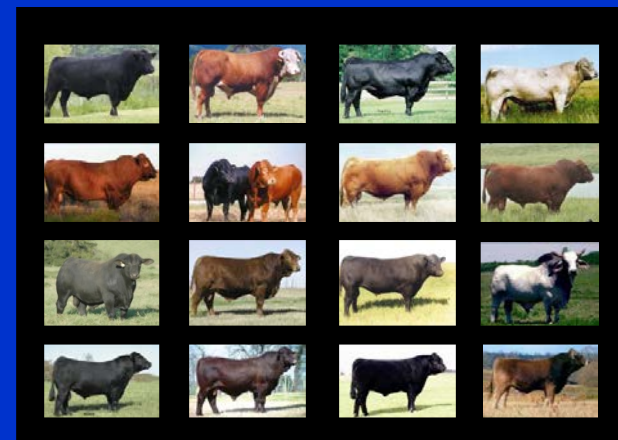


Training, Validation, and Target Populations

Mark Thallman, Kristina Weber,
Larry Kuehn, Warren Snelling,
John Keele, Gary Bennett,
and John Pollak



50,000 Markers on a Chip (50K Chip)



Illumina Infinium Bovine BeadChip

~ 50,000 SNP markers across the bovine genome

-High resolution (1 SNP per 60,000 base pairs)

- Multiple breeds used for SNP discovery

Illumina BeadScan

START → SETUP → TILT → ALIGN → **SCAN** → REVIEW

Green
Red
Overlay

Sertrix ID	Section	G Sat	G P95	G P5	R Sat
1992267171	D	0	4925.09	102.15	0
1992267171	E	0	5001.23	95.73	0
1992267171	F	0	5622.59	91.22	0
1992267171	G	0	6060.61	91.90	0
1992267171	H	0	6396.59	93.55	0
1992267171	I	0	5920.63	92.13	0
1992267171	J	0	6311.91	108.40	0
1992267171	K	0	5796.79	102.44	0
1992267171	L	0	6298.38	103.93	0
1992267172	A	0	4646.20	63.04	0
1992267172	B	0	4649.45	85.21	0
1992267172	C	0	4674.04	90.18	0
1992267172	D	0	4330.30	96.44	0
1992267172	E	0	3923.42	80.57	0
1992267172	F	0	4539.36	83.25	0
1992267172	G	0	5812.91	89.98	0
1992267172	H	0	5247.42	104.53	0
1992267172	I	0	5740.51	88.10	0
1992267172	J	0	4923.59	88.90	0
1992267172	K	0	5804.10	97.25	0
1992267172	L	0	6346.41	91.06	0
1992267352	A	0	4406.97	45.75	0
1992267352	B	0	5097.12	56.95	0
1992267352	C	0	5164.08	42.32	0
1992267352	D	0	5343.97	65.34	0

<< Stop Pause

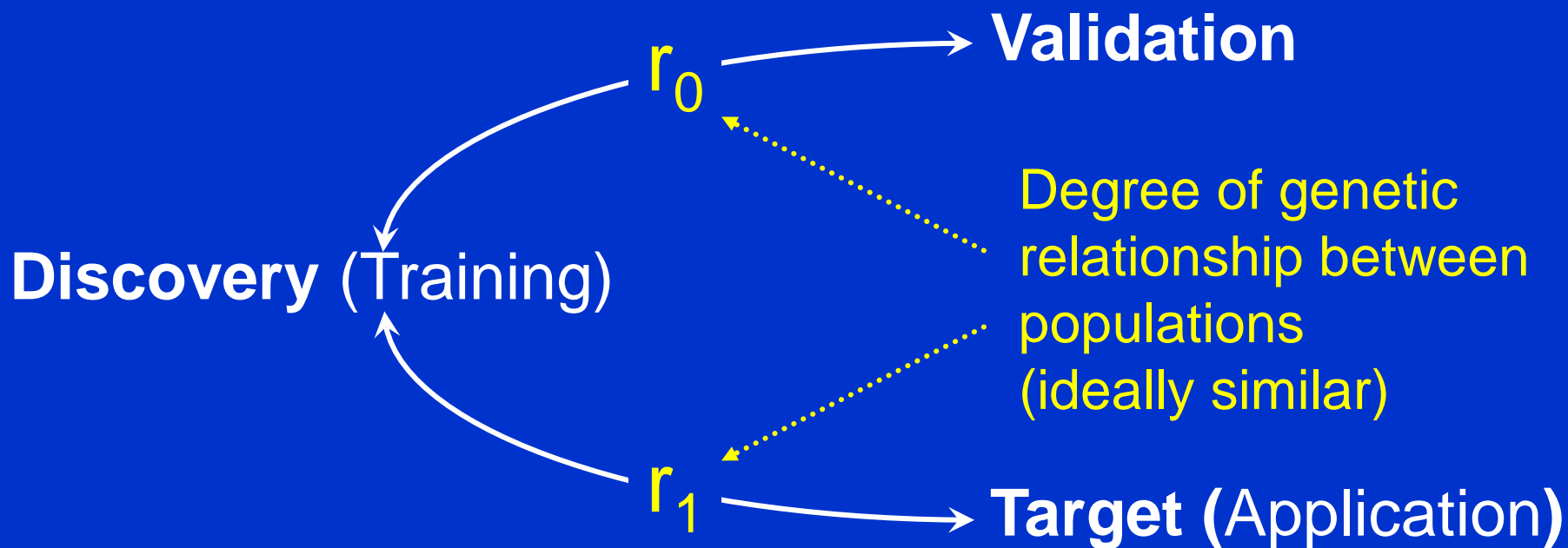
Scanning "1992267352 : L", Capturing Analytic Images

