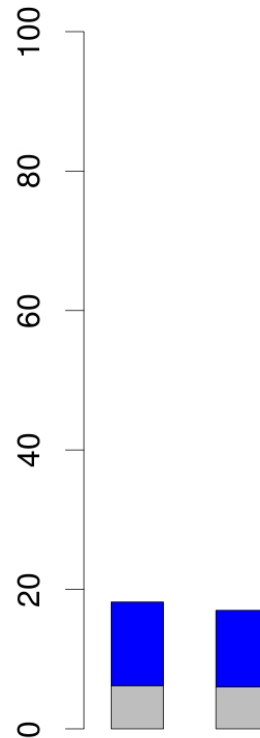
SCS

Mean

R^2_{GEBV} 18%

Average increase

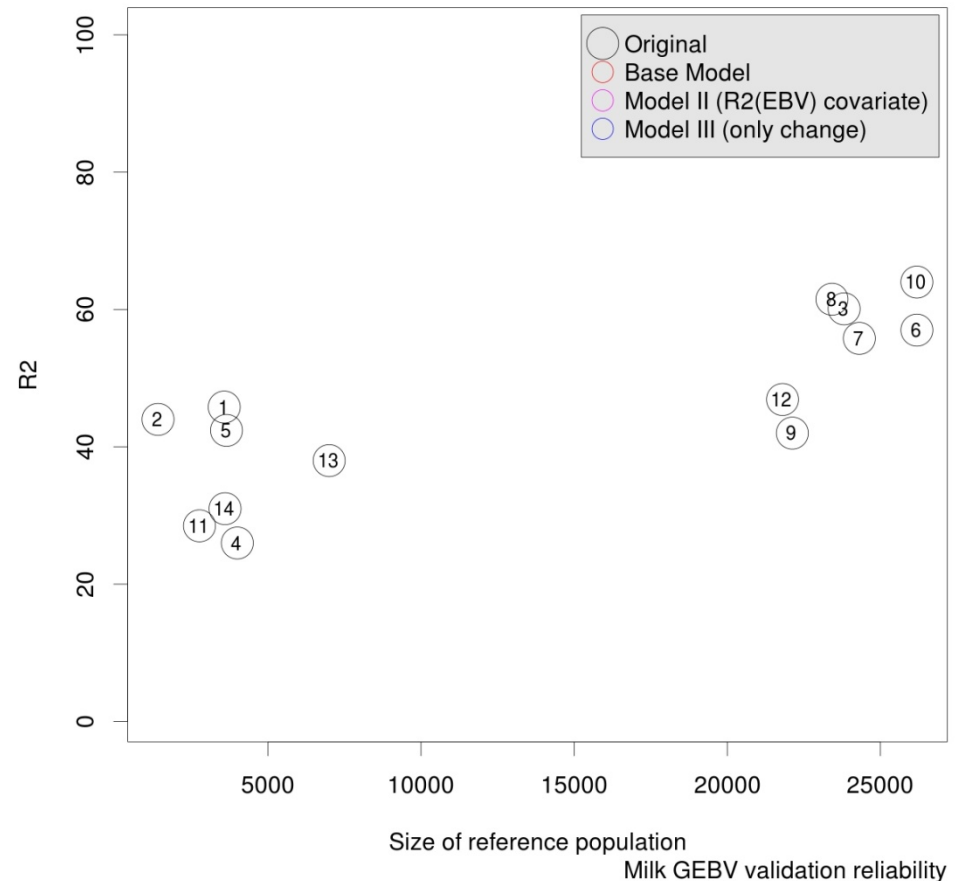
$R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ 12%



