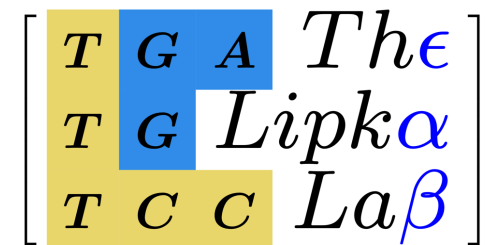


Leveraging bioinformatic breakthroughs into quantitative genetic approaches for crop improvement

I ILLINOIS
Crop Sciences
COLLEGE OF AGRICULTURAL, CONSUMER
& ENVIRONMENTAL SCIENCES



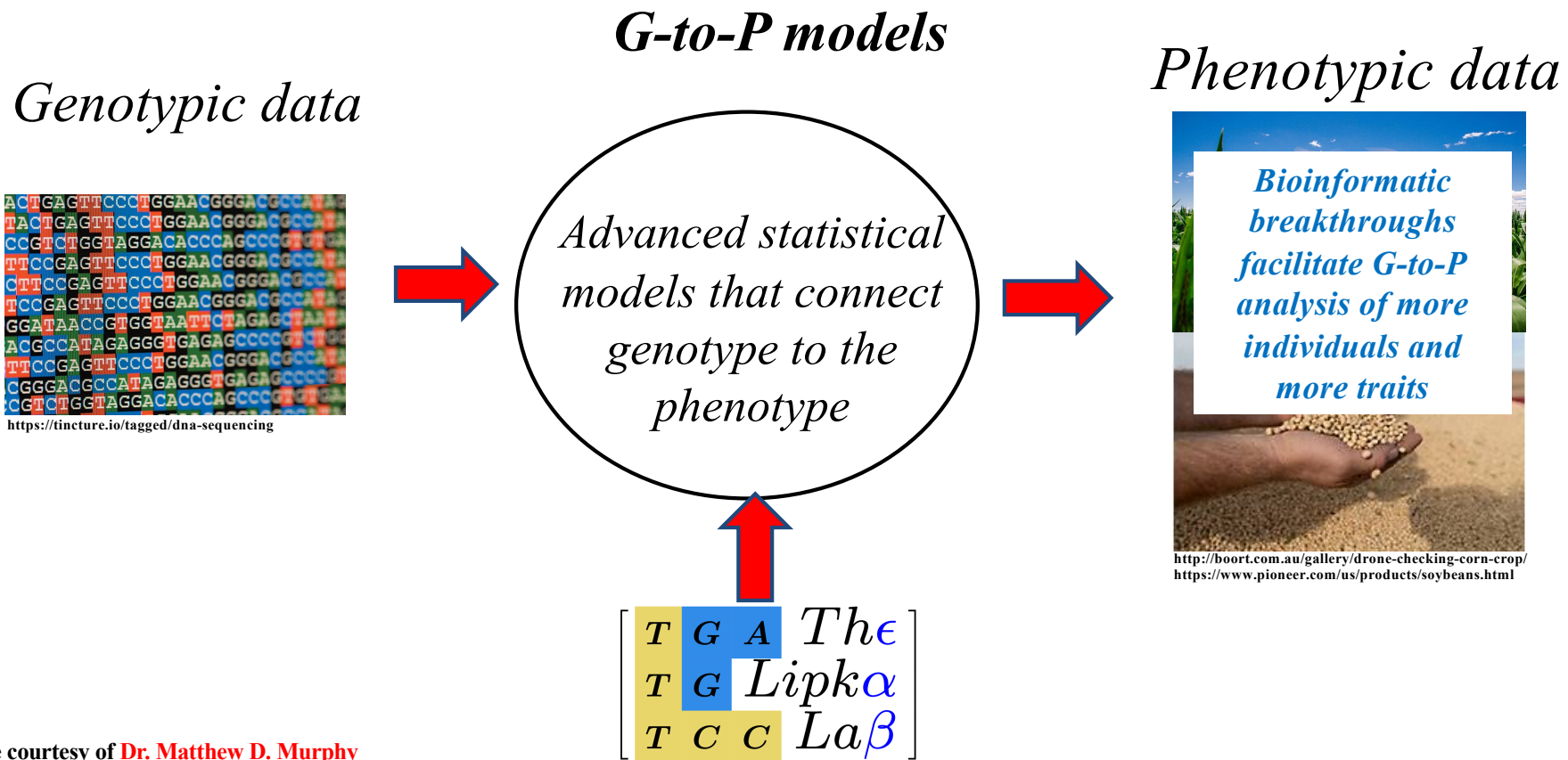
Alexander E. Lipka

Associate Professor,
Department of Crop Sciences,
University of Illinois, USA

Relationship between my research and bioinformatic breakthroughs



Dr. Matthew D. Murphy



Slide courtesy of Dr. Matthew D. Murphy

Bioinformatic breakthroughs facilitate analysis of more sophisticated phenotypes