TX Connected at 10.254.0.1 : 4096 RX Status: Good Green: Gain=1.0 (457v) Filter=100% Red: Gain=1.0 (826v) Filter=100% Elapsed: 2:17:55

STATUS EVENTS LOG Save Path: C:\ImageData

BARC (ARS)
USMARC
University of Missouri
University of Alberta

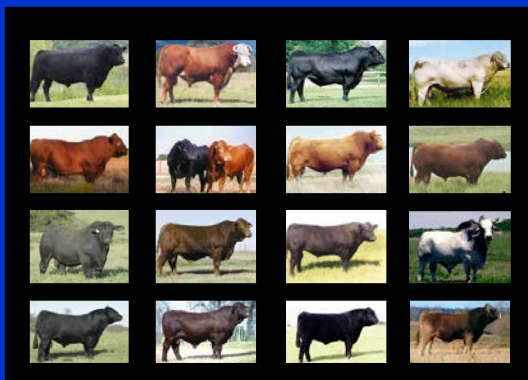
Populations Involved in Genomic Predictions of Economically Important Traits



3 Fundamental Types of Discovery Populations



- Purebreds of a Single Breed



- Purebreds of Multiple Breeds



- Crossbreds

2 Fundamental Types of Discovery Data



- AI Sires with High Accuracy EPDs

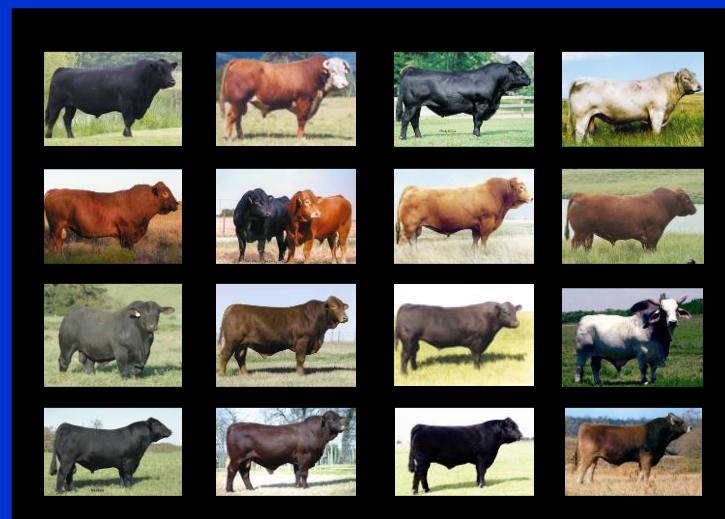
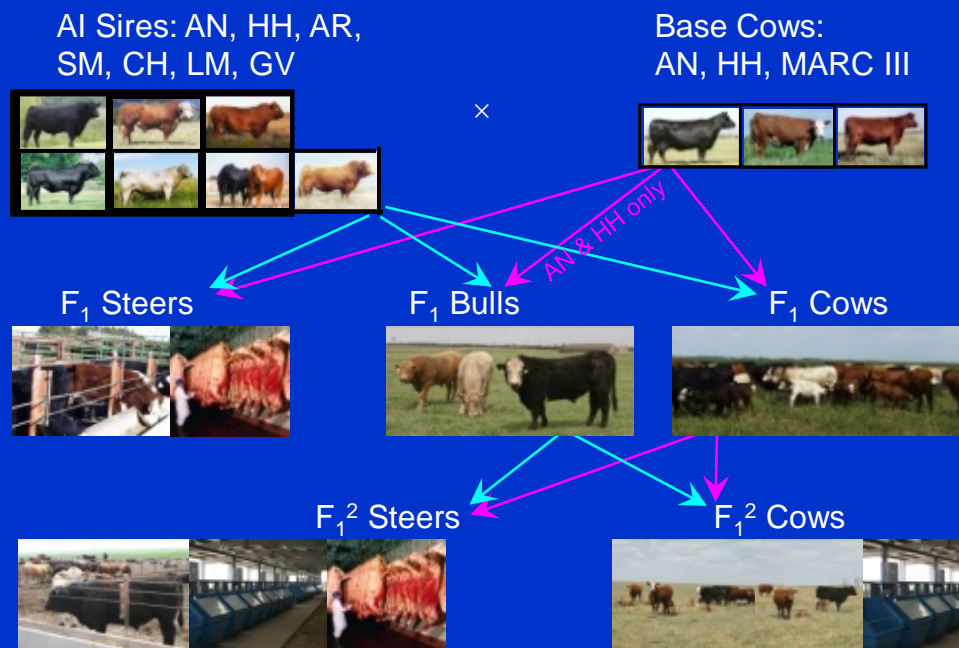


