

Genomic Privacy

- Personalized genomic data generation is booming
- Main focus is on protecting variants
- “Detection of genome in a mixture”
 - Individuals give consent to participate but request anonymity
 - HAPMAP, Personal genome project, 1000 Genomes...
- Larger and more datasets leads to more realistic risks of linking attacks, that may be much more damaging than detection of genome in a mixture attacks

Identifying Participants in the Personal Genome Project by Name

Latanya Sweeney, Akua Abu, Julia Winn

Harvard College
Cambridge, Massachusetts

latanya@fas.harvard.edu, aabu@college.harvard.edu, jwinn@post.harvard.edu

Identifying Personal Genomes by Surname Inference

Melissa Gymrek,^{1,2,3,4} Amy L. McGuire,⁵ David Golan,⁶ Eran Halperin,^{7,8,9} Yaniv Erlich^{1*}

Resolving Individuals Contributing Trace Amounts of DNA to Highly Complex Mixtures Using High-Density SNP Genotyping Microarrays

Nils Homer^{1,2}, Szabolcs Szelinger¹, Margot Redman¹, David Duggan¹, Waibhav Tembe¹, Jill Muehling¹, John V. Pearson¹, Dietrich A. Stephan¹, Stanley F. Nelson², David W. Craig^{1*}

On Sharing Quantitative Trait GWAS Results in an Era of Multiple-omics Data and the Limits of Genomic Privacy

Hae Kyung Im,^{1,*} Eric R. Gamazon,² Dan L. Nicolae,^{2,3,4} and Nancy J. Cox^{2,3,*}

Robust De-anonymization of Large Sparse Datasets

Arvind Narayanan and Vitaly Shmatikov
The University of Texas at Austin



Linking Attack Scenario

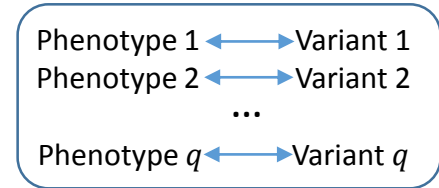
Phenotype dataset
(Public)

Phenotype ID	HIV Status	Phenotype 1	Phenotype 2	Phenotype q	
PID-1	HIV+	0.1	-2.7	...	90.3
PID-2	HIV-	0.5	8.6	...	63.5
⋮	⋮	⋮	⋮	⋮	⋮
PID- n	HIV-	-0.2	5.4	...	50.3

Genotype dataset
(Stolen/Hacked/Queried) 1

Genotype ID	Variant 1	Variant 2	Variant q	
GID-1	0	1	...	1
GID-2	2	1	...	0
⋮	⋮	⋮	⋮	⋮
GID- m	1	2	...	1

Phenotype-Genotype correlation dataset



3

Genotype prediction

Phenotype ID	HIV Status	Predicted variant genotypes			
		Variant 1	Variant 2	Variant q	
PID-1	HIV+	1	0	...	2
PID-2	HIV-	2	2	...	1
⋮	⋮	⋮	⋮	⋮	⋮
PID- n	HIV-	0	1	...	1