Fly drones



<http://boort.com.au/gallery/drone-checking-corn-crop/>

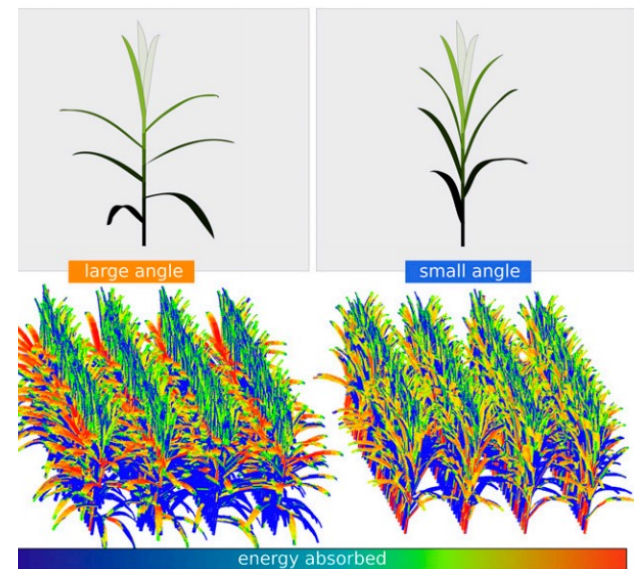
<https://www.pioneer.com/us/products/soybeans.html>

Take images



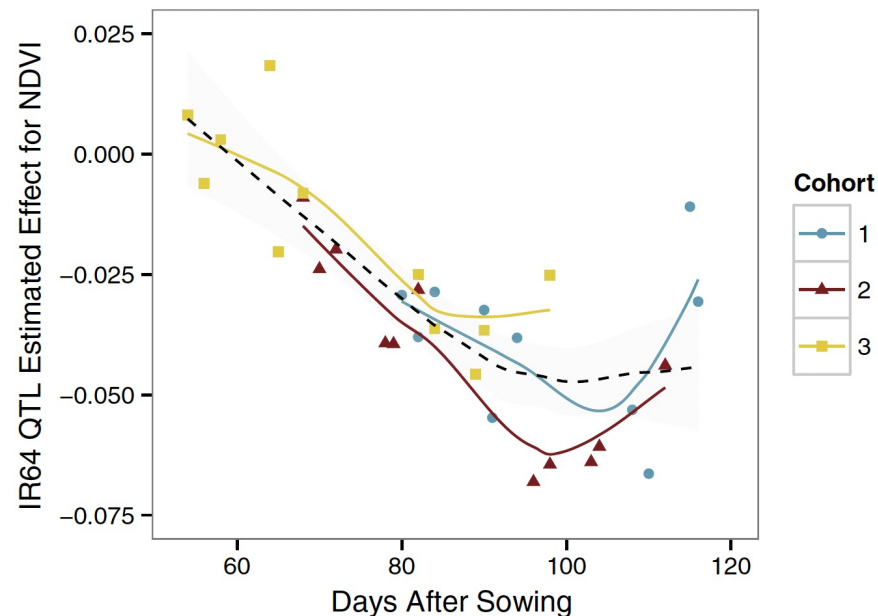
Lipka Lab growout from 2021 UIUC field season

Measure phenotypes



Truong et al., Genetics (2015) Genetics

Bioinformatic breakthroughs make it possible to see how G-to-P relationships change across lifespan



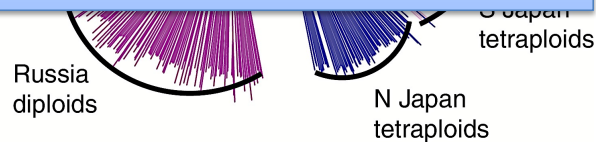
Family of $n = 1,741$ rice RILs

- Trait = Reflectance ratio (NVDI)
- Obtained using multispectral sensors on a tractor