- Individuals with Own Phenotypes

Two Resource Populations at USMARC

USMARC Cycle VII
USMARC Ongoing GPE

2,000 Bull Project



2000 Bull Project

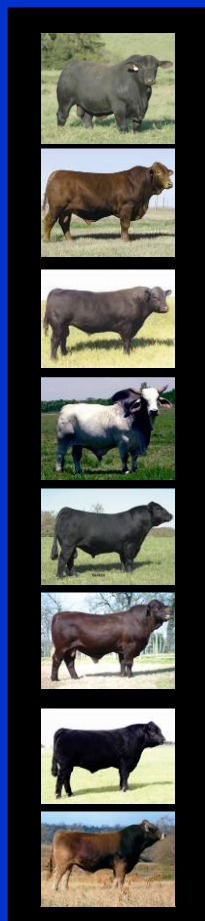


- Collaborative Effort
 - Researchers
 - Breed Associations
- Breed associations provided semen for DNA on influential sires
- USMARC ran the 50K SNP chip on those 2,000 sires
- USMARC provides extensively phenotyped animals for use as training data set

2,000 Bull Project: Number of Sires Sampled



- Angus 402
- Hereford 317
- Simmental 253
- Red Angus 173
- Gelbvieh 136
- Limousin 131
- Charolais 125
- Shorthorn 86



- Brangus 68
- Beefmaster 64
- Maine-Anjou 59
- Brahman 53
- Chiangus 47
- Santa Gertrudis 43
- Salers 42
- Braunvieh 27

2026



Deregression of EPDs

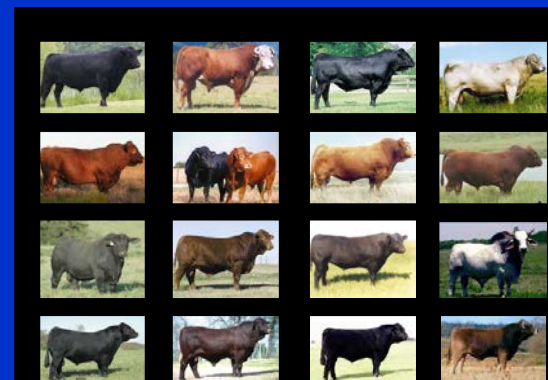


An approach used to scale EPDs so that high and low accuracy animals can be included in the same analysis

- Genetic variances are the same regardless of accuracy
- Residual variances are heterogeneous depending on the accuracy of the animal
- Allows use of EPDs as phenotypes in genomic analyses

