



UNIVERSITY of MARYLAND  
SCHOOL OF MEDICINE

# A scalable Bayesian mixed model approach for GWAS and genomic prediction

**Jicai Jiang**<sup>1</sup>, Li Ma<sup>2</sup>, Paul M. VanRaden<sup>3</sup>, Jeffrey R. O'Connell<sup>1</sup>

<sup>1</sup>University of Maryland School of Medicine, Baltimore, MD

<sup>2</sup>University of Maryland, College Park, MD

<sup>3</sup>AGIL-ARS-USDA, Beltsville, MD

2019 Interbull Meeting, Cincinnati, Ohio

June 22, 2019



# Introduction

# Big-data challenges



**Biobank-scale samples**

**Big  $N > 0.5$  million**



**Sequence variants**

**Big  $M > 10$  million**

# Recent advances

## Variance component estimation

- BOLT-REML (MC-EM-REML)
- RHE-reg (Randomized Haseman-Elston regression)

## GWAS

- BOLT-LMM ( $\chi_{LMM}^2 = c \cdot \chi_{LM}^2$ )
- SAIGE (Like BOLT but better for binary traits)
- GCTA-fastGWA (Sparsifying GRM)

## Genomic prediction

- ssGBLUP
- BayesRv2 (One of the fastest MCMC-based)



## Issues

- $\chi_{LMM}^2 = c \cdot \chi_{LM}^2$  in BOLT and SAIGE for GWAS
  - Correction factor  $c$  may have **big variation**.
  - **Loss of accuracy** in test statistics
- Genomic prediction
  - **Scalable** and **accurate** Bayesian approaches are lacking.

# Methods

# SNP-set Genomic Prediction (SSGP)

SNP grouping

$$\vec{y} = \mathbf{X}\vec{b} + \sum_{h=1}^p \mathbf{K}_h \vec{u}_h + \vec{e}$$

$$b_j \sim N(0, \sigma_b^2), j = 1, \dots, q$$

SNP weighting

$$\vec{u}_h \sim MVN(\mathbf{0}, \mathbf{W}_h \sigma_{u_h}^2), h = 1, \dots, p$$

Prior 1

$$\sigma_{u_h}^2 \sim \text{Inv-Gamma}(a_{u_h}, b_{u_h}), h = 1, \dots, p$$

$$\vec{e} \sim MVN(0, \mathbf{R}\sigma_e^2)$$

$$\sigma_e^2 \sim \text{Inv-Gamma}(c_0, d_0)$$

Prior 2

$$\sigma_{u_h} \sim \text{Half-Cauchy}(A_{u_h}), h = 1, \dots, p$$

No grouping  
Bayes A

No grouping  
Horseshoe

# Mean field approximation

$$P(\mathbf{b}, \mathbf{u}, \sigma^2 | D) \approx Q(\mathbf{b}, \mathbf{u}, \sigma^2) = q(\mathbf{b}) \prod_{h=1}^P q(\mathbf{u}_h) \prod_{h=1}^P q(\sigma_{u_h}^2) q(\sigma_e^2)$$

Consider proximity-based SNP grouping of equal size ( $S$ ) ...

$S=1$ : BayesA

$S=M$ : GBLUP

---

$D_{\text{KL}}(Q||P) \downarrow$  as  $S \uparrow$

Prediction accuracy of  $P$  may  $\uparrow$  as  $S \downarrow$

We want smaller  $D_{\text{KL}}(Q||P)$  and higher-accuracy  $P$ .

A small  $S$  may work well for genomic prediction.



# Variational inference (VI)



**When variance components  
are **not known****

Marker effect estimates are  
**biased**.



**When variance components  
are **known****

VI is equivalent to block-wise  
Gauss-Seidel method.

Marker effect estimates are **BLUP**.

# Association testing

$$V = \prod_{h=1}^P K_h W K_h' \sigma_{u_h}^2 + R \sigma_e^2$$
$$\chi_{LMM}^2 = \frac{(z' V^{-1} \tilde{y})^2}{z' V^{-1} z}$$

- Note  $V^{-1} \tilde{y} = R^{-1} \sigma_e^{-2} \hat{e}$ .
- We compute  $\hat{e}$  by block-wise Gauss-Seidel method.
- We compute  $V^{-1} z$  similarly.

---

$$V_h = K_h W K_h' \sigma_{u_h}^2 + R \sigma_e^2$$

$$\chi_h^2 = \frac{(z' V_h^{-1} \tilde{y})^2}{z' V_h^{-1} z} \text{ for } z \text{ in group } h.$$

- $\chi_{LMM}^2 = c \cdot \chi_h^2$
- $S \in [1000, 5000]$

Fast and accurate approximation!

