



UNIVERSITY OF
GEORGIA

College of Agricultural &
Environmental Sciences

Big Data Genomic Analysis in Dairy Cattle

Daniela Lourenco

A. Cesarani, S. Tsuruta, A. Legarra, E. Nicolazzi, P. VanRaden, I. Misztal

January 14, 2023



The fact

**3.9 million
genotyped animals**

Why is it important?

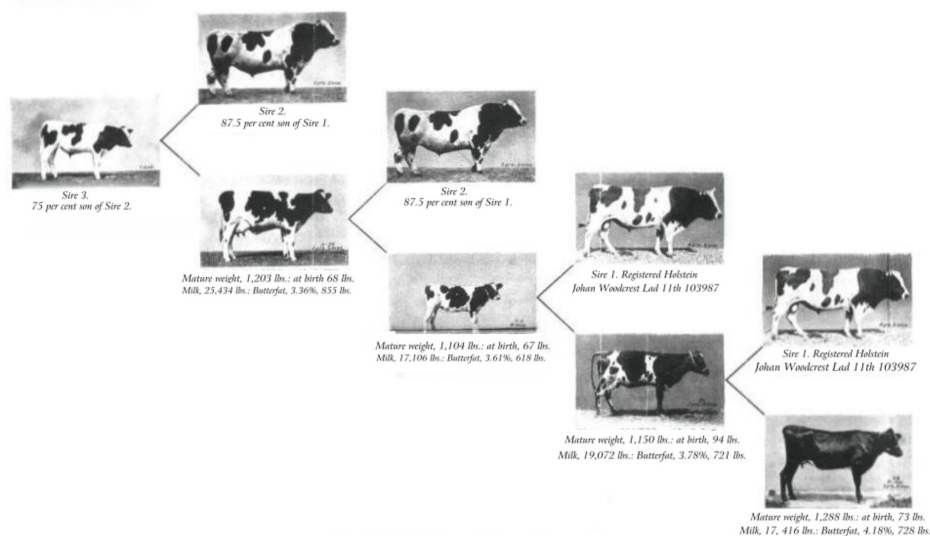
- Impossible in 2015
 - Methods
 - Computing resources vs. algorithms
 - Data availability
- Sends a message
 - Work to accommodate new data
 - Make the most out of available resources

Data

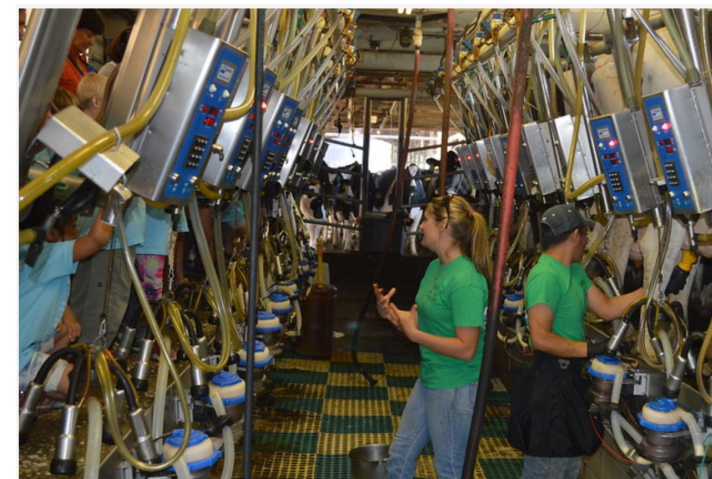
3.9 million

~ 30 million

45 million



VanRaden & Miller



<https://www.usda.gov/media/blog/2020/06/18/data-saydairy-has-changed>

Data

Breed	Phenotypes		Animals	
	N	Cows	Genotypes	Total
All	45M	19.4M	3.9M	29.5M
Ayrshire	116k	47k	9.2k	94k
Brown Swiss	328k	138k	47k	292k
Guernsey	129k	58k	5k	100k
Holstein	40.3M	17.5M	3.4M	26.6M
Jersey	4.1M	1.7M	427k	2.5M