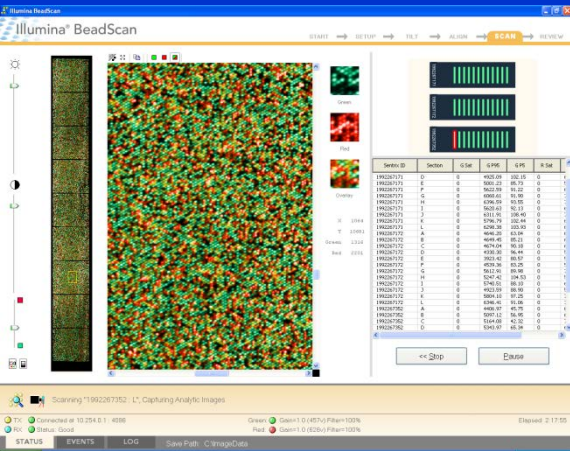


- AI Sires with High Accuracy EPDs



Goals of 2,000 Bull Project

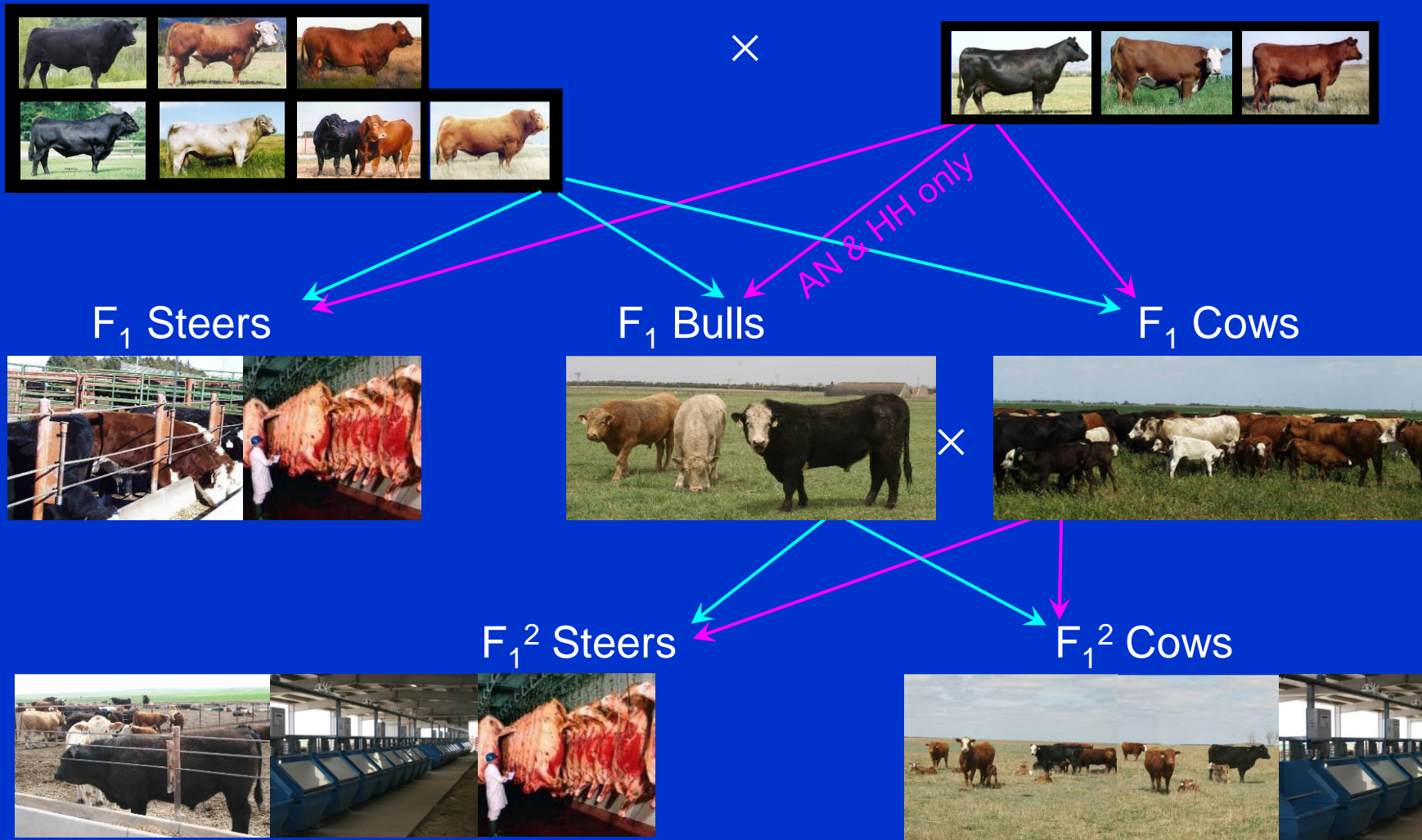
- Demonstrate feasibility and understand challenges of applying whole-genome selection in beef cattle.
- Provide prediction equations for general use
- Provide genomic predictions for the bulls in the project



Training Data: GPE Cycle VII Population

AI Sires: AN, HH, AR,
SM, CH, LM, GV

Base Cows:
AN, HH, MARC III



Training Data: GPE Ongoing Continuous Sampling

AI Sires:

AN, HH, SM, CH, AR, LM, GV, SH,
BN, BM, MA, BR, CI, SG, SA, BV

Dams:

AN, HH, CH, SM,
MARC III, Cycle VII F₁



×



AN, HH, CH, SM

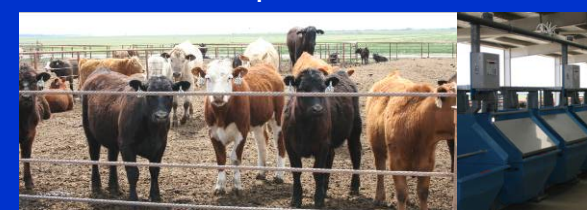
F₁ & BC Steers

F₁ Bulls

F₁ & BC Heifers



×

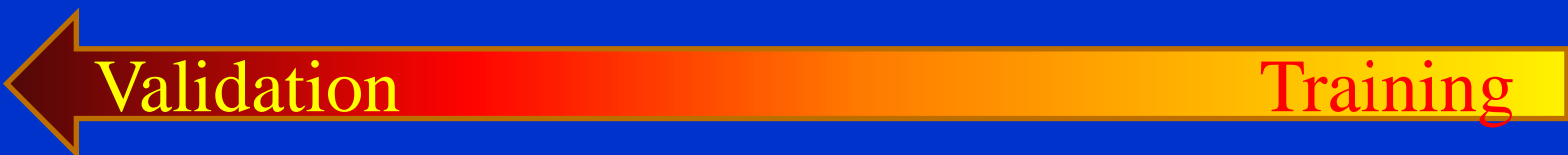
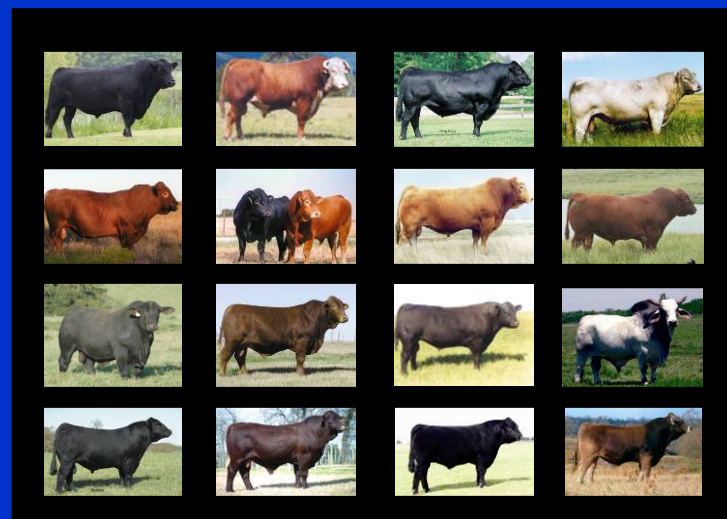
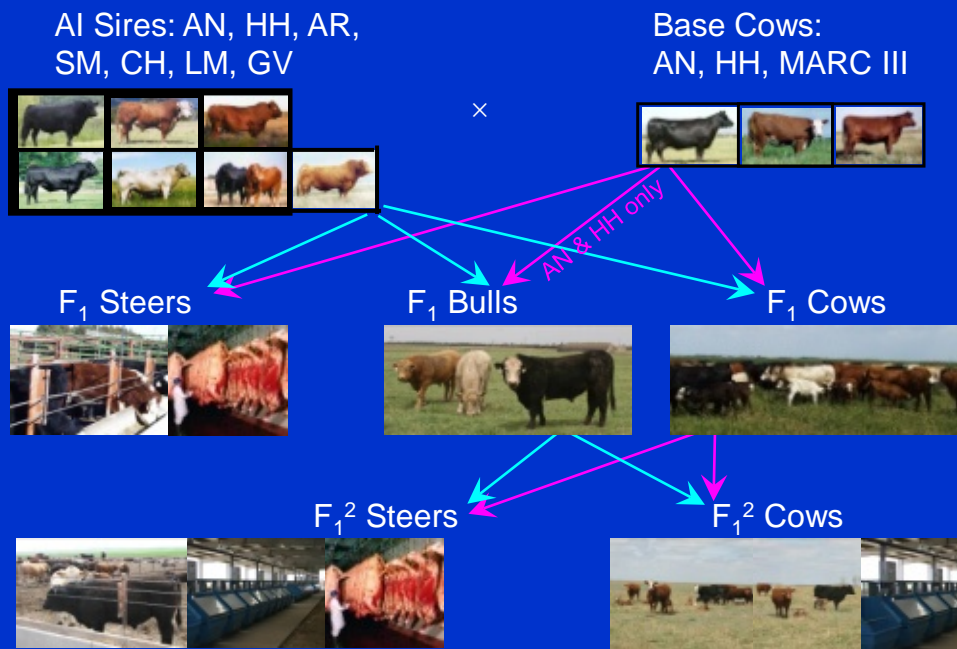