# Time complexity

- Scaling linearly in **group size** ( $S$ ), **animals** ( $N$ ), and **markers** ( $M$ )



# Results

Dairy bull  
data for  
genomic  
prediction

20K old bulls  
as training

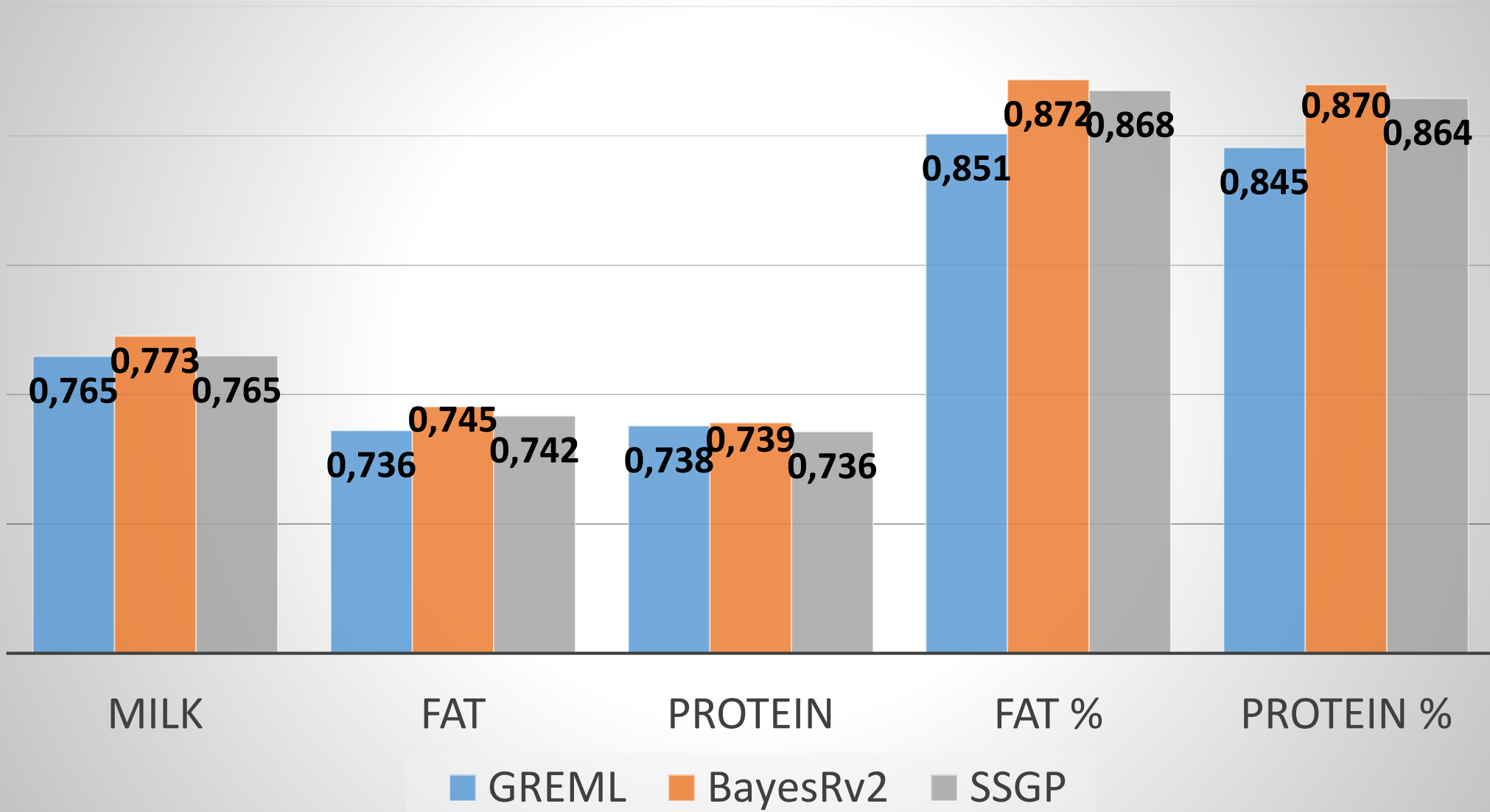
5K young bulls  
as validation

54K SNPs

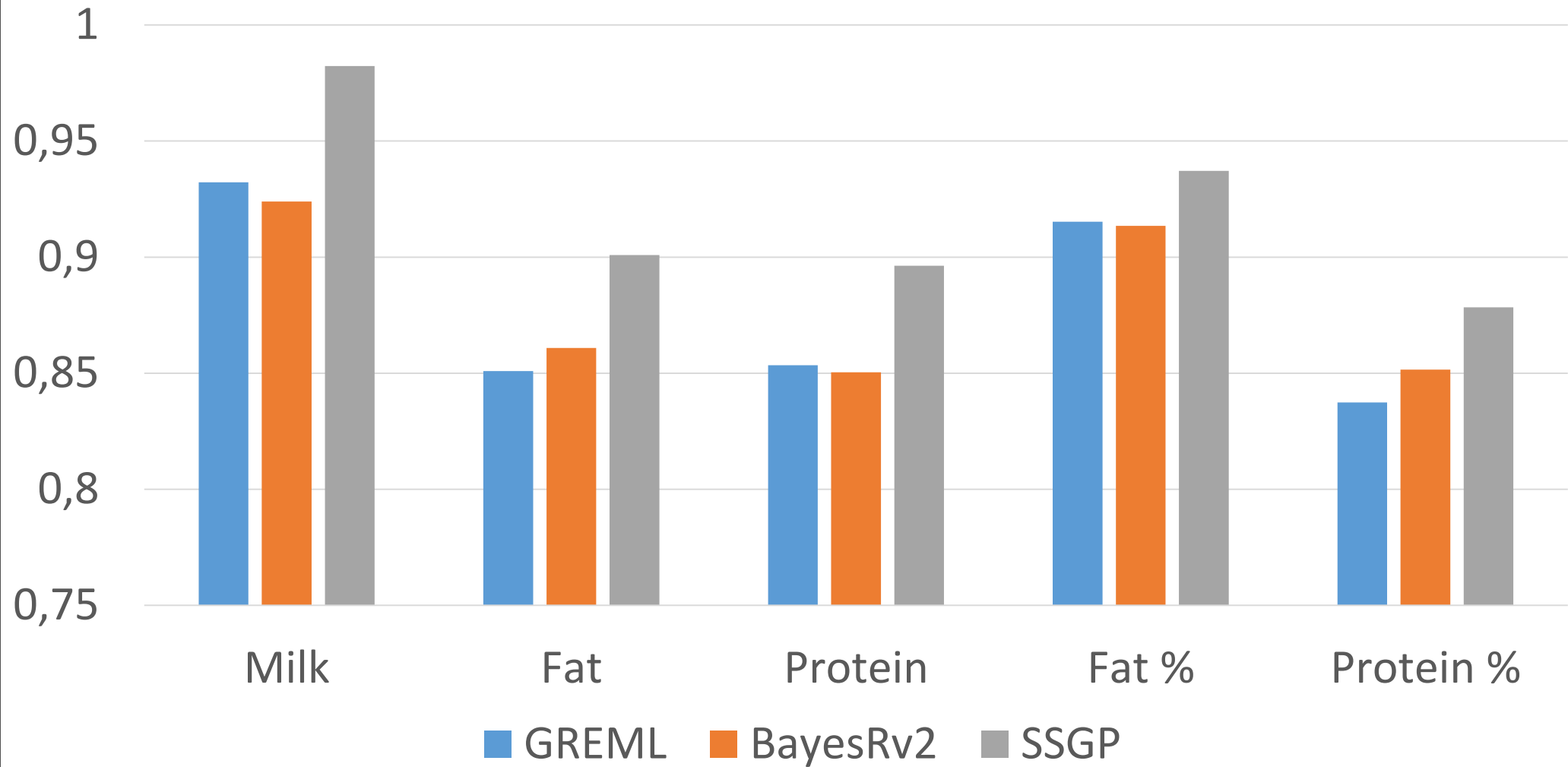
$S=10$  and half-  
Cauchy prior  
used in SSGP