Objectives



Alberto
Cesarani

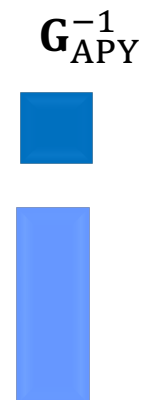
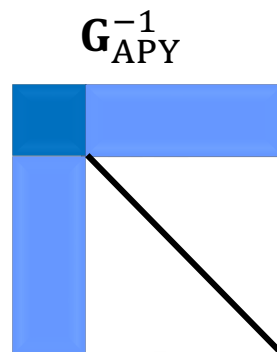
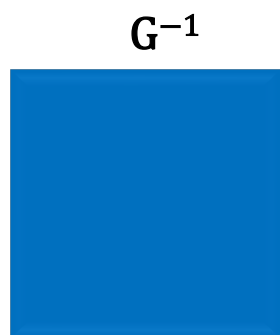
- ssGBLUP multi-breed genomic evaluations for US dairy cattle
- Reliability of GEBV from multi-breed \approx single-breed
- Improve software efficiency \longrightarrow BLUPF90

Single-step (ssGBLUP)

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{W} \\ \mathbf{W}'\mathbf{X} & \mathbf{W}'\mathbf{W} + \mathbf{H}^{-1} \frac{\sigma_e^2}{\sigma_u^2} \end{bmatrix} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{W}'\mathbf{y} \end{bmatrix}$$

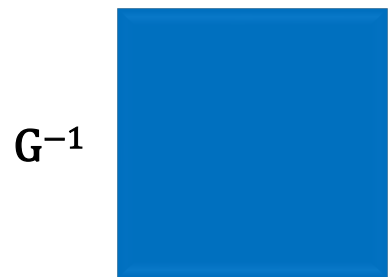
Aguilar et al. (2010)
 Christensen and Lund (2010)

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



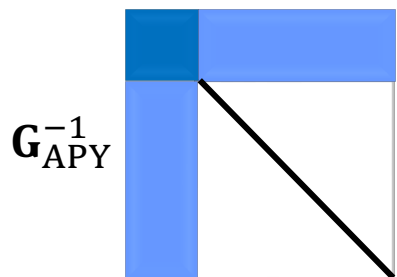
- $\mathbf{G}_{\text{APY}}^{-1}$ sparse
- Efficient computations

APY


 G^{-1}

Dense $\rightarrow u_i | u_1 + u_2 + u_3, \dots, u_{i-1} = \sum_{j=1}^{n-1} p_{ij} u_j + \varepsilon_i$

Genotyped animals \rightarrow Core and Noncore


 G_{APY}^{-1}

Sparse $\rightarrow u_i | u_{c1} + u_{c2} + u_{c3}, \dots, u_{ci} = \sum_{j=1}^c p_{ij} u_j + \varepsilon_i$

$$G_{APY}^{-1} = \begin{bmatrix} G_{CC}^{-1} & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} -G_{CC}^{-1} G_{CN} \\ I \end{bmatrix} M_{NN}^{-1} \begin{bmatrix} -G_{NC} G_{CC}^{-1} & I \end{bmatrix}$$

Misztal et al. (2014)
 Fragomeni et al. (2015)
 Lourenco et al. (2015)
 Masuda et al. (2016)