High quality bioinformatic data helped make exciting research possible in plants

Miscanthus sinensis diversity panel: $n = 538$

Assess the utility of genomic prediction across multiple environments

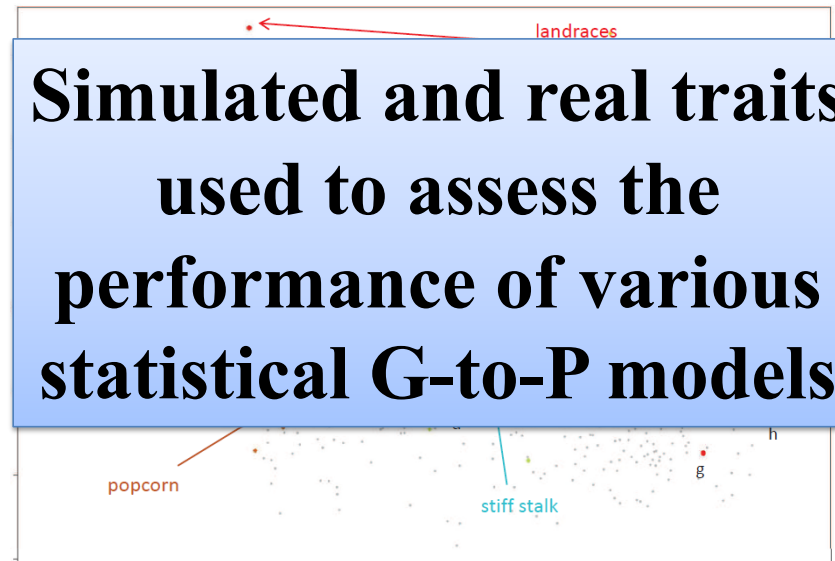


Clark *et al.* (2018)

NCRPIS ("Ames") maize diversity panel: $n = 2,815$

Simulated and real traits used to assess the performance of various statistical G-to-P models

Principal Coordinate 2



Principal Coordinate 1

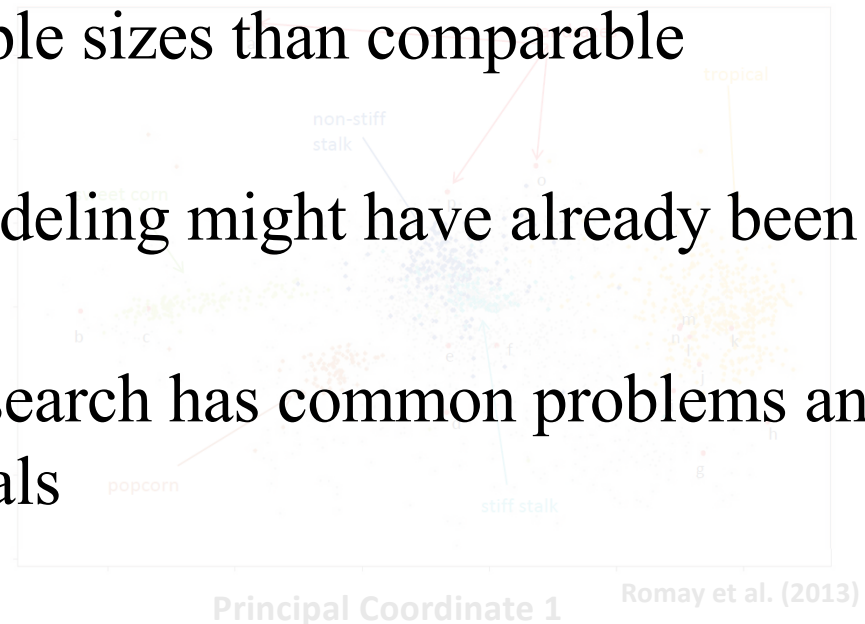
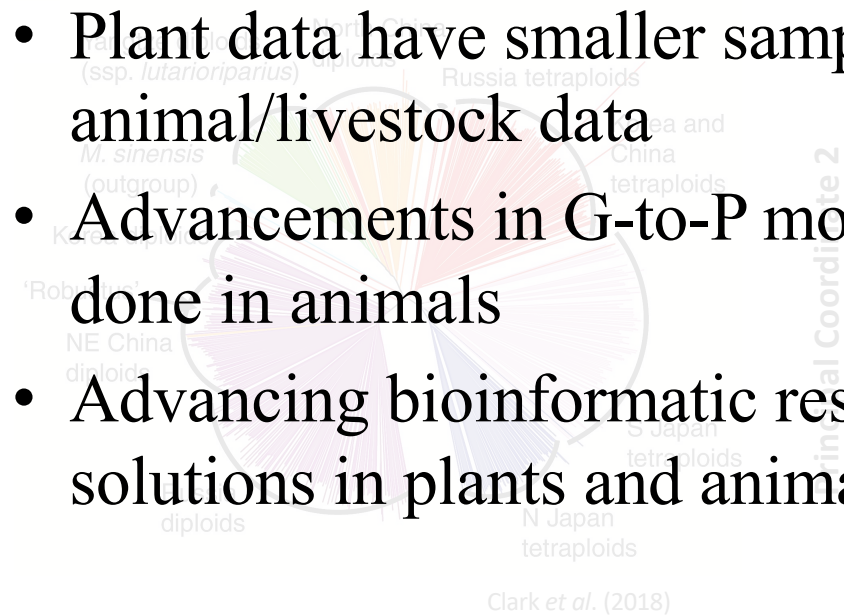
Romay *et al.* (2013)