GCTA-GREML  
and BayesRv2  
as benchmark

# Correlation between EBV and phenotypes



## Slope of phenotypes on EBVs



# Time with single core on Intel Xeon E5-2680



~3 min. by SSGP



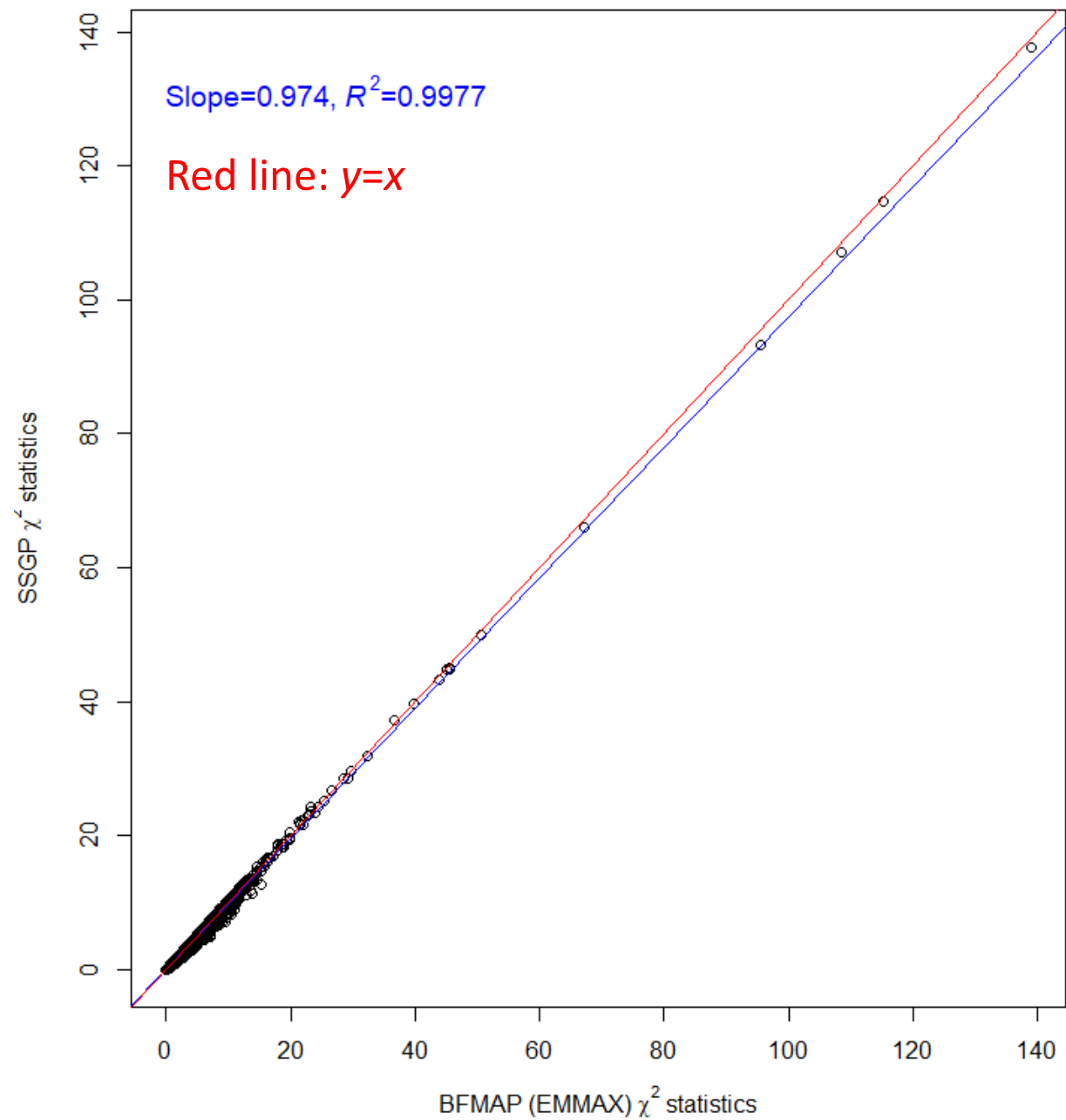
~80 min. by BayesRv2



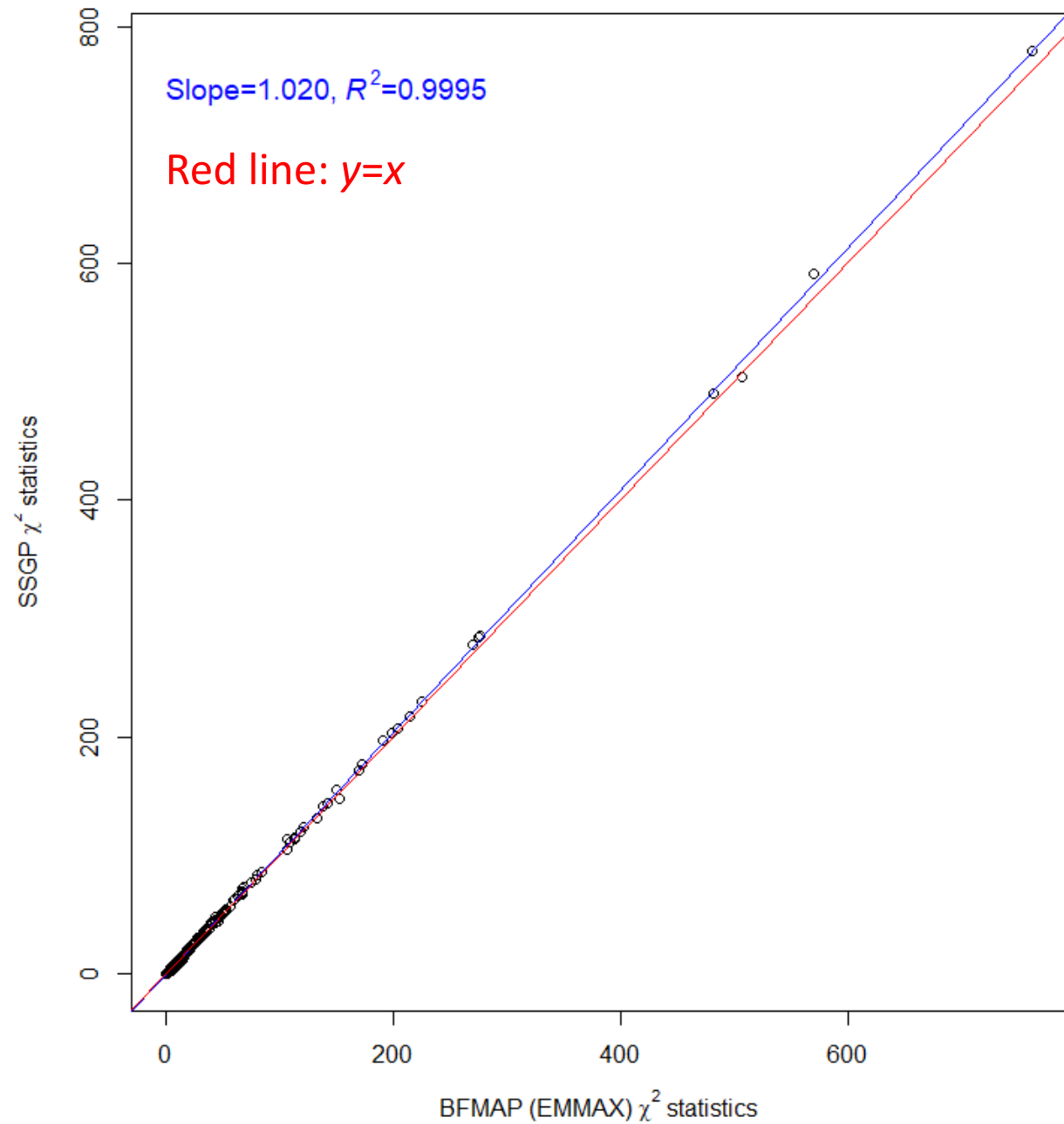
# Cow data for GWAS

- 300K cows with yield deviations
- 60K SNPs
- SAIGE as benchmark
  
- 10K cows randomly sampled from 300K
  - 100 replicates
  - BFMAP (like EMMAX but 15X faster) as benchmark
  - Slope and  $R^2$  of  $\text{Im}(\chi_{\text{SSGP}}^2 \sim \chi_{\text{BFMAP}}^2 - 1)$  for 60K SNPs