SCS GEBV validation reliability
 Grey is EBV-PA validation reliability

R^2_{GEBV} vs. reference population Holstein

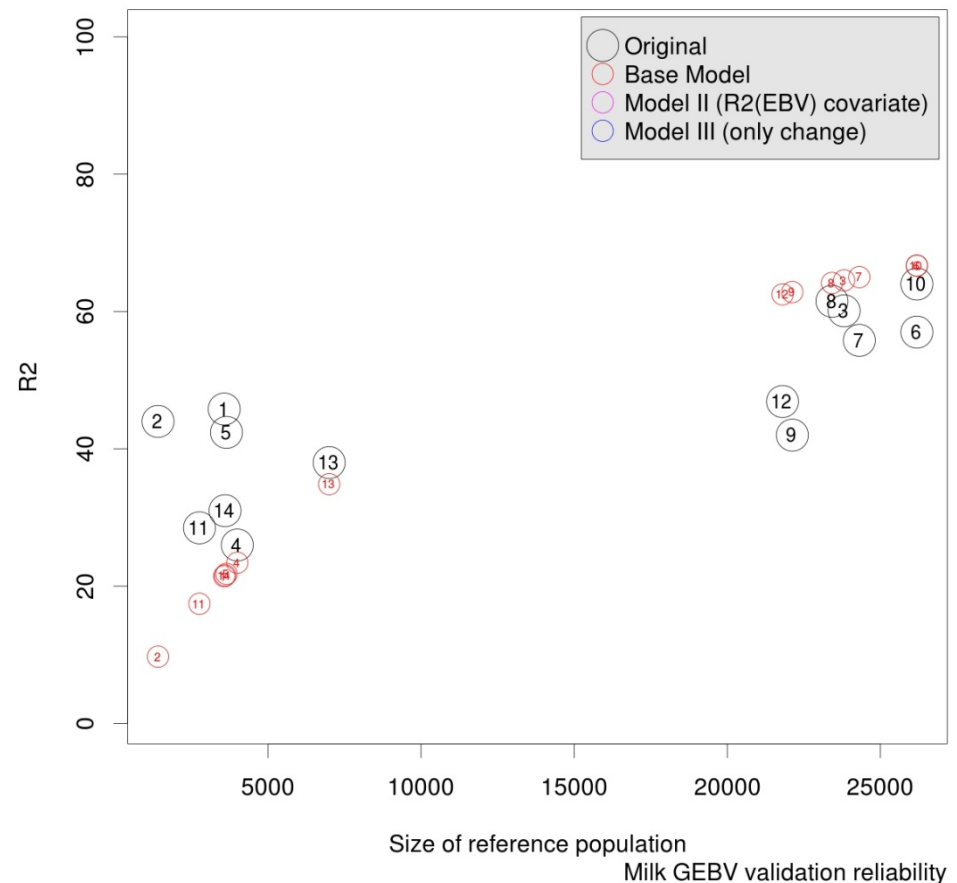
- Milk
- Clear difference between single populations and populations in alliances
 - In nref size
 - Not as clear in R^2_{GEBV}



R^2_{GEBV} vs. reference population Holstein

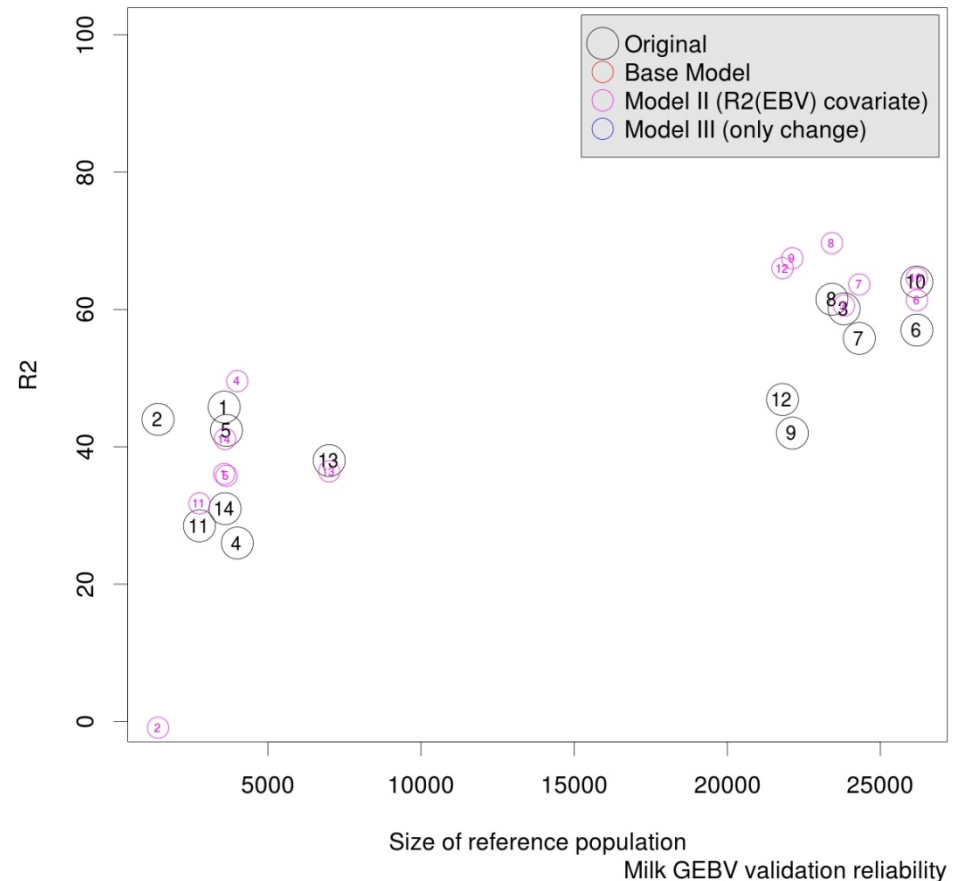


- Milk
- Base model fitted
- Underestimation of R^2 in small pop and overestimation in large pop