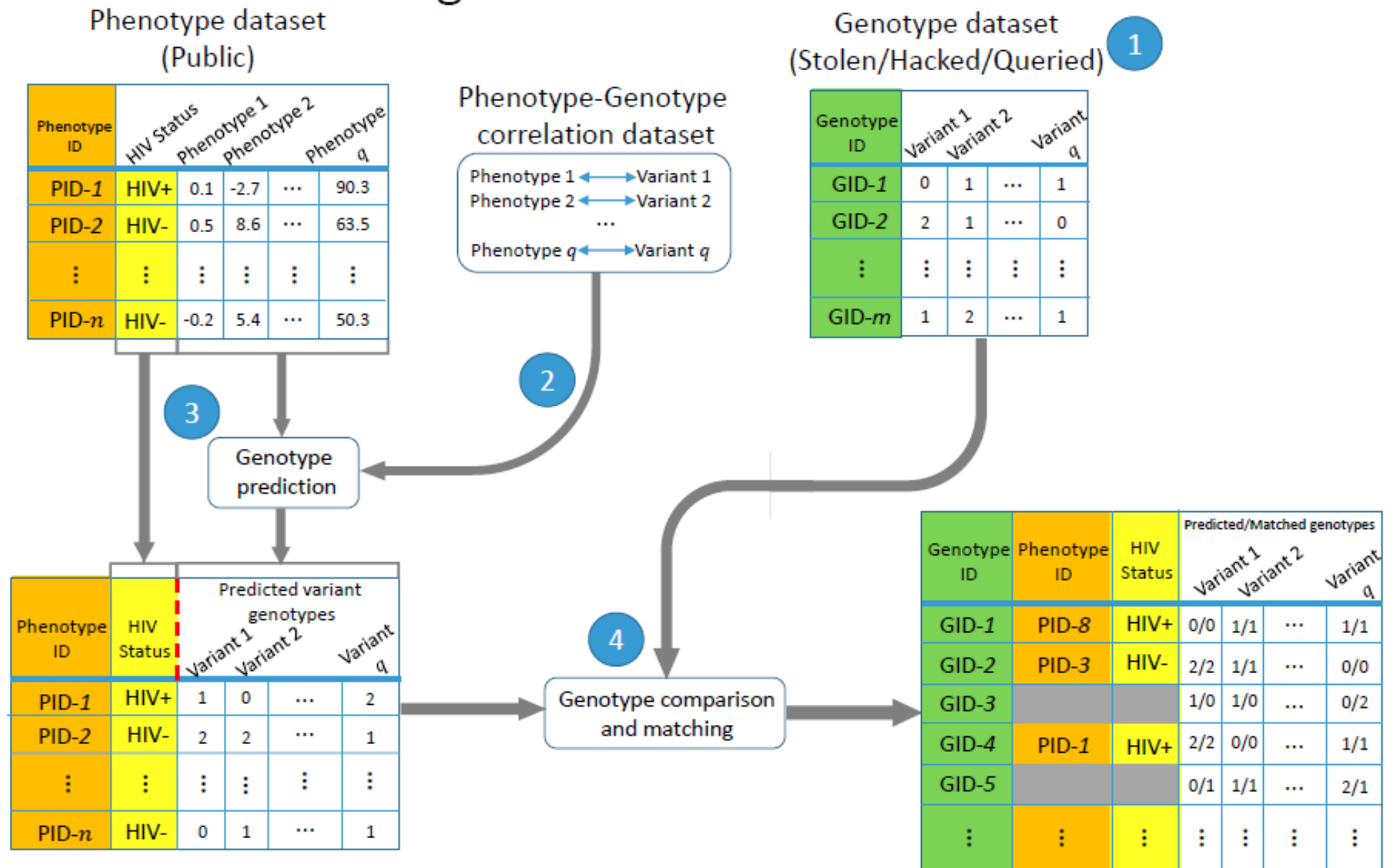
2

4

Genotype comparison and matching

Genotype ID	Phenotype ID	HIV Status	Predicted/Matched genotypes			
			Variant 1	Variant 2	Variant q	
GID-1	PID-8	HIV+	0/0	1/1	...	1/1
GID-2	PID-3	HIV-	2/2	1/1	...	0/0
GID-3			1/0	1/0	...	0/2
GID-4	PID-1	HIV+	2/2	0/0	...	1/1
GID-5			0/1	1/1	...	2/1
⋮	⋮	⋮	⋮	⋮	⋮	⋮

Linking Attack Scenario

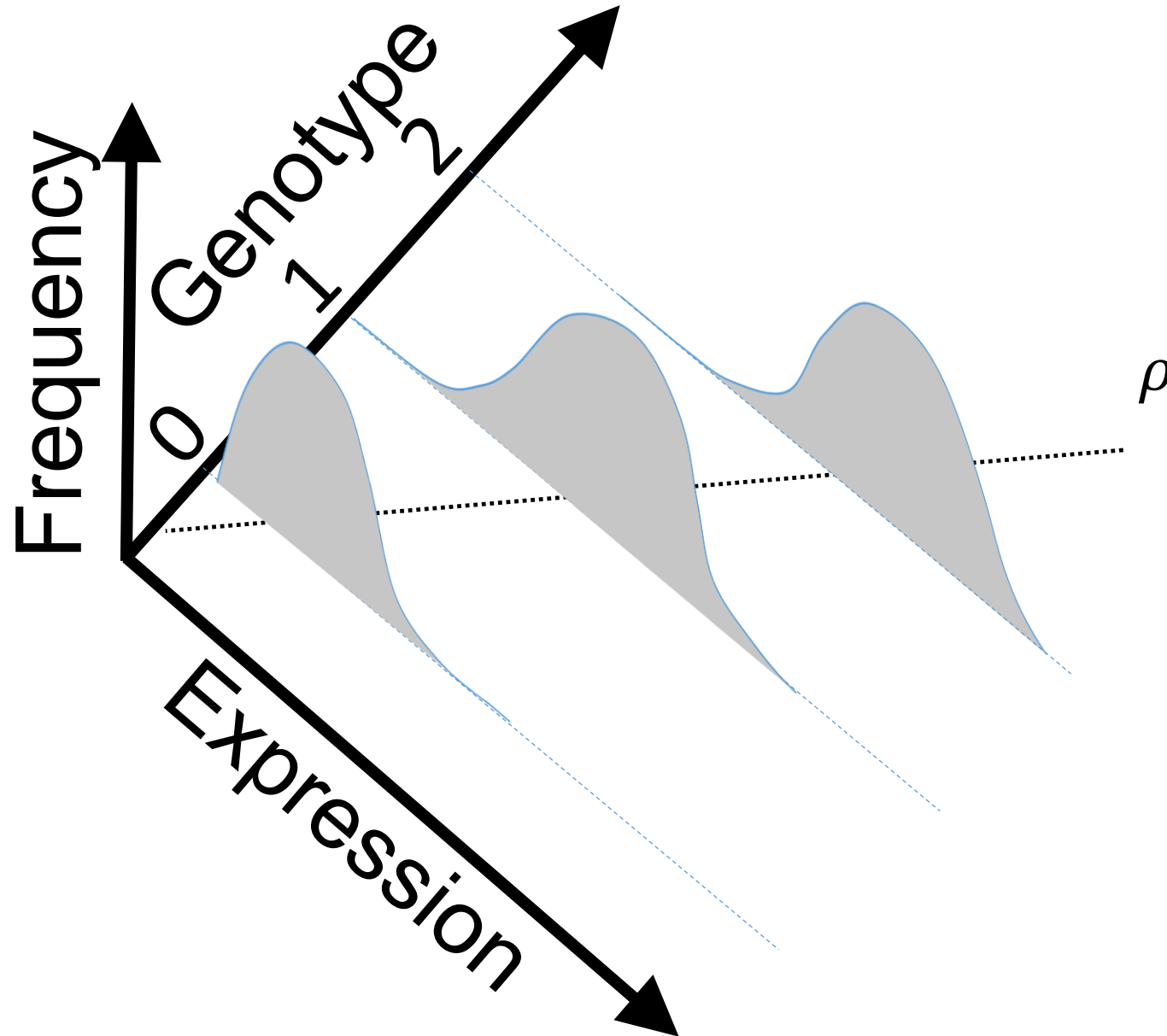


Representative Expression, Genotype, eQTL Datasets

- mRNA sequencing for 462 individuals
 - Publicly available Quantification for protein coding genes
- Approximately 3,000 cis-eQTL (FDR<0.05)
- Genotypes are available from the 1000 Genomes Project