core animals = # eigenvalues of \mathbf{G} explaining $\geq 98\%$ variance

- Cattle 10k to 15k
- Pigs and chicken 4k to 8k

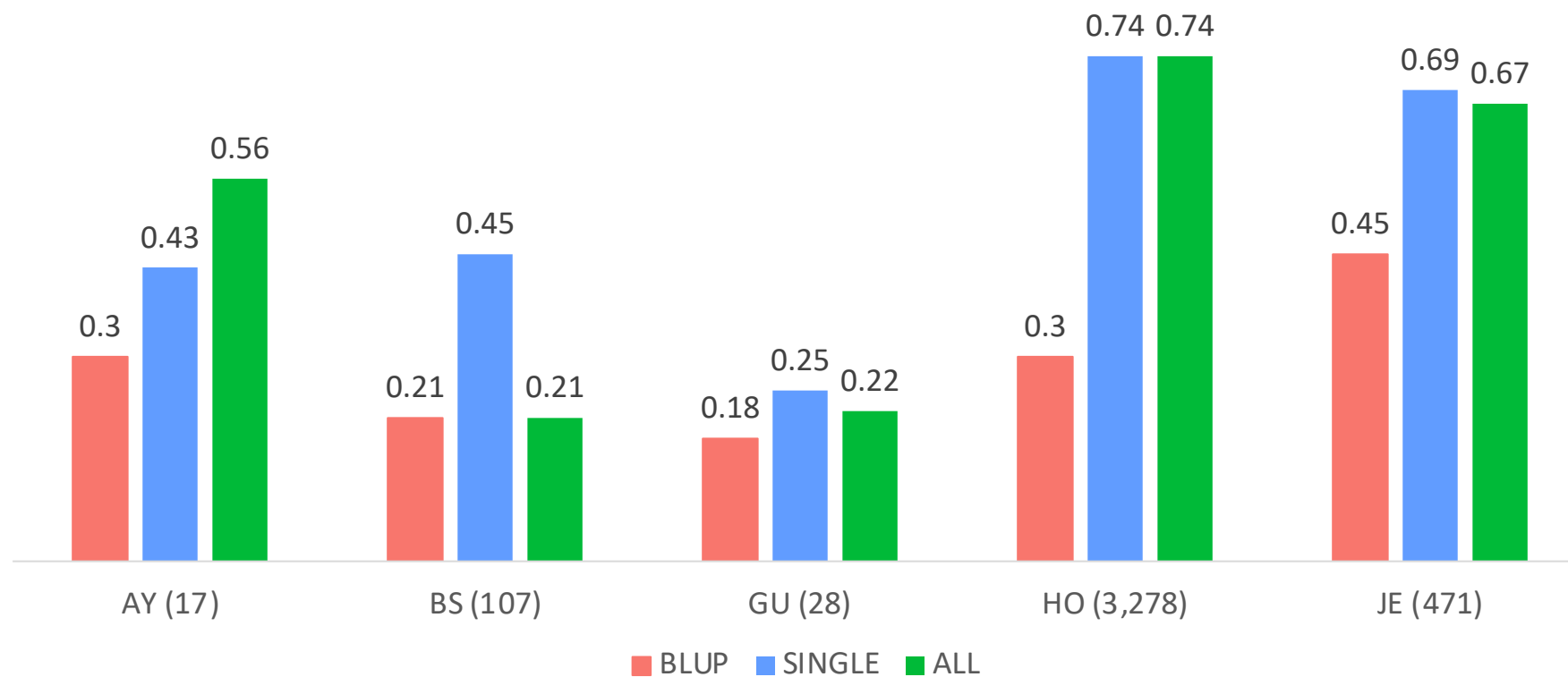
Evaluations

- **SINGLE BREED**
 - each breed separately
 - HO and JE: APY ssGBLUP with 15k random core animals
- **ALL BREEDS**
 - five breeds together
 - breed-specific effects
 - 15k random core animals in APY

Milk (MY), fat (FY), and protein (PY) yields recorded from January 2000 to June 2020