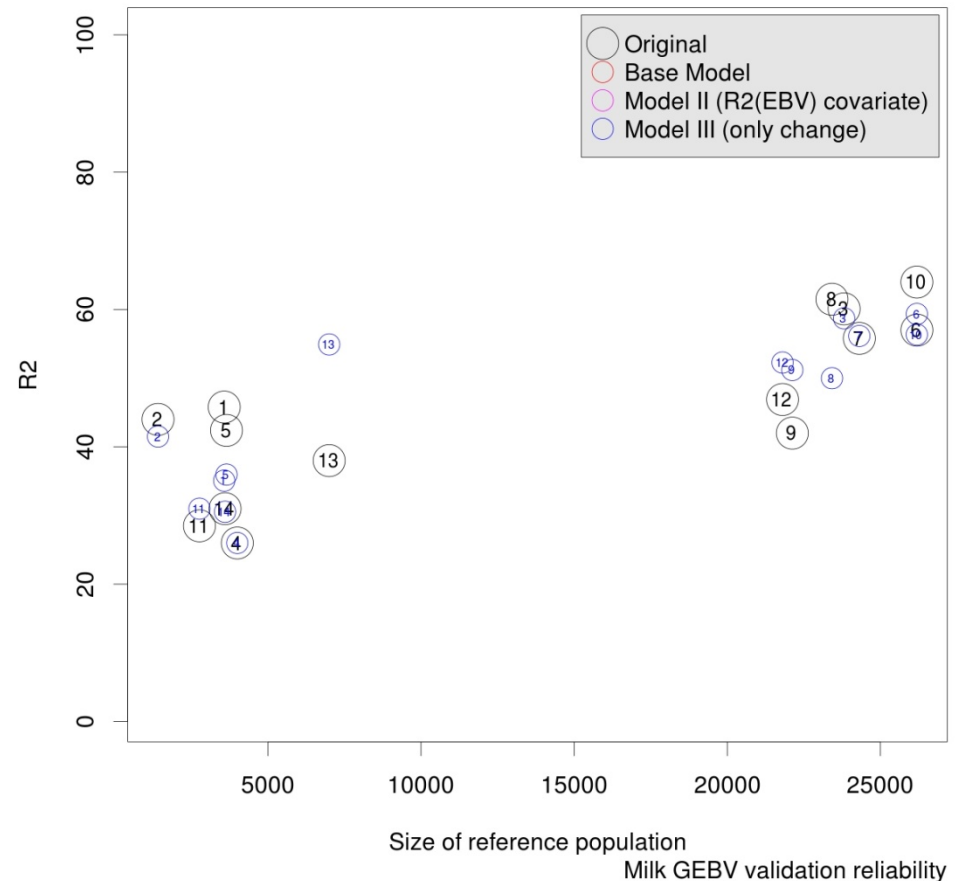
R^2_{GEBV} vs. reference population Holstein

- Milk
- When R^2 is predicted with a model that has R^2_{EBV-PA} as covariate:
 - no underprediction
 - less overprediction



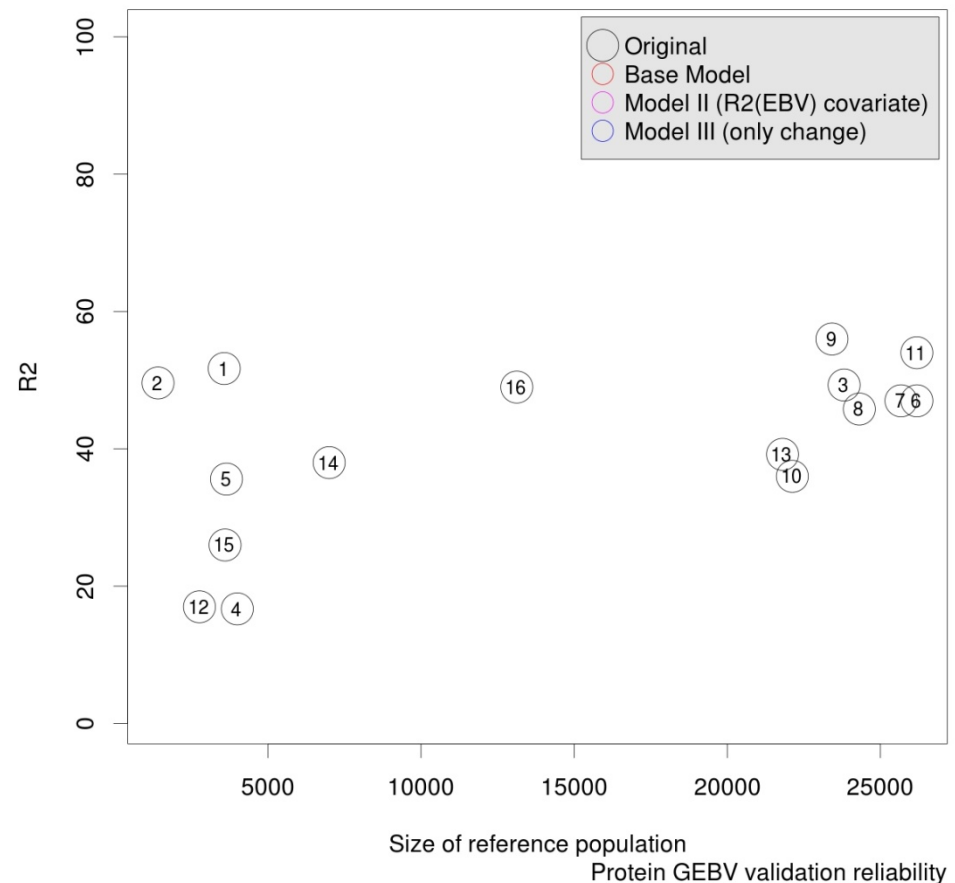
R^2_{GEBV} vs. reference population Holstein