Expression and Genotype Distribution

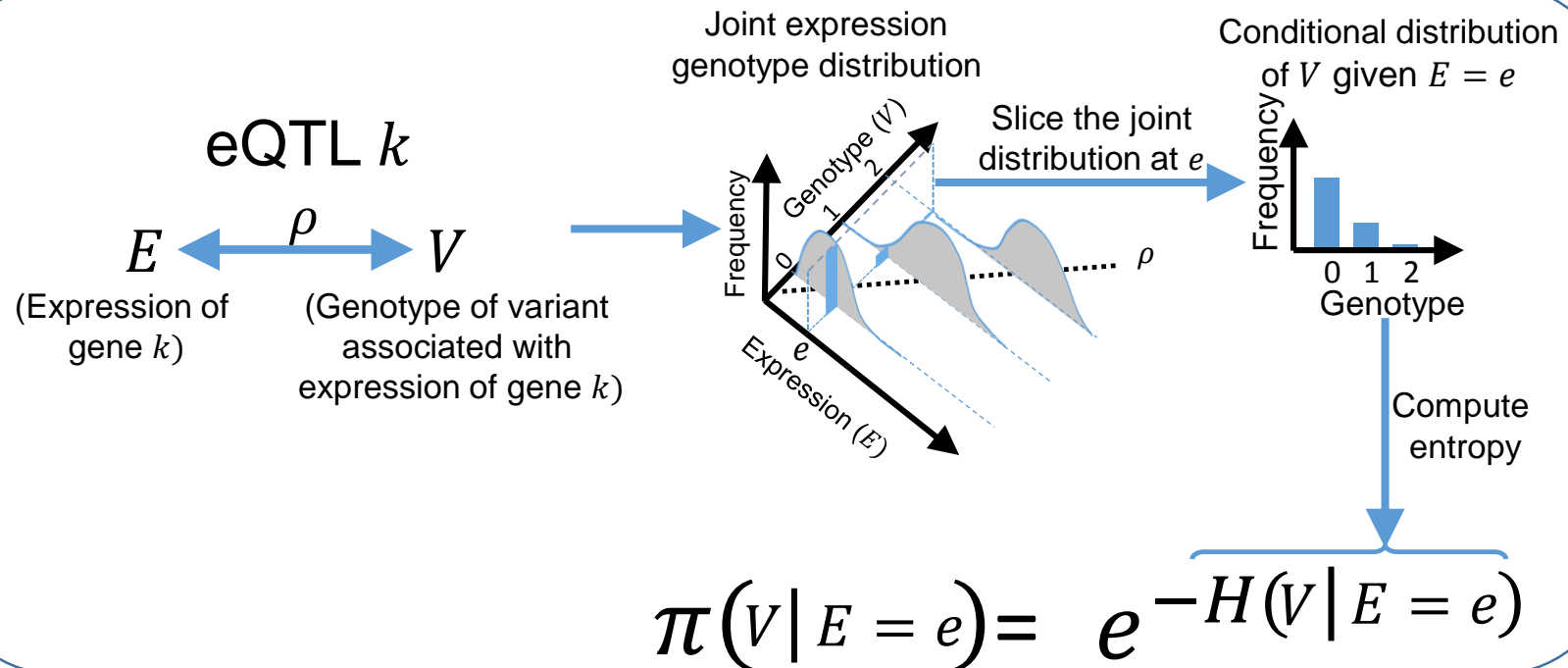


Information Content and Predictability

$$ICI \left(\begin{array}{l} \text{Individual has variant} \\ \text{genotypes } g_1, g_2, \dots, g_n \\ \text{for variants } V_1, V_1, \dots, V_n \end{array} \right) = \log \left(\frac{1}{\text{Frequency of } V_1 \text{ genotype}} \right) + \log \left(\frac{1}{\text{Frequency of } V_2 \text{ genotype}} \right) + \dots + \log \left(\frac{1}{\text{Frequency of } V_n \text{ genotype}} \right)$$