AG2PI has put my lab's research experience into a broader context

Miscanthus sinensis diversity panel: $n = 538$

NCRPIS ("Ames") maize diversity panel: $n = 2,815$

- Plant data have smaller sample sizes than comparable animal/livestock data
- Advancements in G-to-P modeling might have already been done in animals
- Advancing bioinformatic research has common problems and solutions in plants and animals



AGP2I Thinking Big: Advancing Genomics Research

Current Challenges and Future of Agricultural Genomes to Phenomes in the U.S.

Authors:

Christopher K. Tuggle^{1*#}, Jennifer L. Clarke^{2#}, Brenda M. Murdoch^{3#}, Eric Lyons^{4#}, Nicole M. Scott^{1#}, Bedrich Beneš⁵, Jacqueline D. Campbell⁶, Sruti Das Choudhury², Henri Chung¹, Courtney L. Daigle⁷, Jack C. M. Dekkers¹, Joao R. R. Dórea⁸, David S. Ertl⁹, Max Feldman¹⁰, Breno O. Fragomeni¹¹, Janet E. Fulton¹², Carmela R. Guadagno¹³, Darren E. Hagen¹⁴, Andrew S. Hess¹⁵, Luke M. Kramer¹, Carolyn J. Lawrence Dill¹, Alexander E. Lipka¹⁶, Thomas Lübberstedt¹, Fiona M. McCarthy⁴, Stephanie D. McKay⁷, Seth C. Murray⁷, Penny K. Riggs⁷, Troy N. Rowan¹⁸, Moira J. Sheehan¹⁹, Juan P. Steibel¹, Addie M. Thompson²⁰, Kara J. Thornton²¹, Curtis P. Van Tassell²², Patrick S. Schnable^{1*}

More genomes need to be available

- **More individuals need to be sequenced**
 - Overreliance on a single reference genome

Sequencing at both genomic and epigenomic levels are needed

- **What this will facilitate**
 - Studying genetic diversity of relevant breeding material
 - Studying structural variation
 - Study history of genetic architecture

Basic science needs to be translated to applications

- **Basic genomic research has advanced**
 - Better understanding of non-additive effects
 - Better understanding of GxE
- **Improving analytical tools can expediate applications**
 - Practical
 - Usable
 - Understandable to users in multiple disciplines