Milk replicate with the lowest  $R^2$



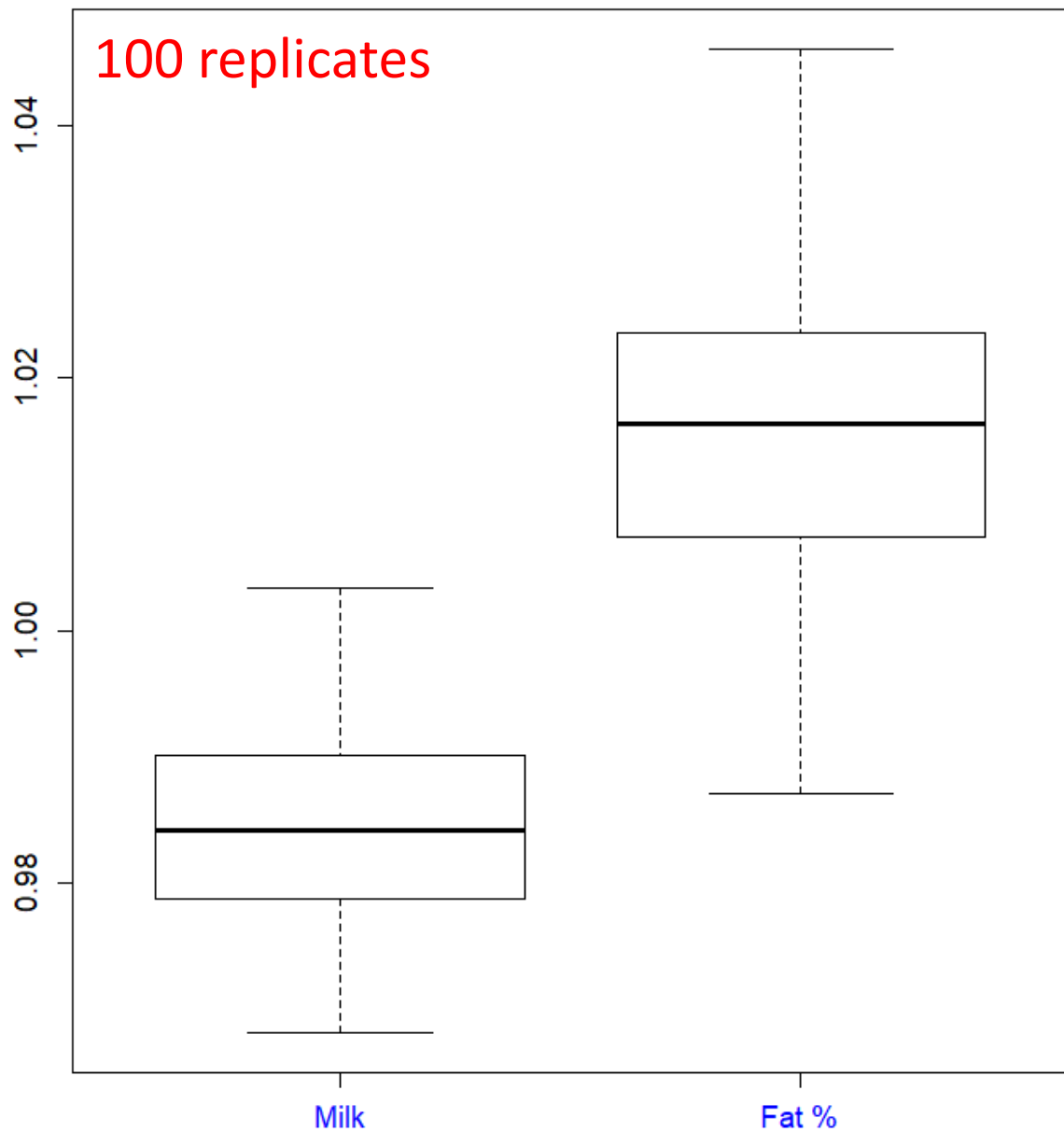
Fat % replicate with the lowest  $R^2$



# Slope

SSGP  $\chi^2$  vs BFMAP (EMMAX)  $\chi^2$