$g_1 = 2$ $g_2 = 1$ $g_n = 2$

V_1 genotype frequencies V_2 genotype frequencies V_n genotype frequencies

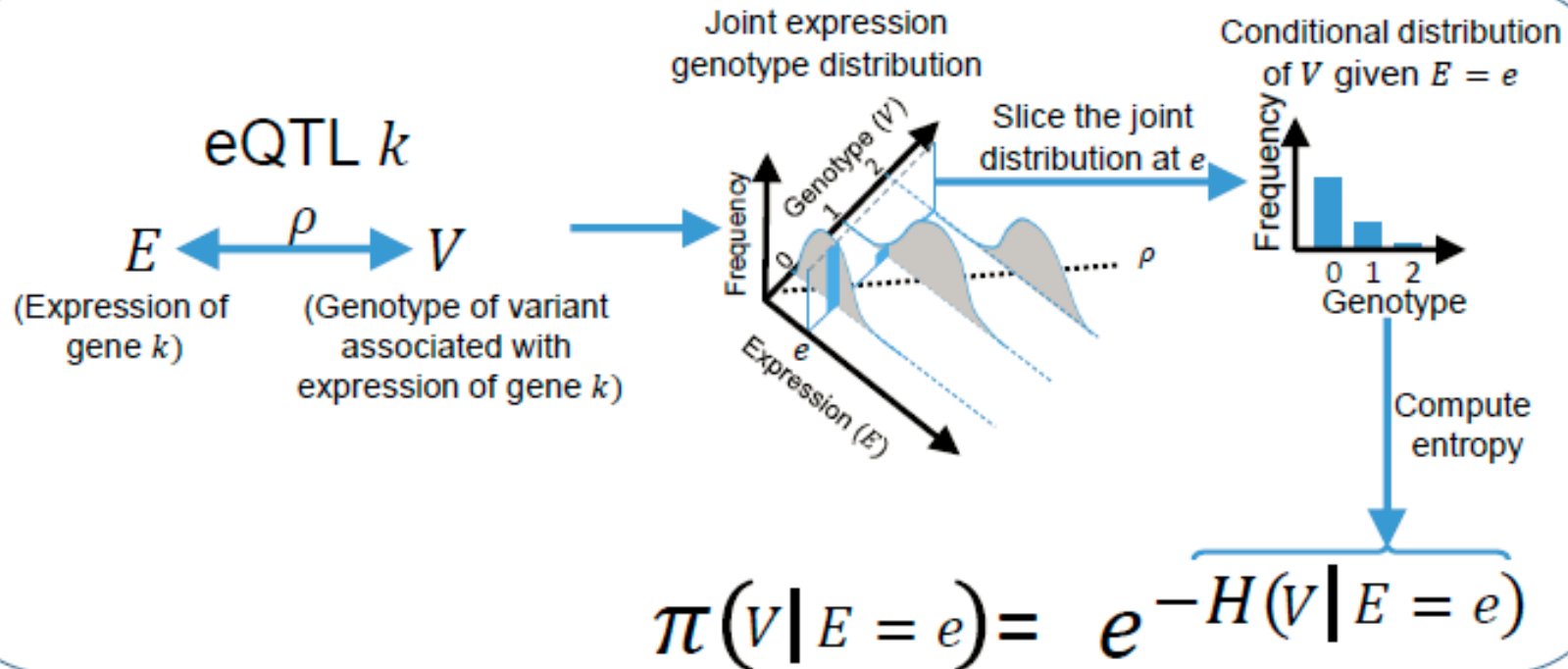


Information Content and Predictability

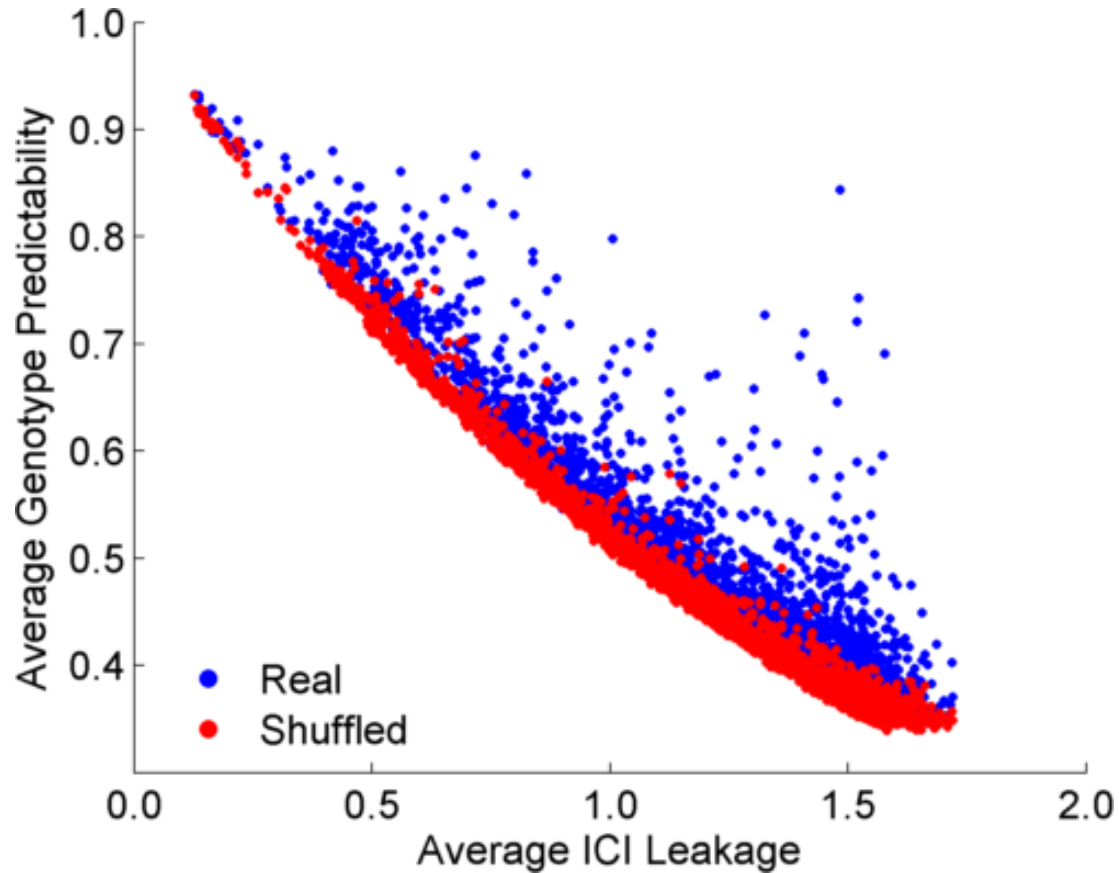
$$ICI \left(\begin{array}{l} \text{Individual has variant} \\ \text{genotypes } g_1, g_2, \dots, g_n \\ \text{for variants } V_1, V_1, \dots, V_n \end{array} \right) = \log \left(\frac{1}{\text{Frequency of } V_1 \text{ genotype}} \right) + \log \left(\frac{1}{\text{Frequency of } V_2 \text{ genotype}} \right) + \dots + \log \left(\frac{1}{\text{Frequency of } V_n \text{ genotype}} \right)$$