R² for bulls

Protein



Core animals in ssGBLUP ALL

AY = 32

BS = 182

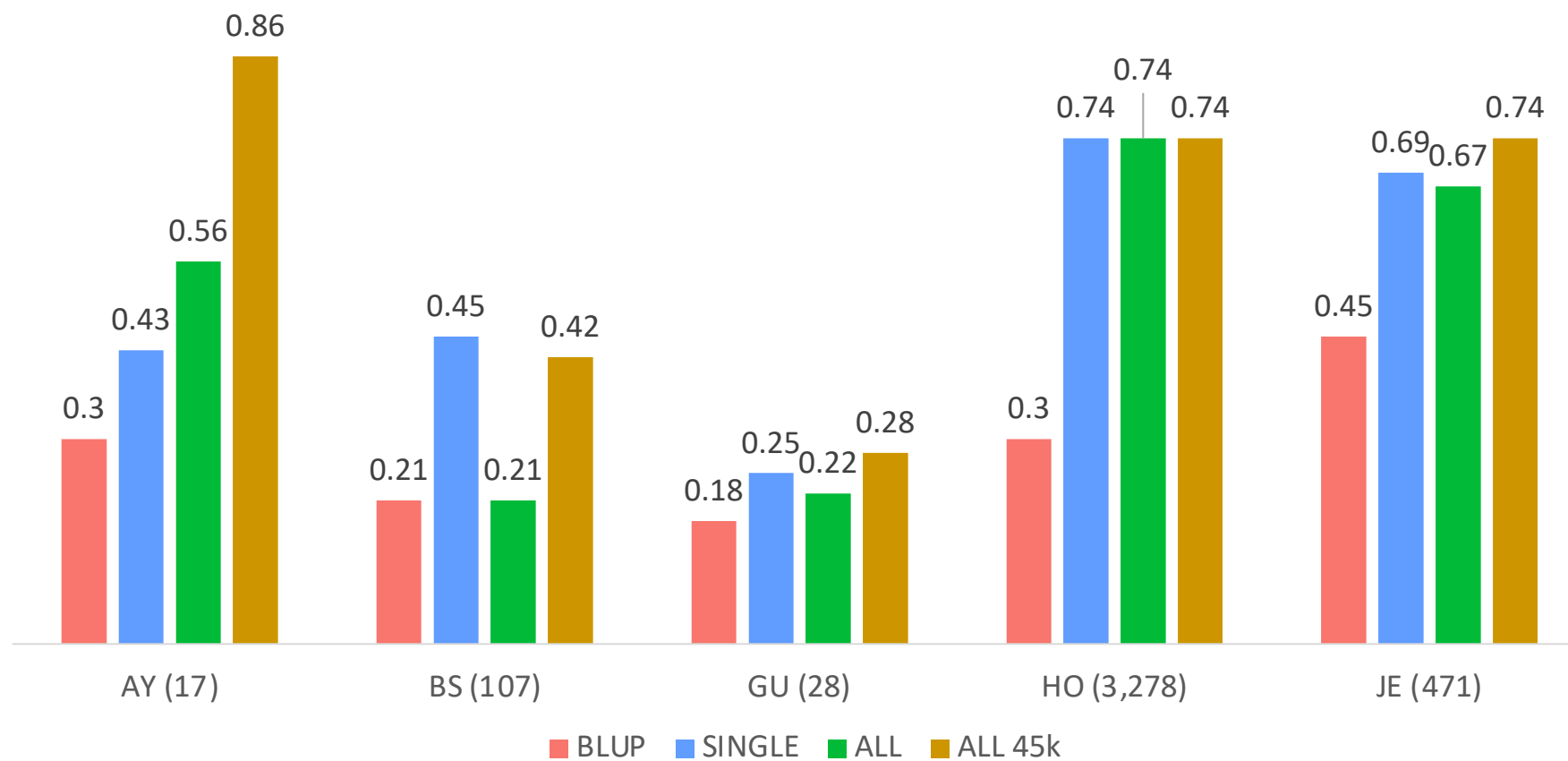
GU = 17

HO = 13k

JE = 1.7k

R² for bulls

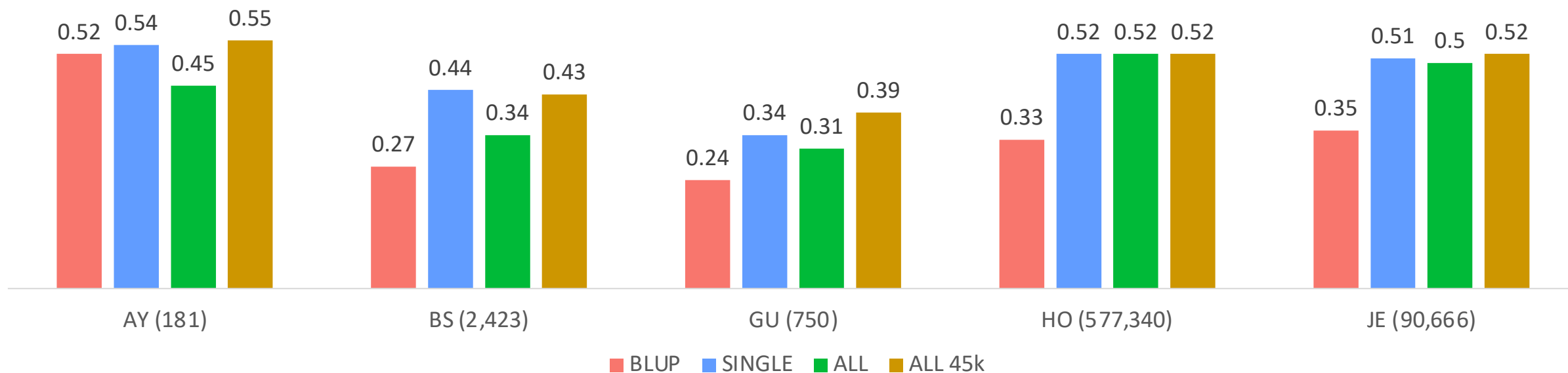
Protein



ALL 45k
 AY = 5k
 BS = 5k
 GU = 5k
 HO = 15k
 JE = 15k

Predictivity for cows

Predictivity for cows - Protein



Computing time

	Rounds		Sec / round		Time	
	BLUP	ssGBLUP	BLUP	ssGBLUP	BLUP	ssGBLUP
AY	504	863	0.08	0.08	< 1 min	~ 1 min
BS	364	867	0.18	0.45	1 min	~ 6 min
GU	345	757	0.07	0.07	< 1 min	< 1 min
HO	457	473	21.25	56.31	2.7 h	7.4 h
JE	586	432	2.00	5.58	~ 20 min	~ 40 min
ALL		1,142		64.84		~ 20 h
ALL_45k	643	1,763	27.01	130.68	4.8 h	~ 64 h

- 2.5 days for solutions
- 5 days for computing G_{APY}^{-1} and A_{22}^{-1} in ALL_45k