- Milk
- When $R^2_{GEBV} - R^2_{EBV-PA}$ is predicted
 - no underprediction
 - no clear overprediction



R^2_{GEBV} vs. reference population Holstein

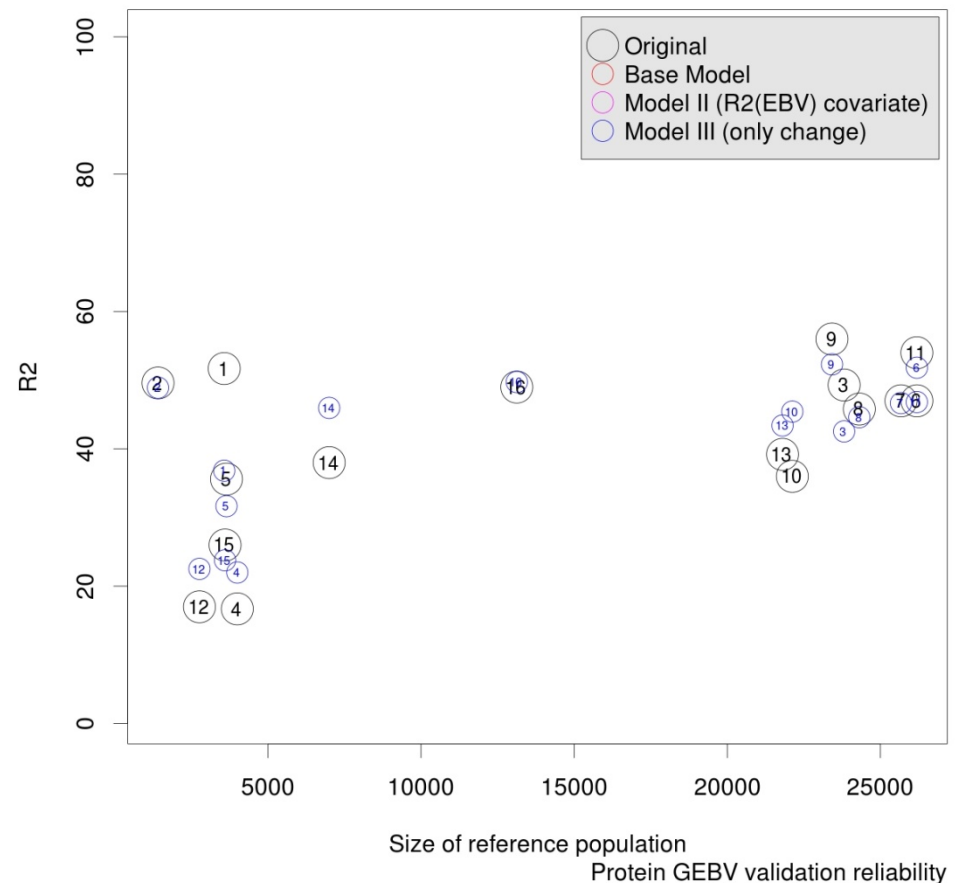
- Protein
- More variability than in milk
 - especially in small pop
- In large pop values are lower than in milk



R^2_{GEBV} vs. reference population Holstein

- Protein
- When $R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ is predicted
 - no underprediction
 - no clear overprediction