$g_1 = 2$ $g_2 = 1$ $g_n = 2$

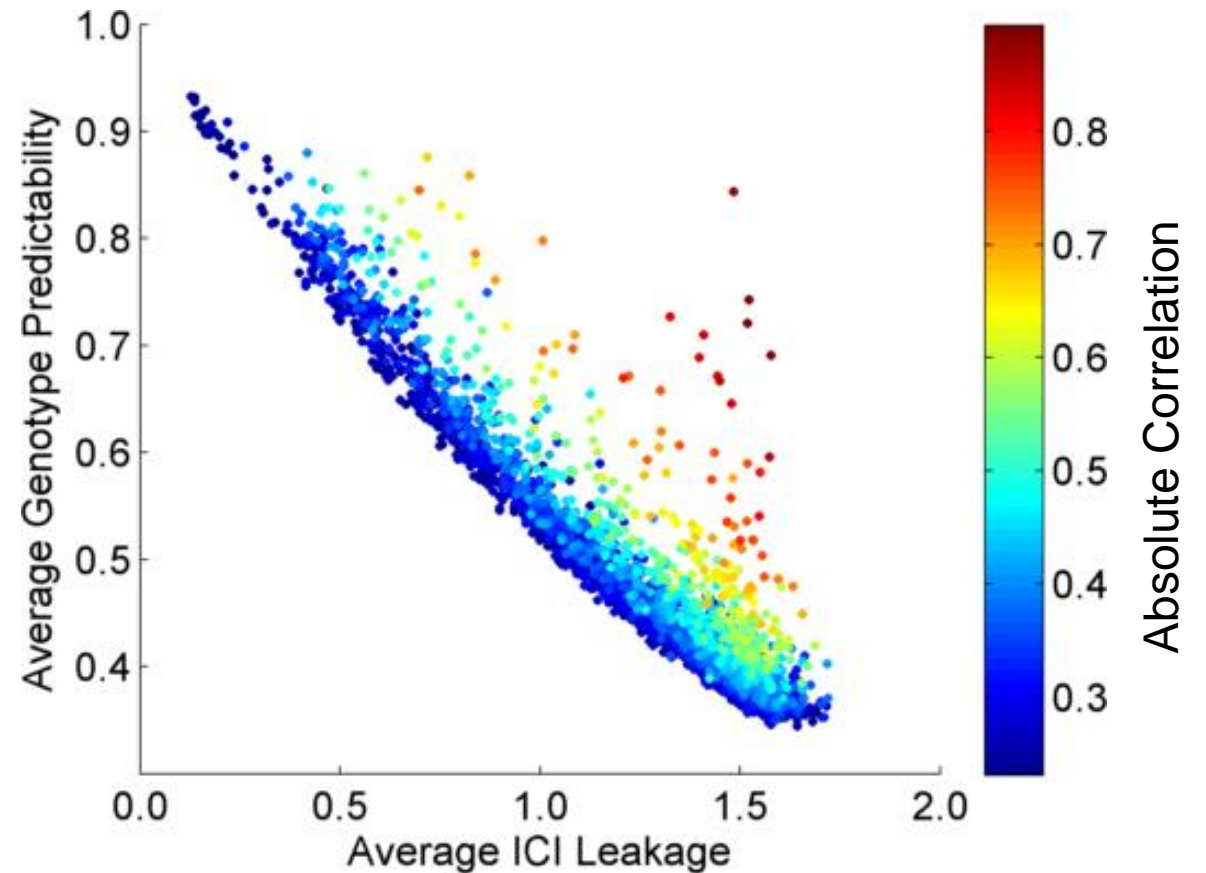
V_1 genotype frequencies V_2 genotype frequencies V_n genotype frequencies



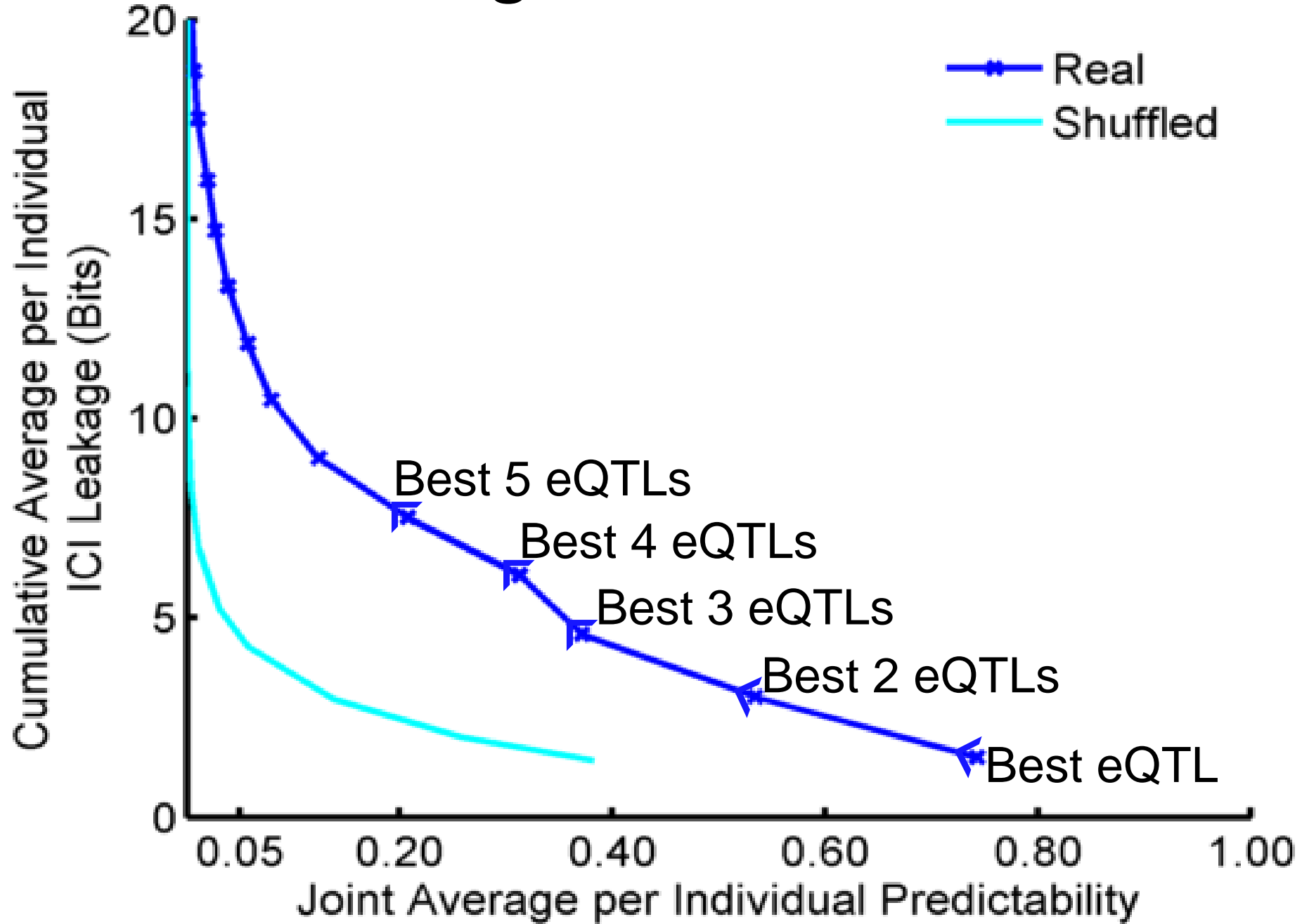
Per eQTL and ICI Cumulative Leakage versus Genotype Predictability