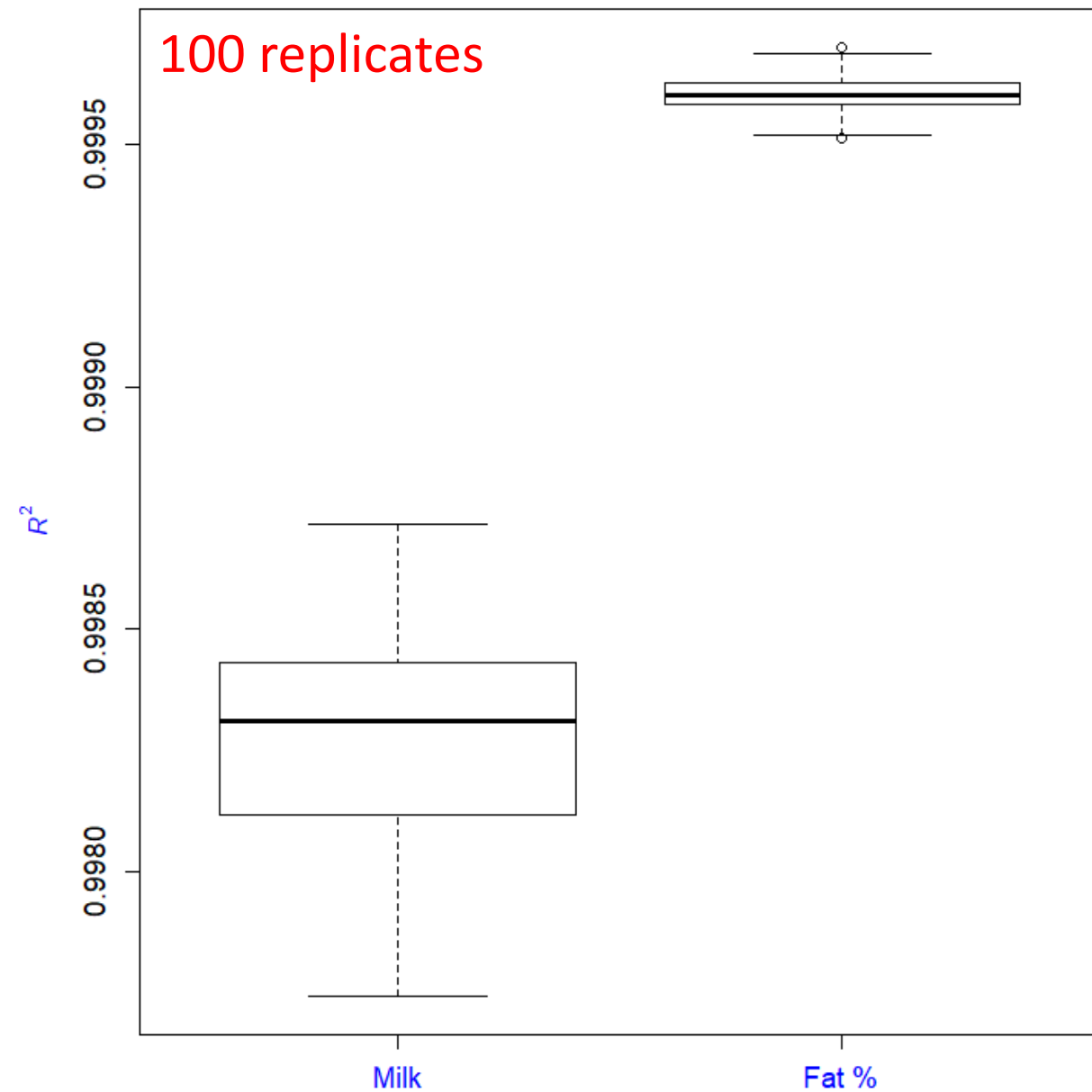
100 replicates



# $R^2$

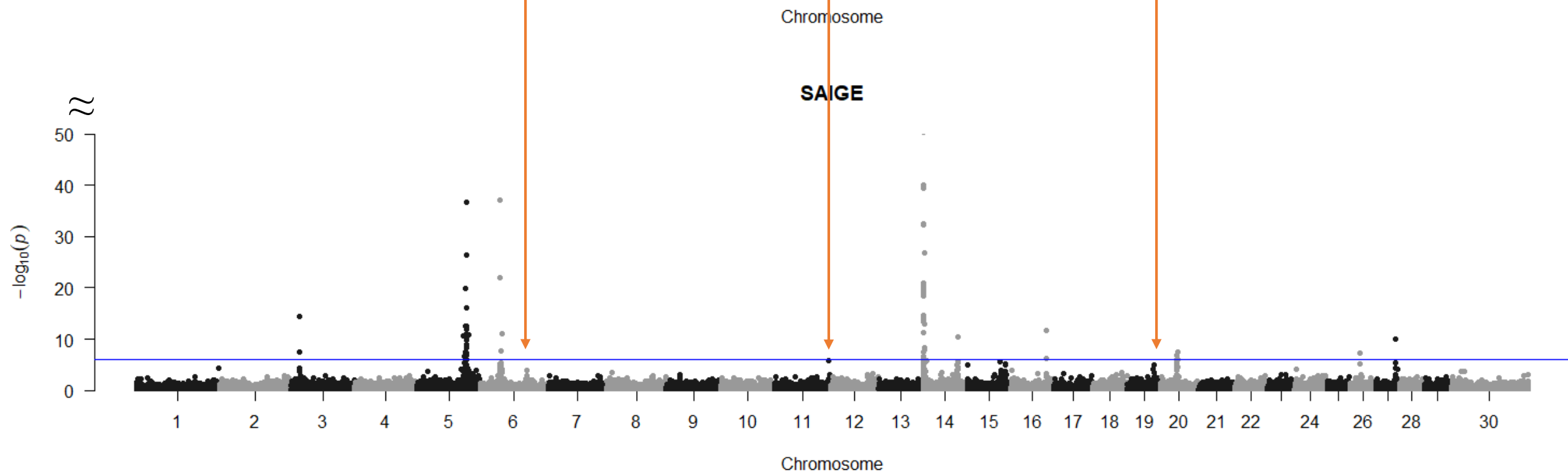
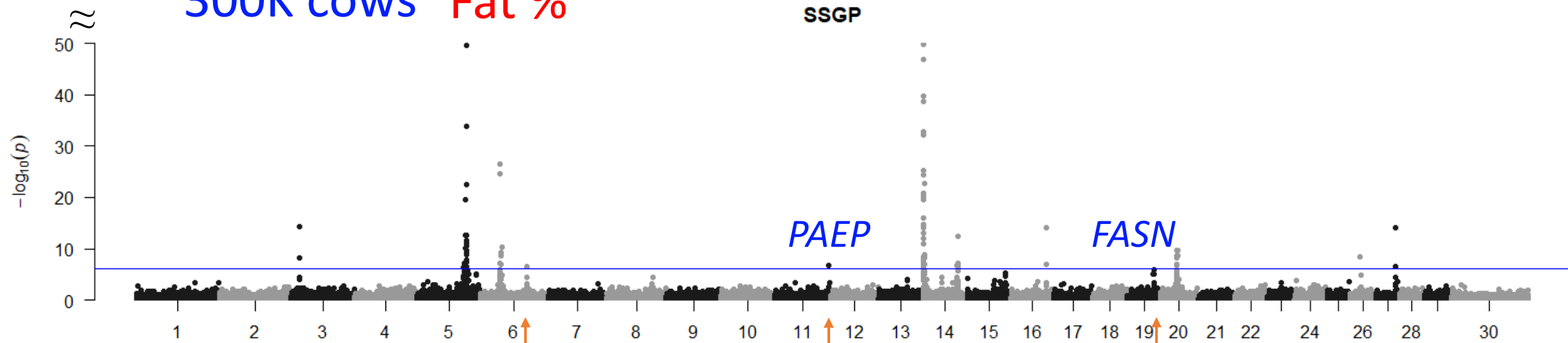
SSGP  $\chi^2$  vs BFMAP (EMMAX)  $\chi^2$