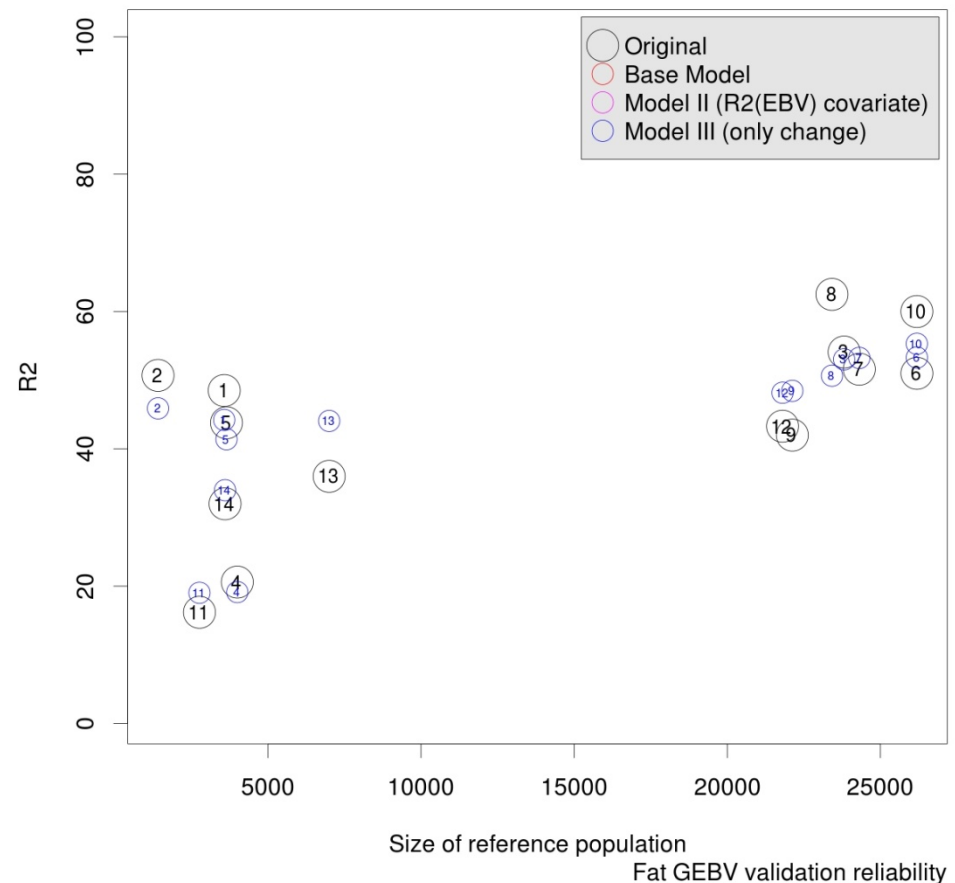
VERY GOOD FIT



R^2_{GEBV} vs. reference population Holstein

- Fat
- Again more variability than in milk
 - especially in small pop!
- $R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$
Shown

VERY GOOD FIT



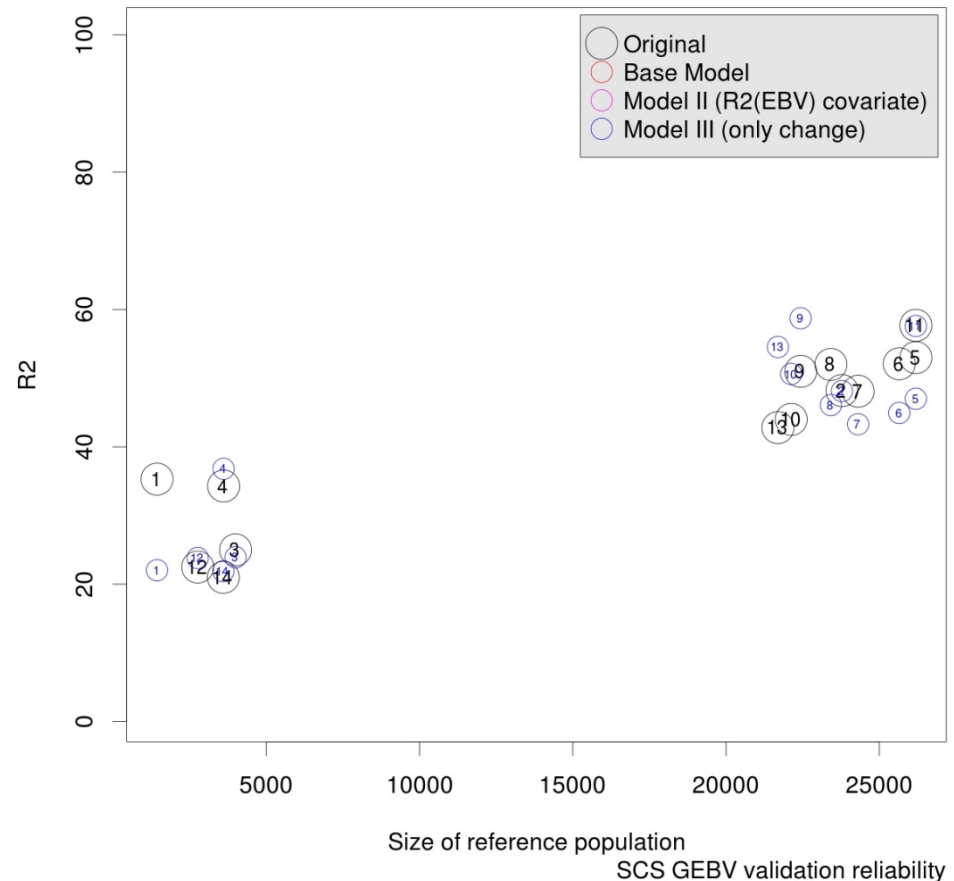
R^2_{GEBV} vs. reference population Holstein