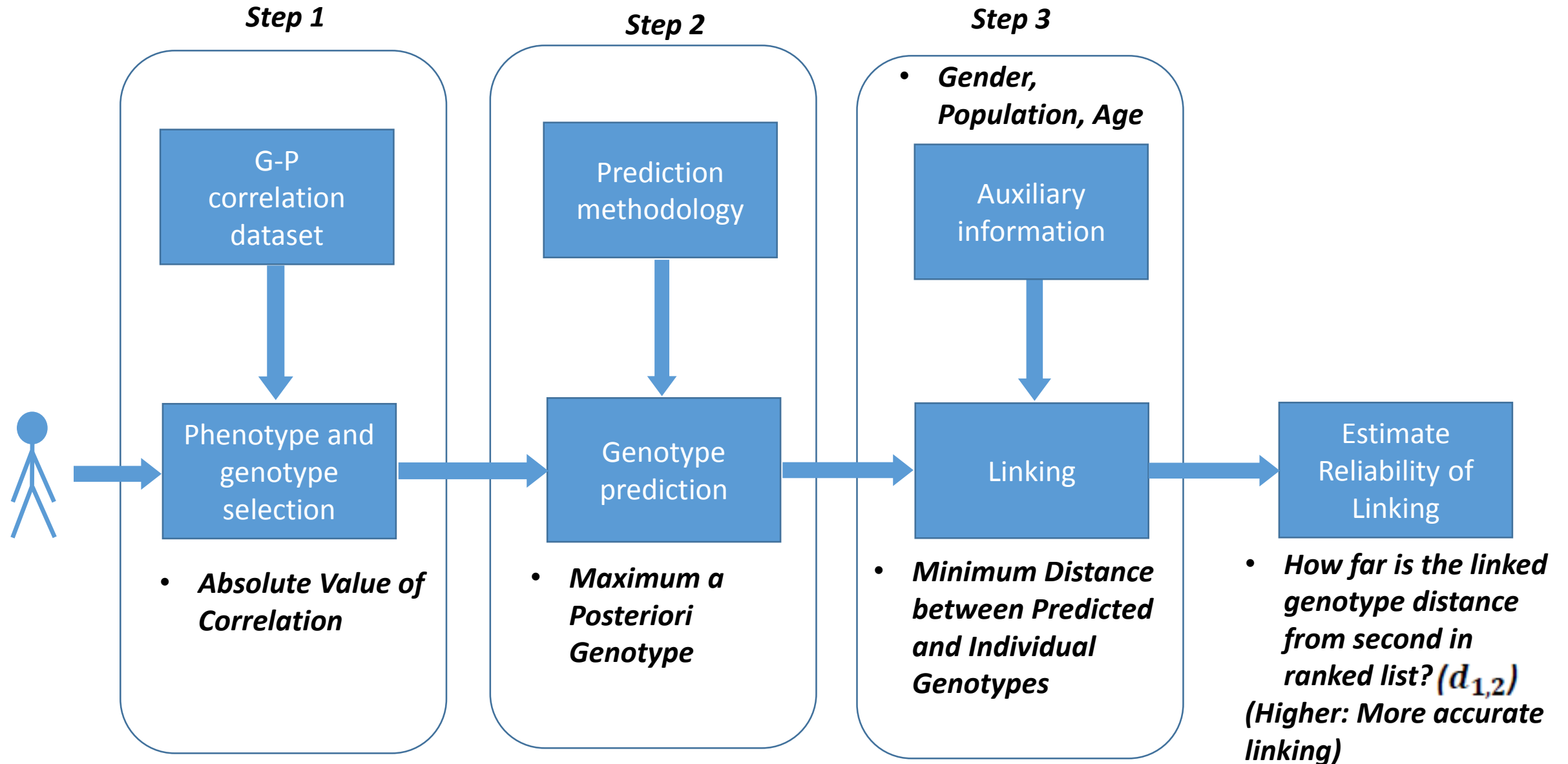
Colors by absolute correlation

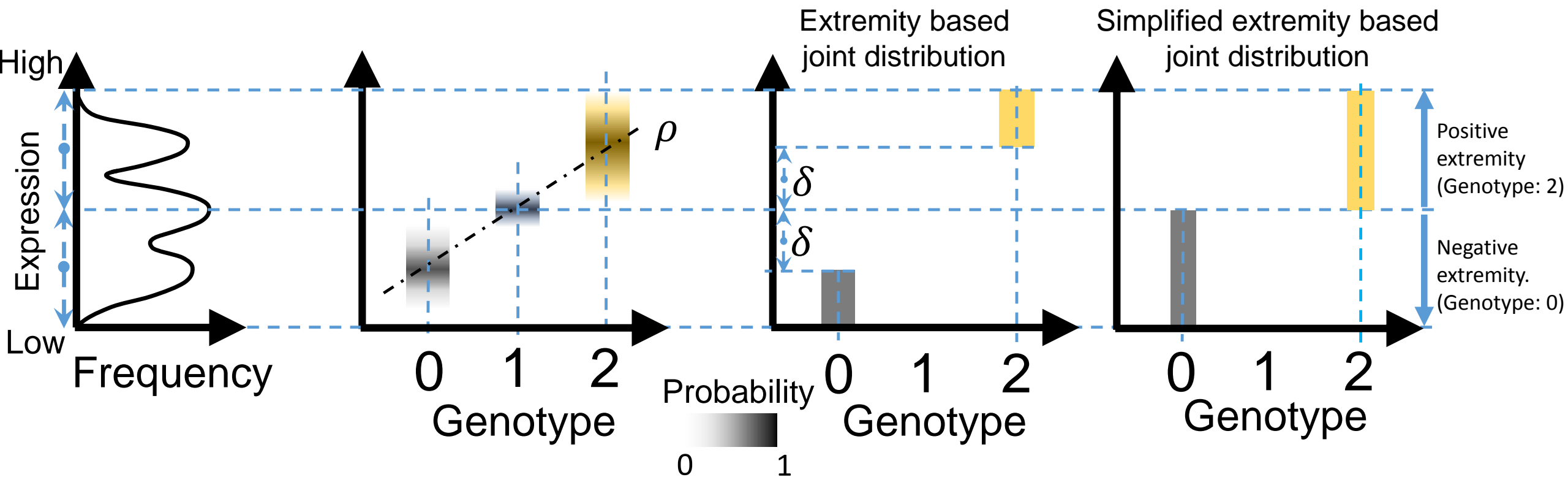


Cumulative Leakage versus Joint Predictability