100 replicates



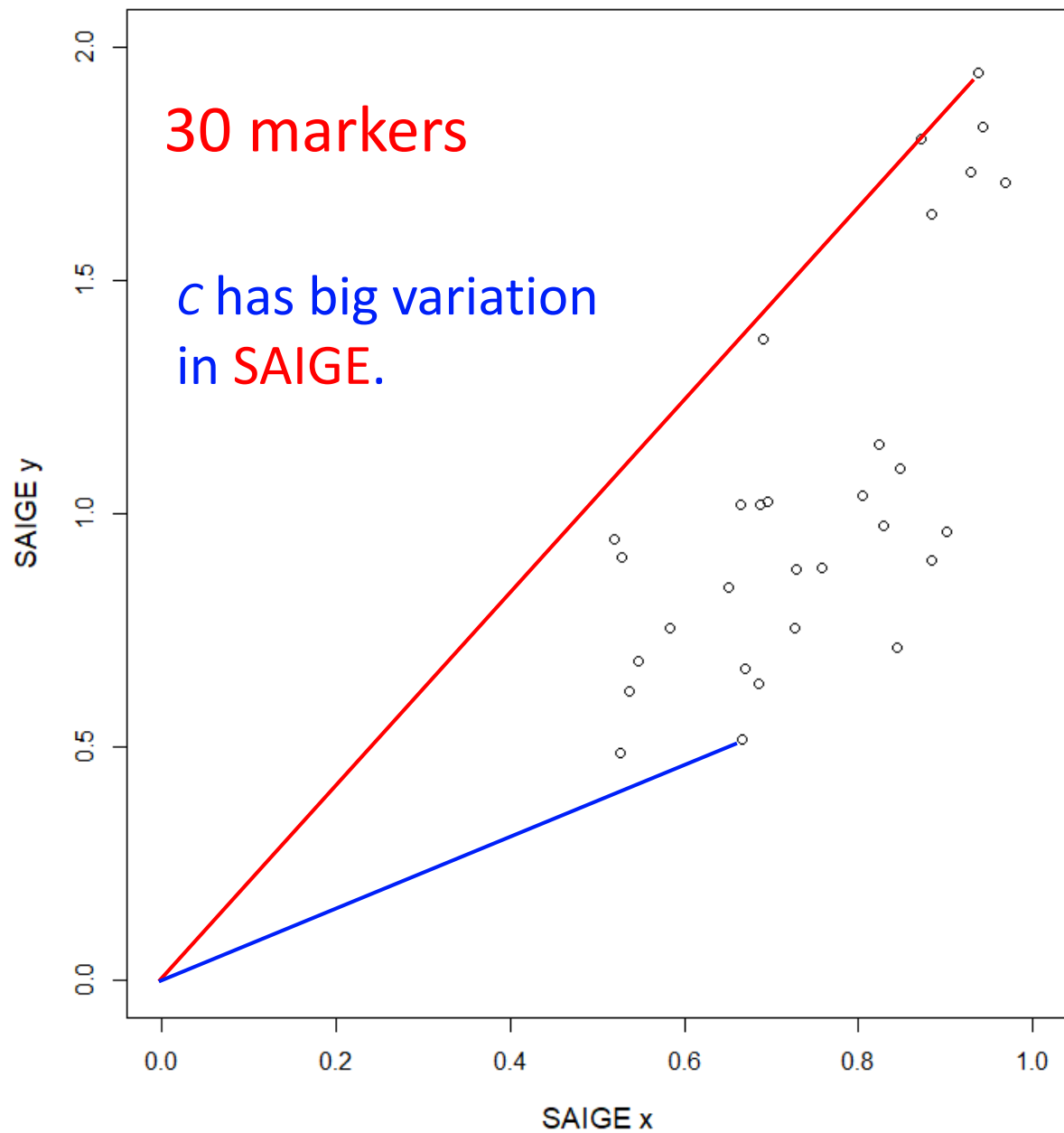
300K cows Fat %

Fat % GWAS with 300K cows  
SSGP



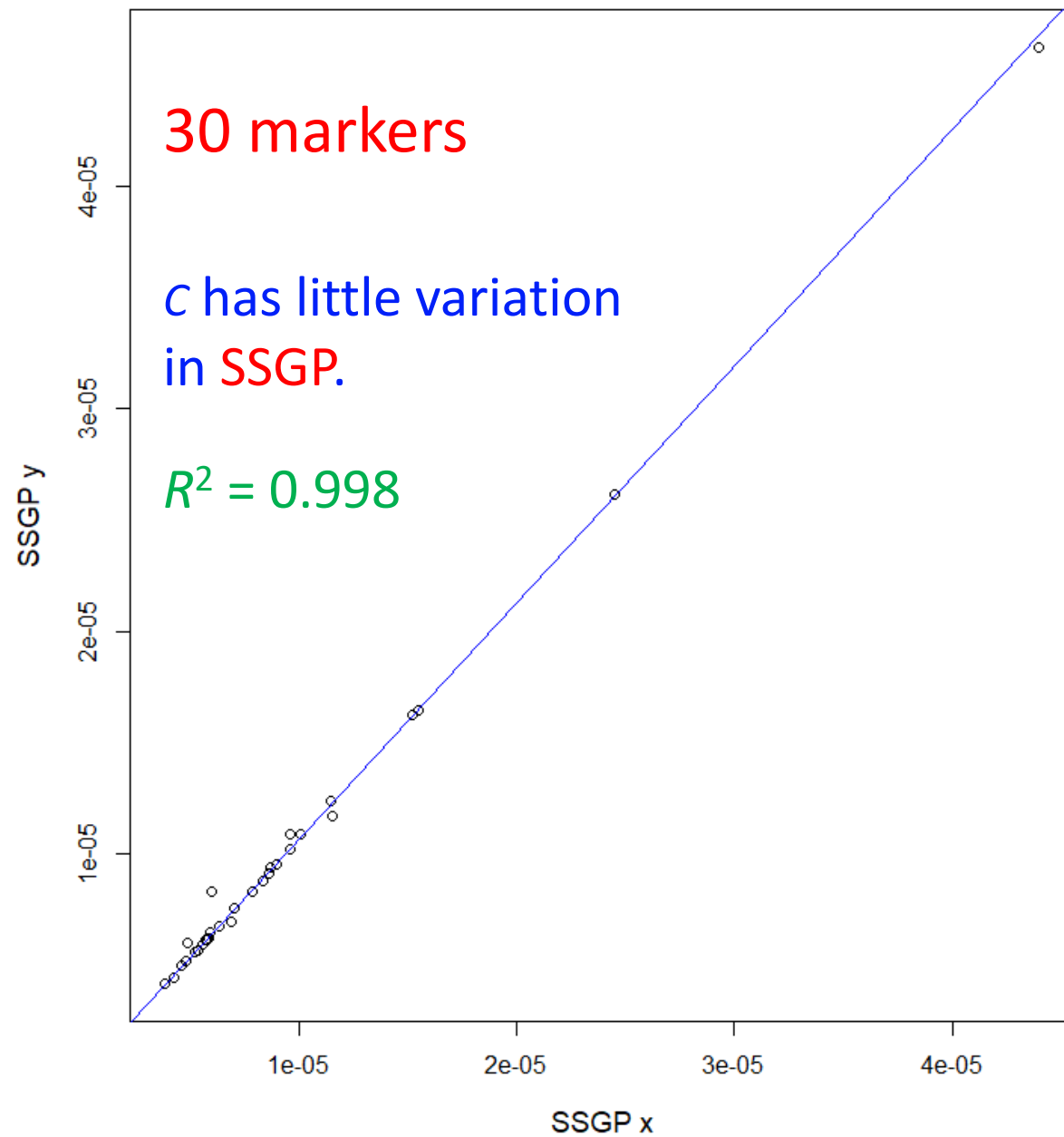
# SAIGE

$$\text{SAIGE } \hat{c} = \text{mean}(y/x)$$



# SSGP

$$\text{SSGP } \hat{c} = \text{slope of } \text{Im}(y \sim x - 1)$$



# Time on MacBook Pro (Intel i9) for 300K cows



~2.7 hours for GWAS by SSGP



>200 hours by SAIGE

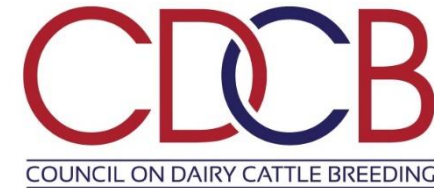
Estimating the correction factor is time-consuming.

# Summary

- SSGP can be applied to various types of samples.
  - **Mostly unrelated**, like UK Biobank
  - **Highly related**, like dairy cattle
  - **Admixed samples**
- SSGP is **accurate** for GWAS and for genomic prediction.
- SSGP is **fast**.
  - 1 million animals and 60K SNPs: <10 hours for GWAS and <5 hours for computing SNP effects on standard hardware.
- SSGP can be applied to **sequence GWA**.
  - **Reasonable** increase in computation compared to **linear regression**

## Acknowledgements

- Funding NHLBI U01 HL137181-01



---

## Software

- SSGP
  - <https://sites.google.com/view/ssgp>
- BFMAP
  - <https://jiang18.github.io/bfmap/>
  - GWA is currently not available in the online version.

## Contact

- Jicai Jiang
  - [jicai.jiang@gmail.com](mailto:jicai.jiang@gmail.com)