- SCS

- Not much variability
Clear effect of Nref size

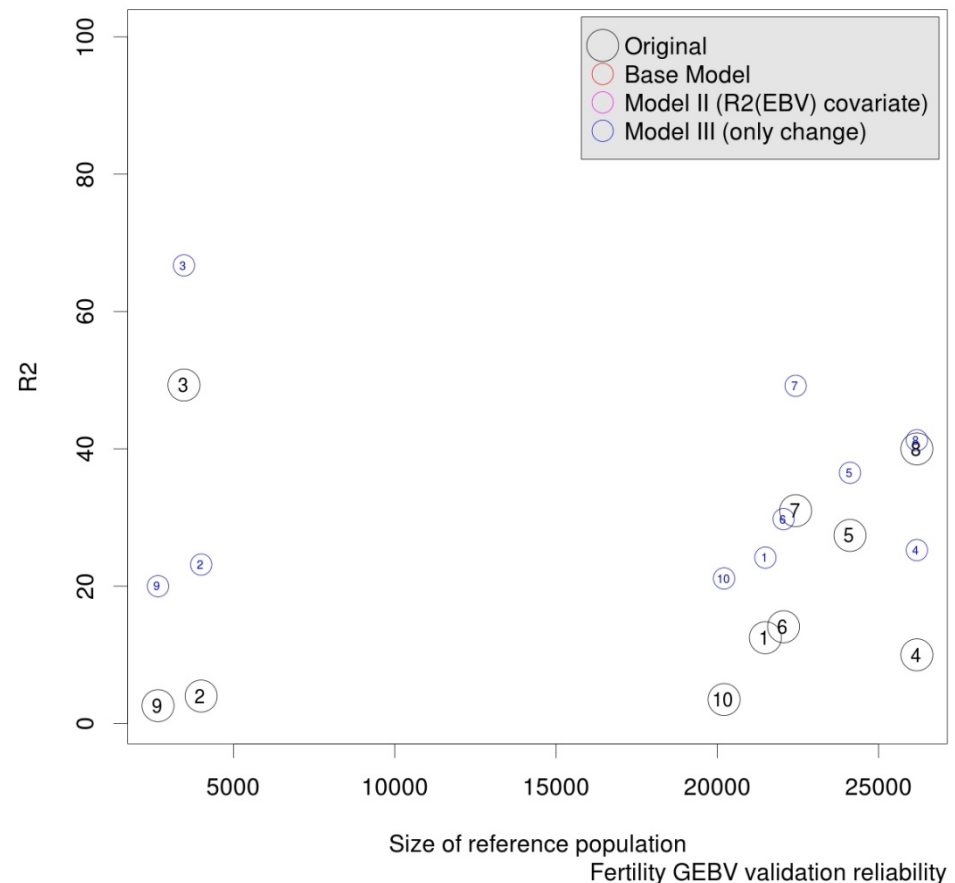
- $R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$
Shown

Reasonable GOOD FIT



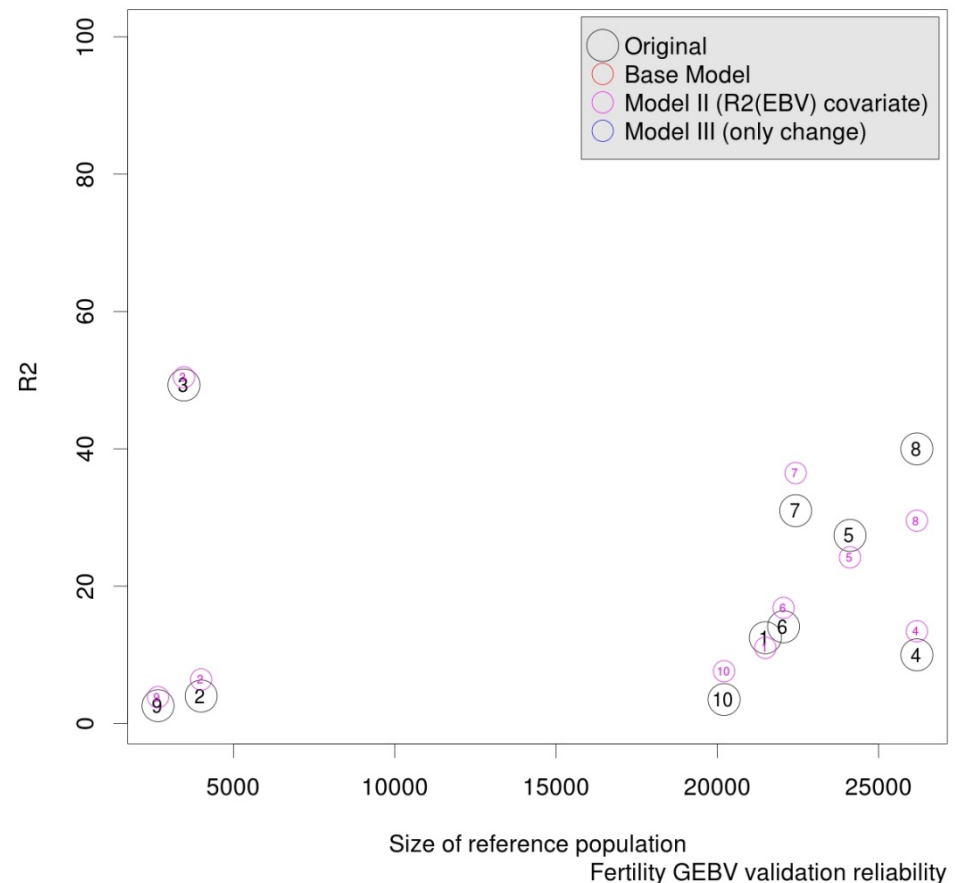
R^2_{GEBV} vs. reference population Holstein