A cohesive agricultural genomics community is needed

- **Support needed in the following areas**

- Scientific
- Funding
- Human resources

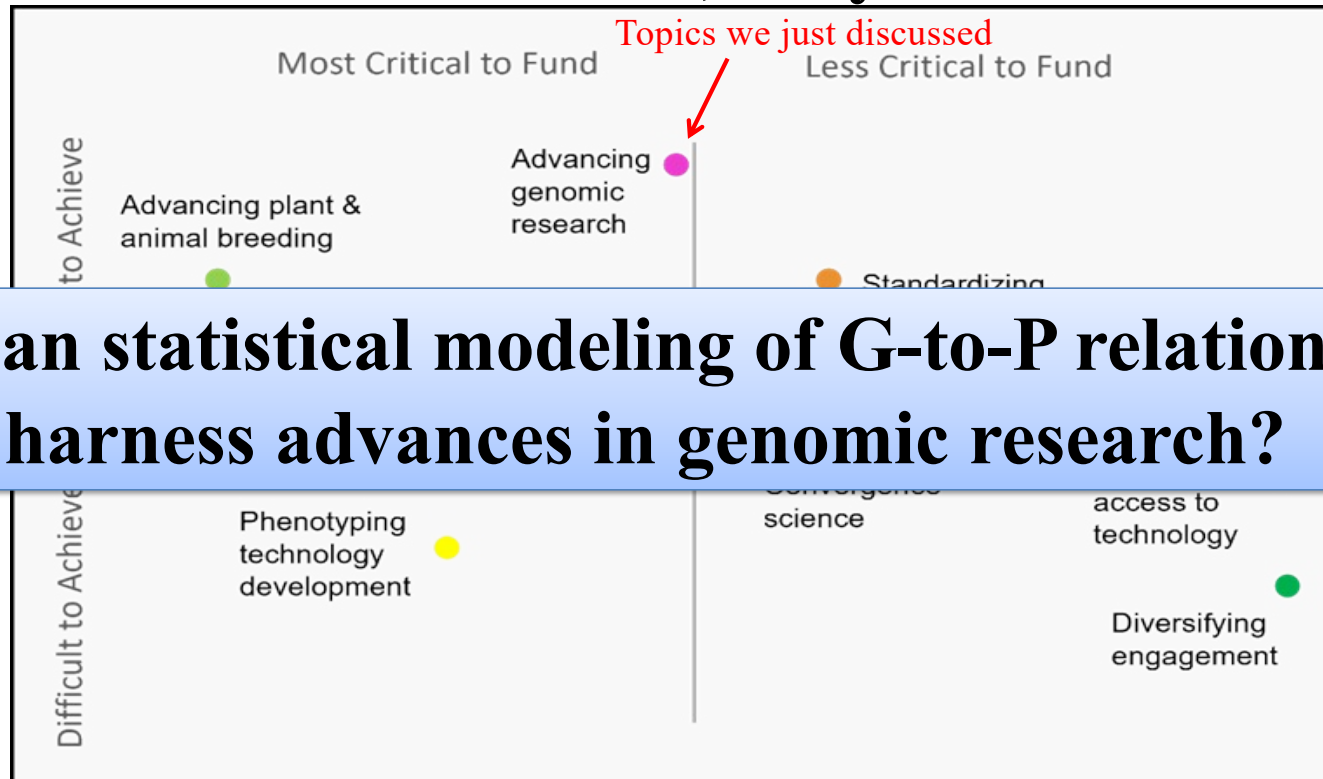
- **Benefits**

- Scientists can focus on science
- Innovative scientific thinking is encouraged
- Latest approaches can be used

Generate high-resolution and more diverse -omics data are needed

- **Species agnostic**
- **Serve as reference data sets**
 - Decrease data collection costs
- **Engaging scientists in industry ensures:**
 - Relevance
 - Accessibility to diverse stakeholders

Advancing genomic research: Critical to fund, easy to achieve



How can statistical modeling of G-to-P relationships harness advances in genomic research?

More multi-trait analyses are needed

$$\begin{array}{c} \text{Vector of} \\ \text{multiple traits} \end{array} \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ \cdot \\ Y_k \end{pmatrix} = \begin{array}{c} \text{Grand mean} \\ \downarrow \\ \mathbf{1}\boldsymbol{\mu} \end{array} + \begin{array}{c} \text{Fixed effects:} \\ \text{account for} \\ \text{population} \\ \text{structure} \end{array} \begin{array}{c} \text{Marker effect} \\ \downarrow \\ \mathbf{Q}\boldsymbol{\beta} \end{array} + \mathbf{X}\boldsymbol{\alpha} + \begin{array}{c} \text{Random effects:} \\ \text{account for familial} \\ \text{relatedness} \end{array} \begin{array}{c} \mathbf{Z}\mathbf{u} \end{array} + \begin{array}{c} \text{Random error vector} \\ \leftarrow \\ \boldsymbol{\varepsilon} \end{array}$$

\mathbf{Q} , \mathbf{X} , and \mathbf{Z} are design matrices

$$\mathbf{u} \sim \text{MVN}(\mathbf{0}, 2K\sigma_G^2)$$

K = kinship matrix

$$\boldsymbol{\varepsilon} \sim \text{MVN}(\mathbf{0}, I\sigma_E^2)$$

Multiple loci need to be considered in one model (MSTEP shown here)

$$Y = XB + E$$

Diagram illustrating the MSTEP model equation: $Y = XB + E$.

- Y : Vector of multiple traits (indicated by a red arrow).
- X : Design matrix including multiple markers (indicated by a blue circle around X).
- B : Matrix of additive effects for each marker and each trait (indicated by a blue circle around B).
- E : Random matrix of residuals (indicated by a red arrow).