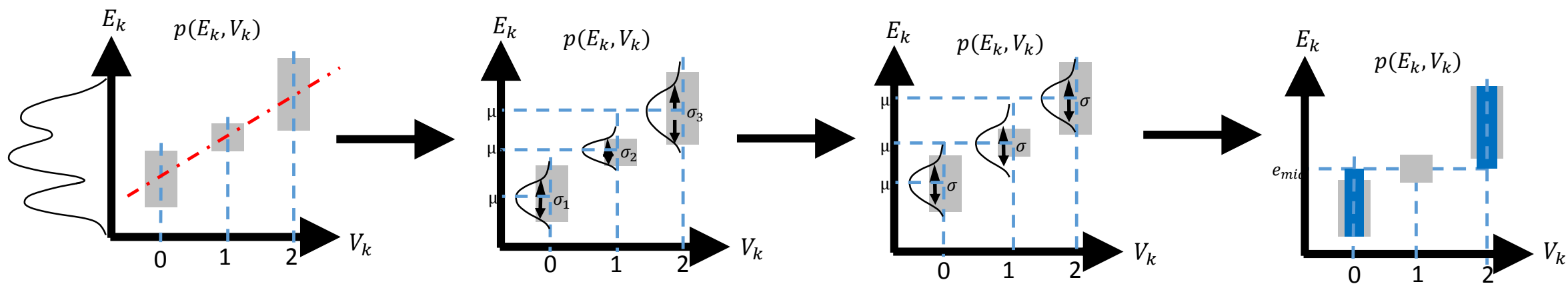


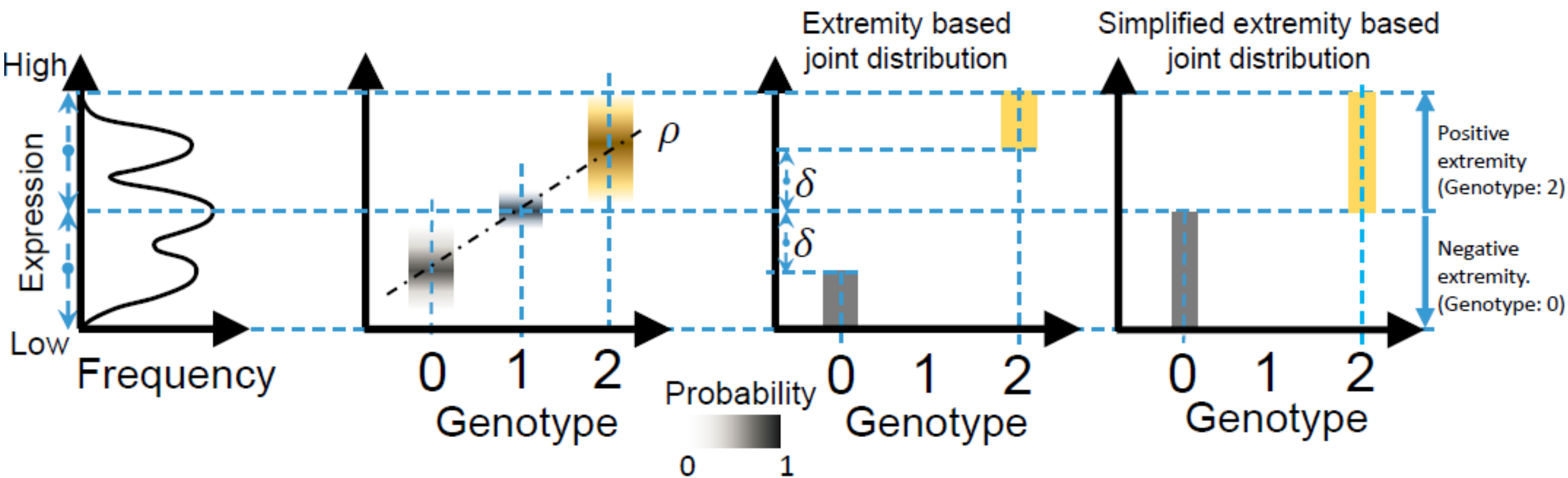
Steps in Instantiation of a (Mock) Linking Attack