- Fertility
- Much more variability than in production traits
 - especially in small pop
- Both Nref groups have values lower than production traits
- $R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$
Fits very poorly



R^2_{GEBV} vs. reference population Holstein

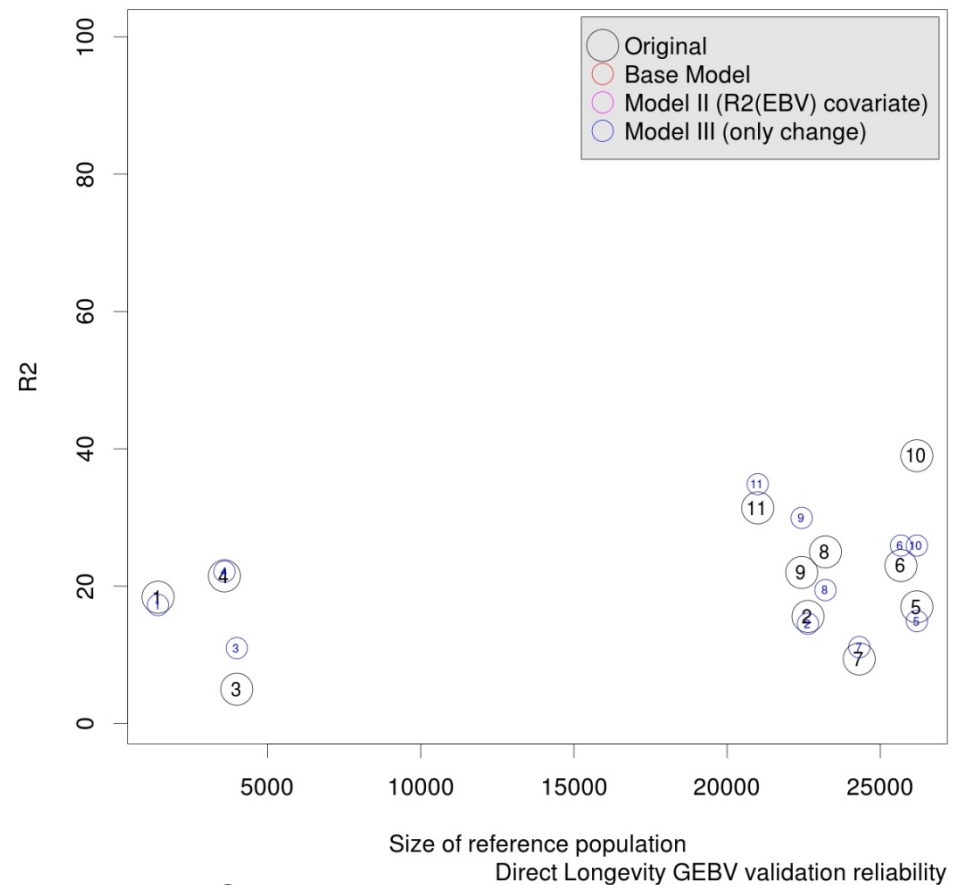
- Fertility
- Much more variability than in production traits
 - especially in small pop
- Both Nref groups have values lower than production traits
- Model III with estimate of covariable for $R^2_{\text{EBV-PA}}$ is much better