- Determining the optimal model:
 - AIC, BIC, mBIC
 - Permutation procedure

Non-additive effects need to be modeled

Grand Mean

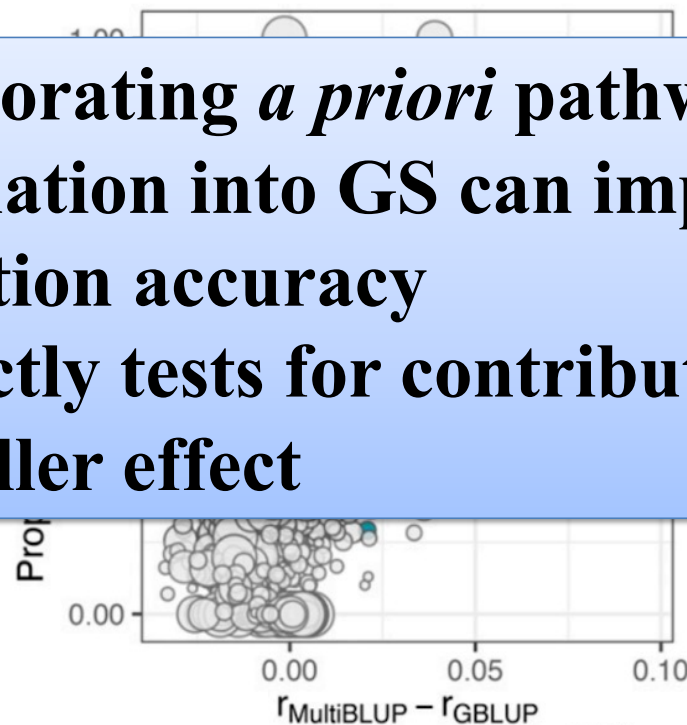
Random error term

**Stepwise epistatic model selection =
Stepwise Procedure for constructing an
Additive and Epistatic Multi-Locus model
(SPAEMML)**

- I is a subset of markers with additive effects in model
- U is a subset of markers with two-way epistatic effects in model
- Determining the optimal model:
 - AIC, BIC, mBIC
 - Permutation procedure

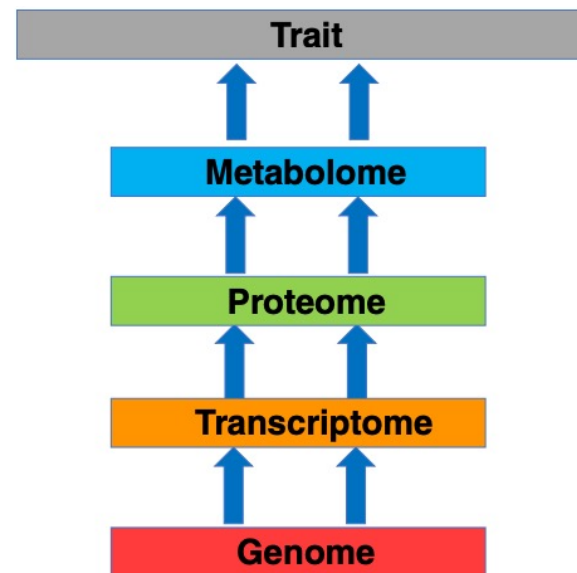
Contribution of non-statistically significant loci needs to be quantified

- Incorporating *a priori* pathway information into GS can improve prediction accuracy
- Indirectly tests for contributions of genes of smaller effect



Multi-kernel performs better than standard

**-omic levels connecting intermediate steps between
“G” and “P” need to be included in the model**



$$Y = \mu \mathbf{1} + Z \mathbf{u}_G + Z \mathbf{u}_T + Z \mathbf{u}_P + Z \mathbf{u}_M + \varepsilon$$

$$\mathbf{u}_G \sim MVN(\mathbf{0}, G\sigma_g) \quad \mathbf{u}_T \sim MVN(\mathbf{0}, T\sigma_T)$$

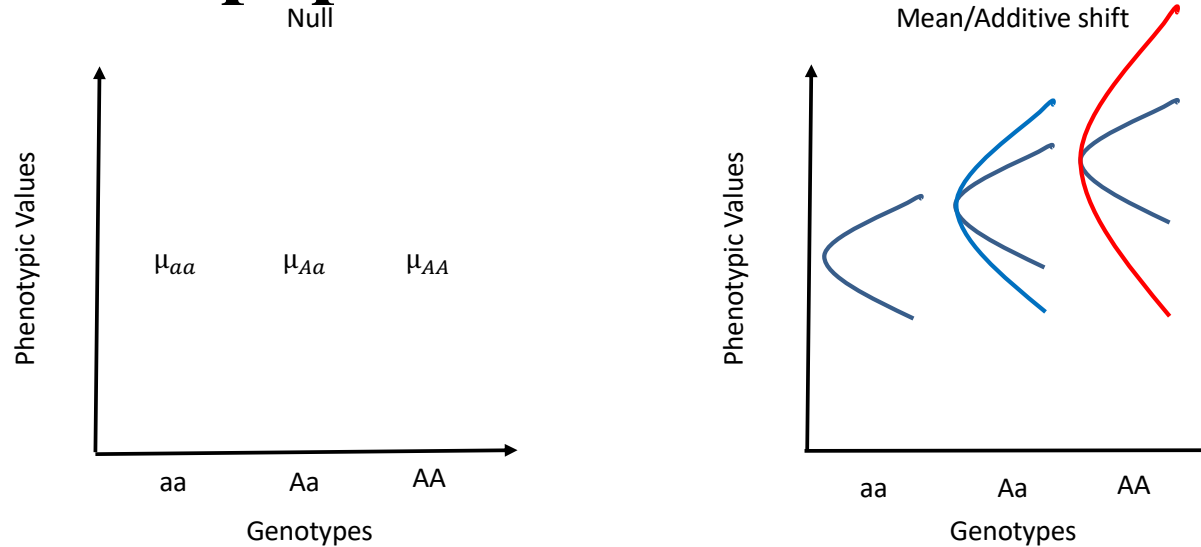
$$\mathbf{u}_P \sim MVN(\mathbf{0}, P\sigma_P) \quad \mathbf{u}_M \sim MVN(\mathbf{0}, M\sigma_m)$$