Cross-validation

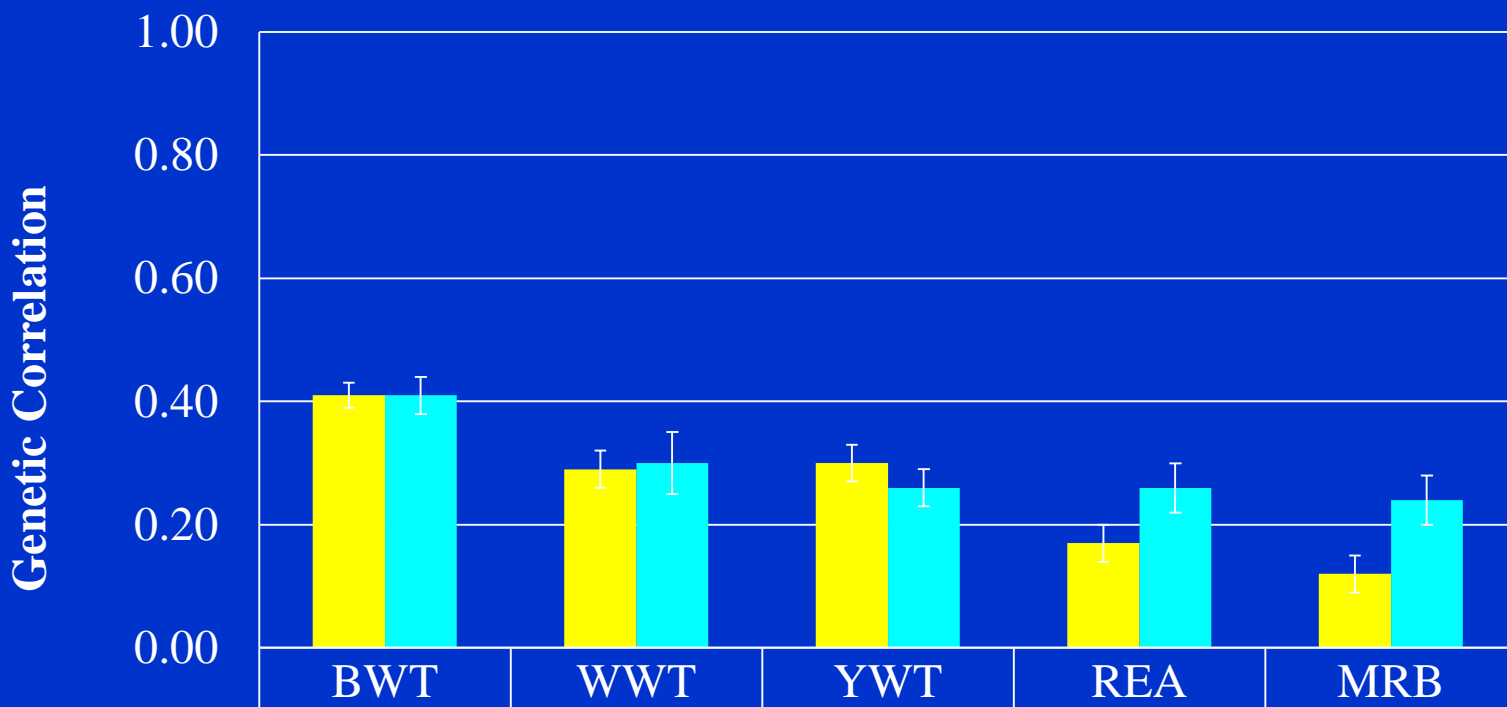


USMARC Cycle VII
USMARC Ongoing GPE

2,000 Bull Project



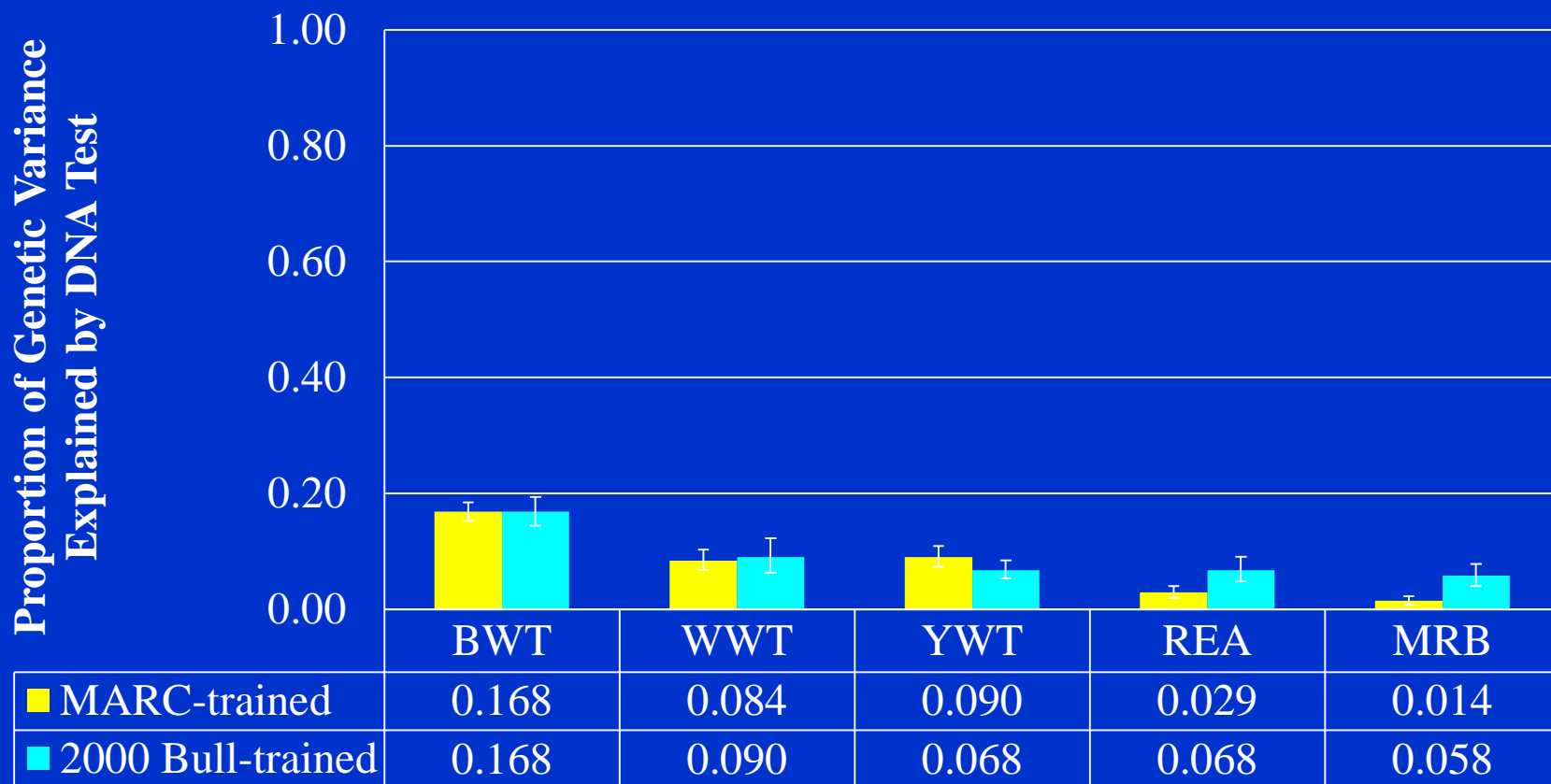
Genetic Correlations in Cross-validation