R^2_{GEBV} vs. reference population Holstein

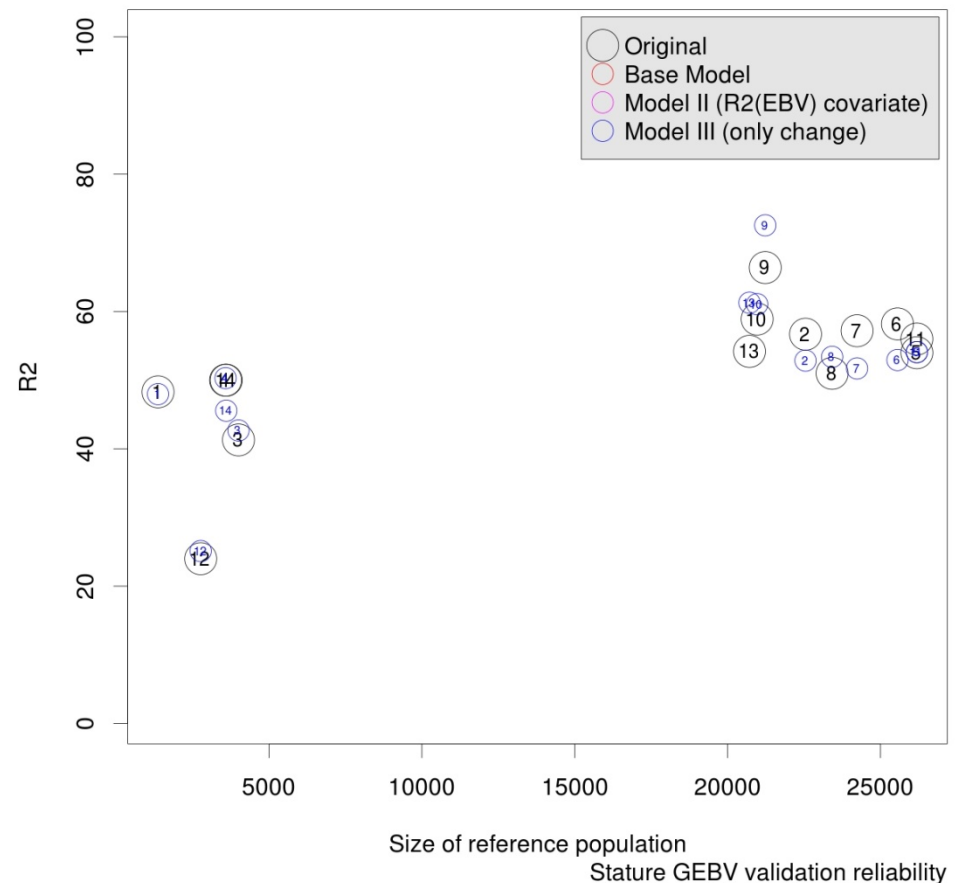


- Longevity
- Values of R^2 are low to very low
- Fit for $R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ is quite nice



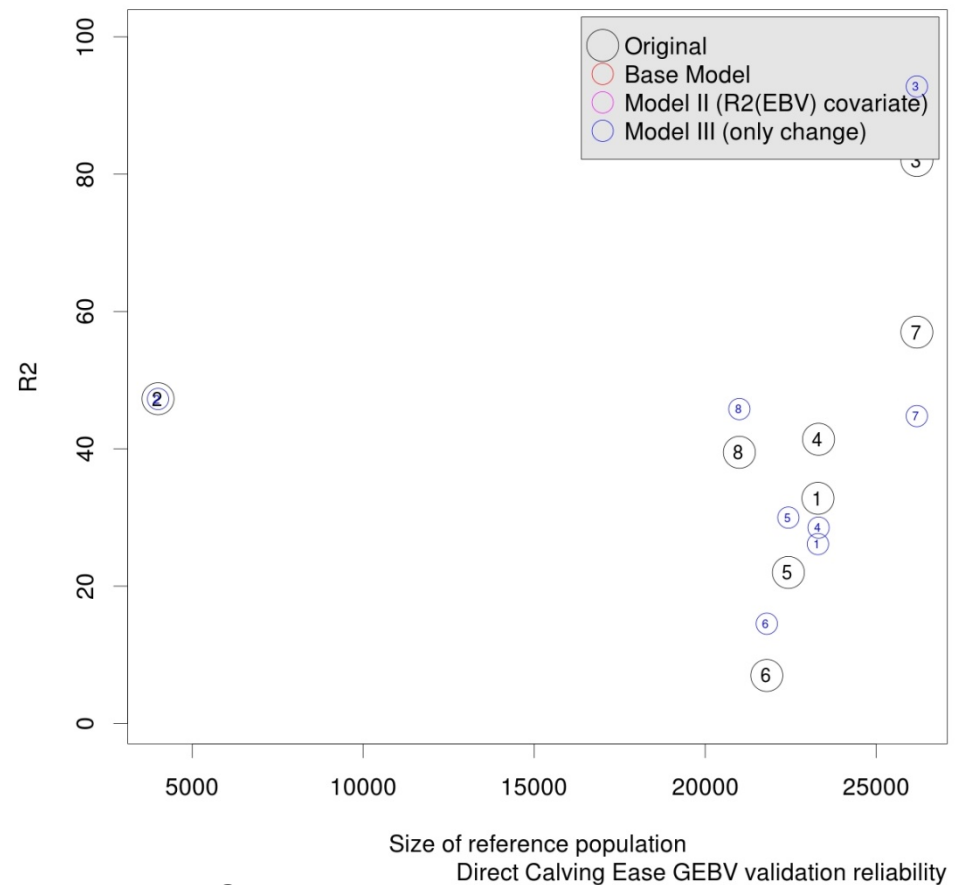
R^2_{GEBV} vs. reference population Holstein

- Stature
- Values of R^2 are pretty much in same level as w. production
 - Not excessive variability either
- Fit for $R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$ is quite nice



R^2_{GEBV} vs. reference population Holstein

- Calving Ease
- Values of R^2 are low and very variable
 - Population 6 has a R^2 of 6%
- Fit for $R^2_{\text{GEBV}} - R^2_{\text{EBV-PA}}$:
 - Fits well to point 2 in low nref
 - For the large Nref the covariate model is maybe better



Conclusions

- GEBV R^2 data does not fit directly to theoretical accuracy prediction model
 - Large variation noise by populations
 - Maybe different models (also in validation bull data)
 - This can be somewhat modeled via R^2_{EBV-PA}
- Clearly lower R^2 with low heritability traits
 - Also more variable
 - ==> Genomic evaluation can be used to improve fertility
- Would be reasonable to require more just non-zero genomic gain.

Maybe $\Delta 20\%$ i.e. $R^2_{GEBV} > 1.2 * R^2_{EBV-PA}$

THANK YOU