Updates in \mathbf{A}_{22} for blending

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



Matias Bermann

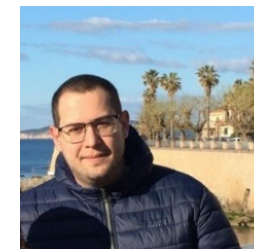
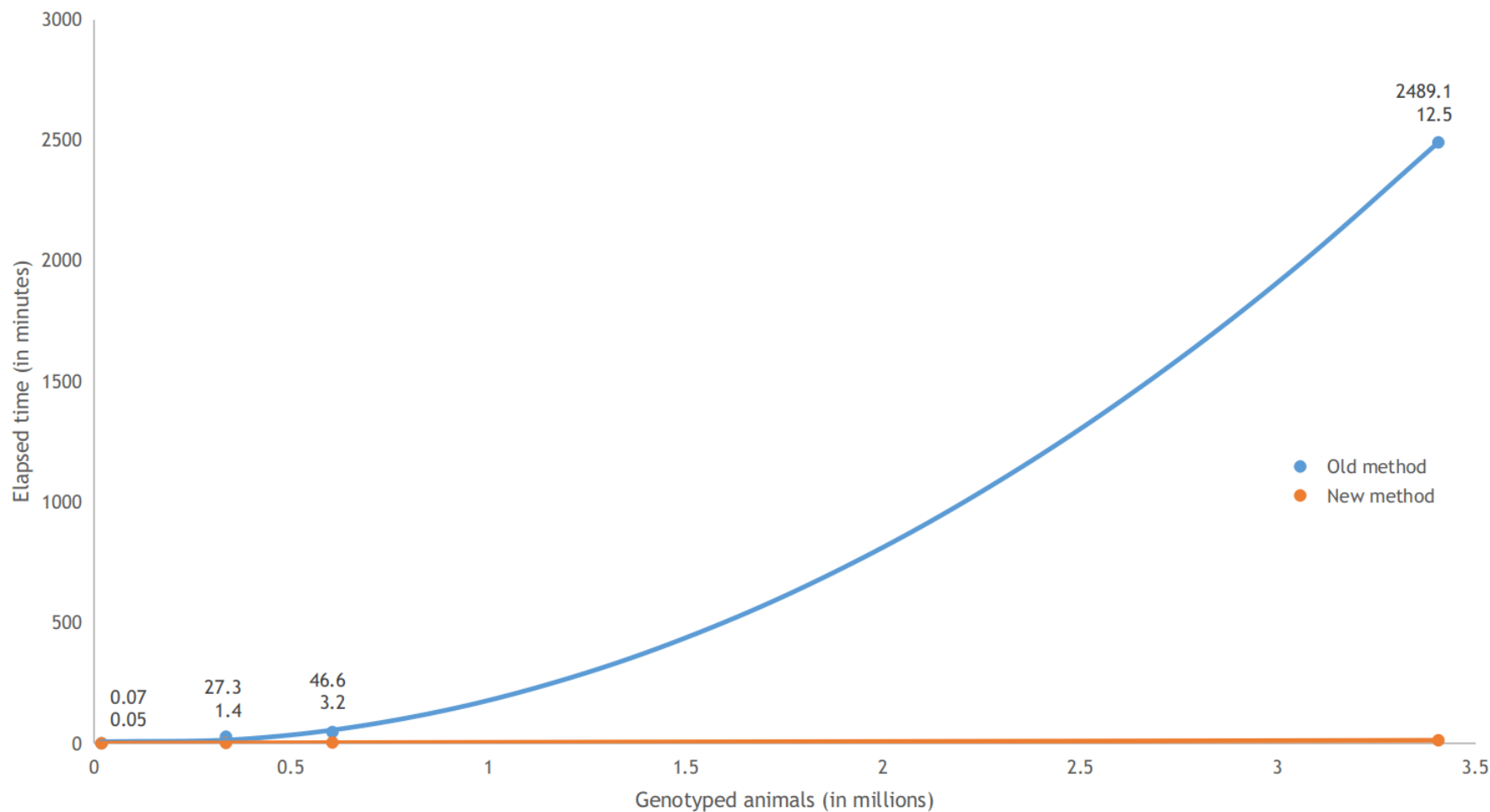
\mathbf{A}_{22}^{-1} components: $\mathbf{A}_{22}^{-1} = \mathbf{A}^{22} - \mathbf{A}^{21}(\mathbf{A}^{11})^{-1}\mathbf{A}^{12}$

APY \mathbf{G}^{-1} : $\mathbf{G}_{\text{APY}}^{-1} = \begin{bmatrix} \mathbf{G}_{\text{CC}}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \begin{bmatrix} -\mathbf{G}_{\text{CC}}^{-1}\mathbf{G}_{\text{cn}} \\ \mathbf{I} \end{bmatrix} \mathbf{M}_{\text{nn}}^{-1} \begin{bmatrix} -\mathbf{G}_{\text{nc}}\mathbf{G}_{\text{CC}}^{-1} & \mathbf{I} \end{bmatrix}$

- Blending: $\mathbf{G} = 0.95 \mathbf{G}^* + 0.05 \mathbf{A}_{22}$
 - Colleau (2002)
 - 3.4M genotyped Holsteins 15k core: 41 hours
 - Rearranging Colleau for core and noncore: 12.5 minutes

Updates in A_{22} for blending

Elapsed time in minutes for old and new method



Alberto
Cesarani

- From 5 days to 8 hours to compute G_{APY}^{-1} and A_{22}^{-1} in ALL_45k

Next challenge



Take home message

Large-scale single-step genomic evaluations are feasible for dairy cattle

- All species
- Reasonable computing time
- Single- or multi-breed populations
 - Respect the dimensionality within each breed

Large-scale genomic evaluations are easier said than done

- Always challenged by the increasing amount of data
- Efficient algorithms and computing power

Acknowledgments



Grant no. 2020-67015-31030

Dairy producers who supplied data through their participation in the Dairy Herd Improvement program and Dairy Records Processing Centers that edited and relayed information on to the Council of Dairy Cattle Breeding