■ MARC-trained	0.41	0.29	0.30	0.17	0.12
■ 2000 Bull-trained	0.41	0.30	0.26	0.26	0.24

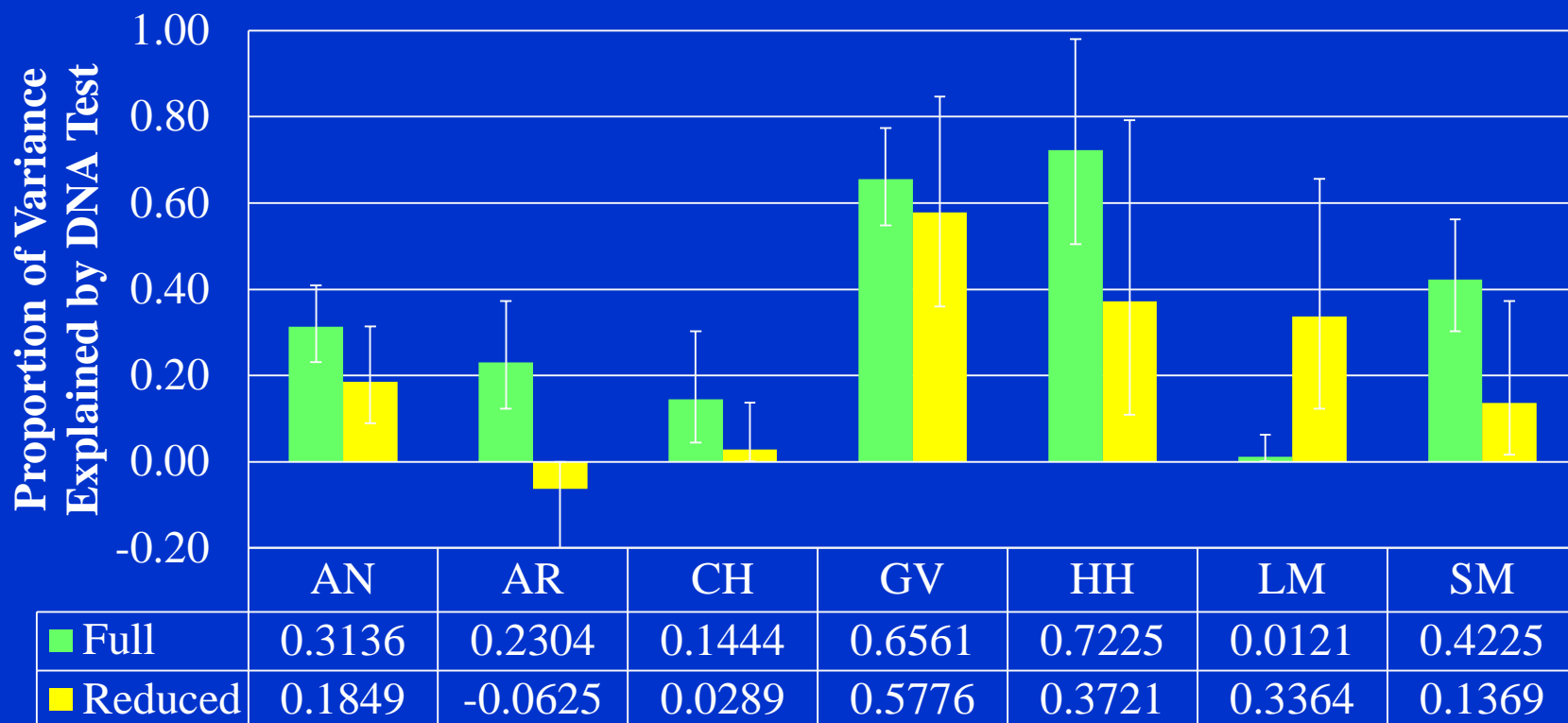
*2,000 bull predictions excluded the sires of the MARC validation populations

Proportion of Genetic Variance Explained in Cross-validation



*2,000 bull predictions excluded the sires of the MARC validation populations

Proportion of Variation in Weight Traits Project from Training on 2,000 Bulls



*Full = Prediction of sires including the 2,000 bulls

*Reduced = Predictions excluding the 2,000 bulls

What Can We Do to Improve Prediction Accuracy?