$$\varepsilon \sim MVN(\mathbf{0}, I\sigma_\varepsilon^2)$$



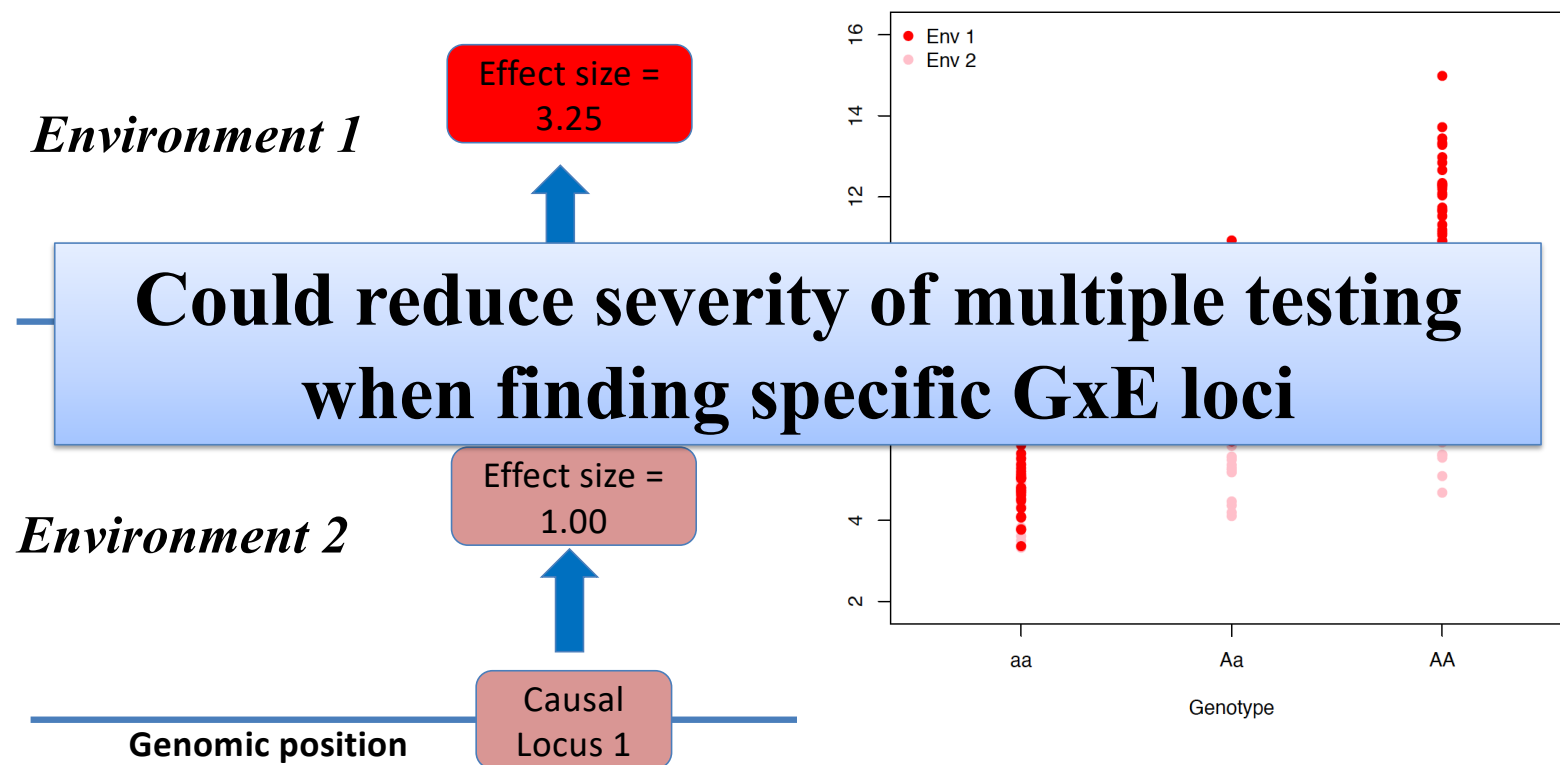
Dr. Matthew D.
Murphy

Move beyond testing for differences in population mean trait values



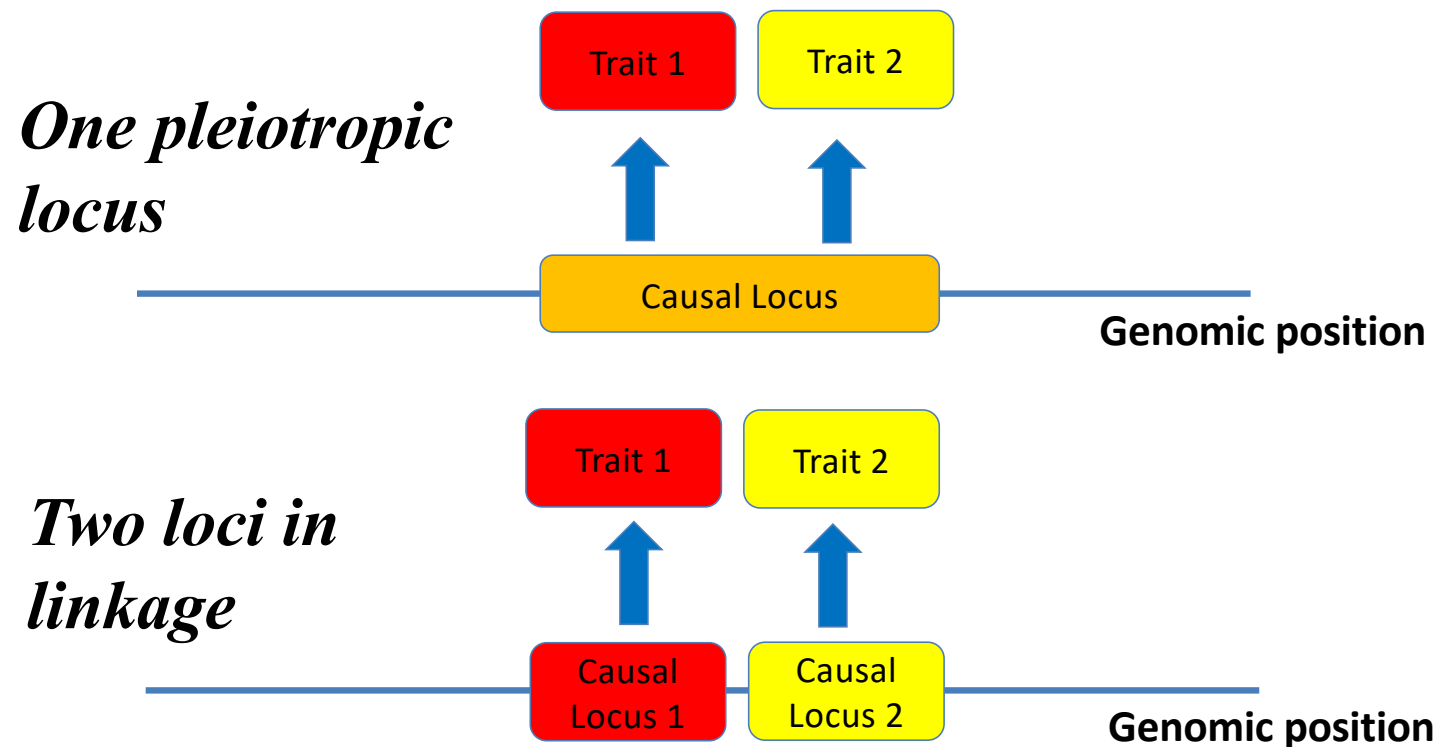
Variance Quantitative Trait Loci (vQTL)
Variance GWAS (vGWAS)

GxE interactions could appear as a vQTL



Bioinformatic breakthroughs can assist with testing scientific hypothesis on genetic architecture

Simple example: Pleiotropy versus linkage?



Bioinformatic breakthroughs can assist with testing scientific hypothesis on genetic architecture

Complex example: Evidence for omnigenic and/or other genetic architectures?