- Add more phenotypes (animals)
- Increase marker density
- Incorporate individual animal DNA sequence on influential sires
- Do a better job of using the information we have (better statistical analyses)

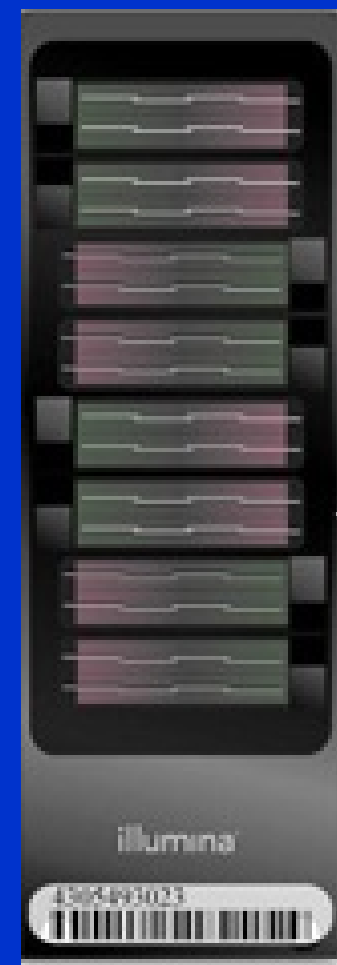
What Can We Do to Improve Prediction Accuracy?

- Add more phenotypes (animals)
- Increase marker density
- Incorporate individual animal DNA sequence on influential sires
- Do a better job of using the information we have (better statistical analyses)

We need to find ways to use information from all available discovery populations, regardless of the target breed(s).

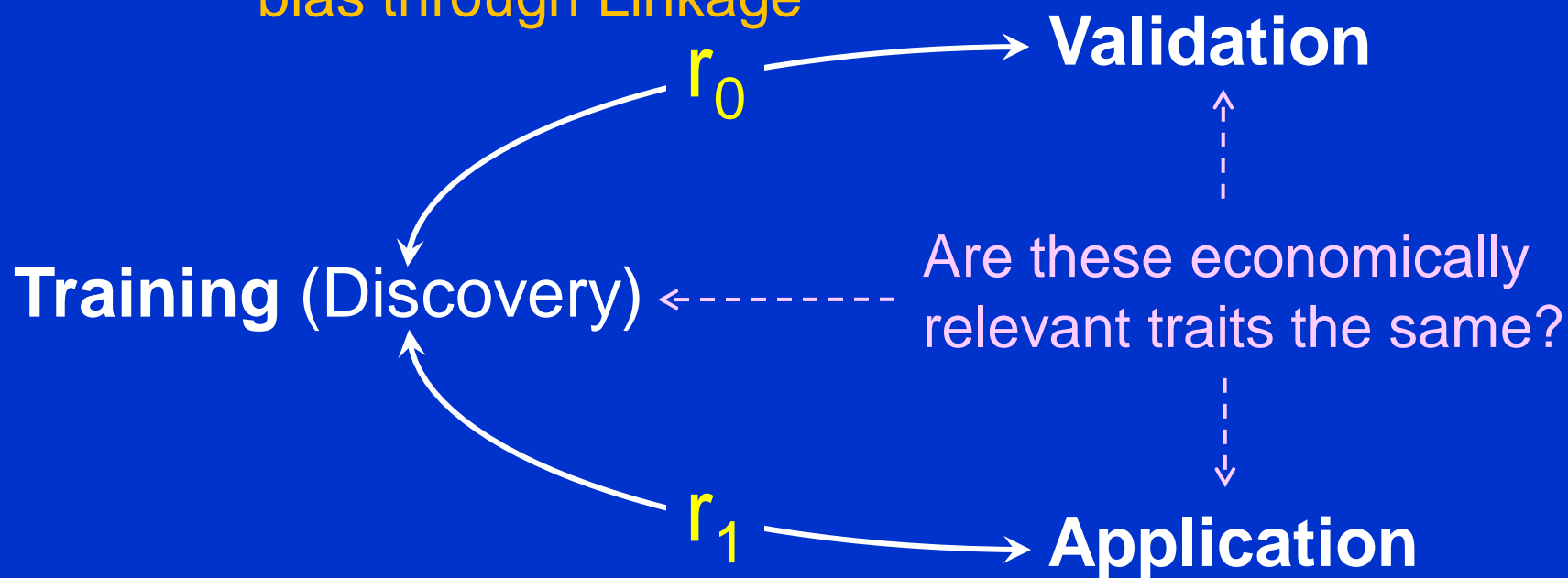
770K Chip

- Illumina BovineHD
- 770,000 Single Nucleotide Polymorphisms (SNP)
- Much higher SNP density than 50K chip
- Should allow predictions to be less breed-specific
- We are in the process of having this new chip run on > 300 sires that have substantial progeny at USMARC



Populations Involved in Whole Genome Selection

This relationship contributes to discovery bias through Linkage



This relationship affects the accuracy of prediction, but the effects erode over time

Prospects for Moving Beyond Validation Populations

- Conceptually, it should be possible to combine discovery and validation into a single step with a single population.
- There is still considerable work to be done to make this practical.
- This concept assumes that the accuracy of the genomic part of each individual's genetic evaluation should vary depending on its genetic relationship to the discovery population.

Prospects for Moving Beyond Validation Populations

- This concept requires that raw genotypes be available on the training/validation population as well as the target population.
- For traits that are routinely recorded in the target population, phenotypes should be continuously integrated into the training/validation population.

Conclusions

- Within-breed predictions based on the 50K work well.
- Training on multiple purebred populations is more effective than training on only a single, small purebred population.
- With increasing marker density, crossbred populations will likely become increasingly important components of training.