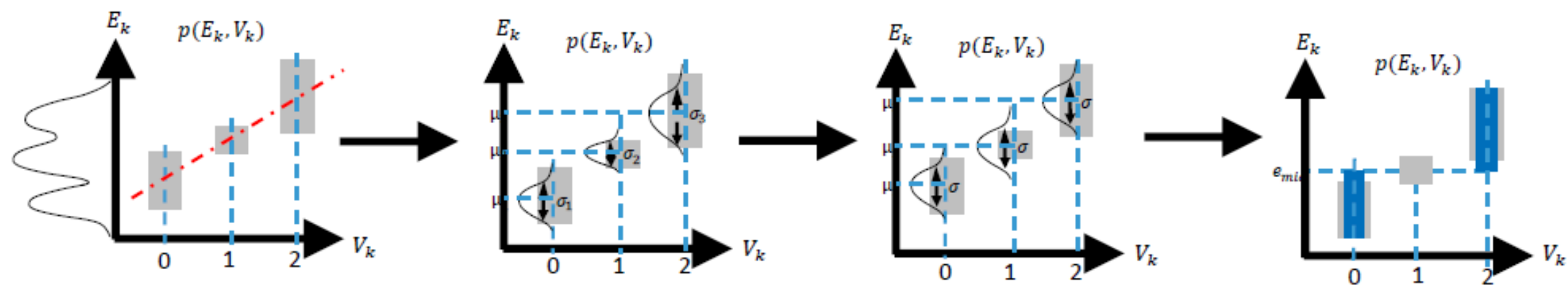


Levels of Expression-Genotype Model Simplifications:

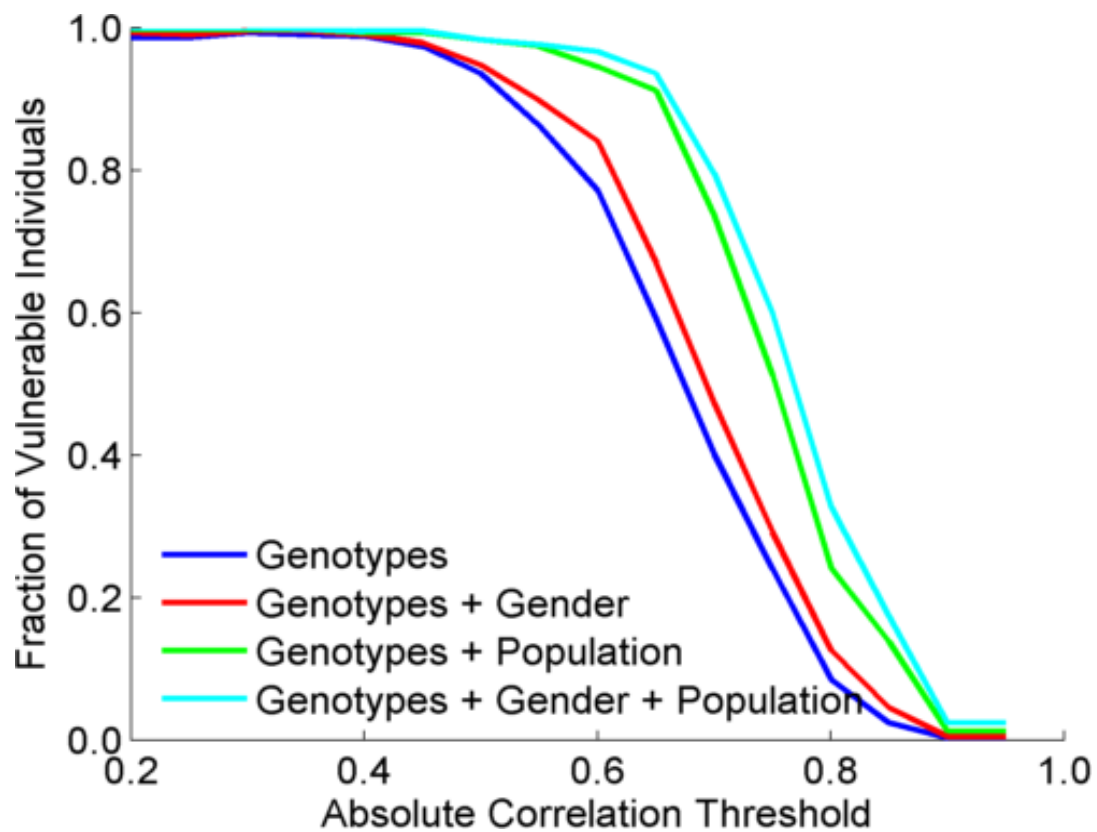




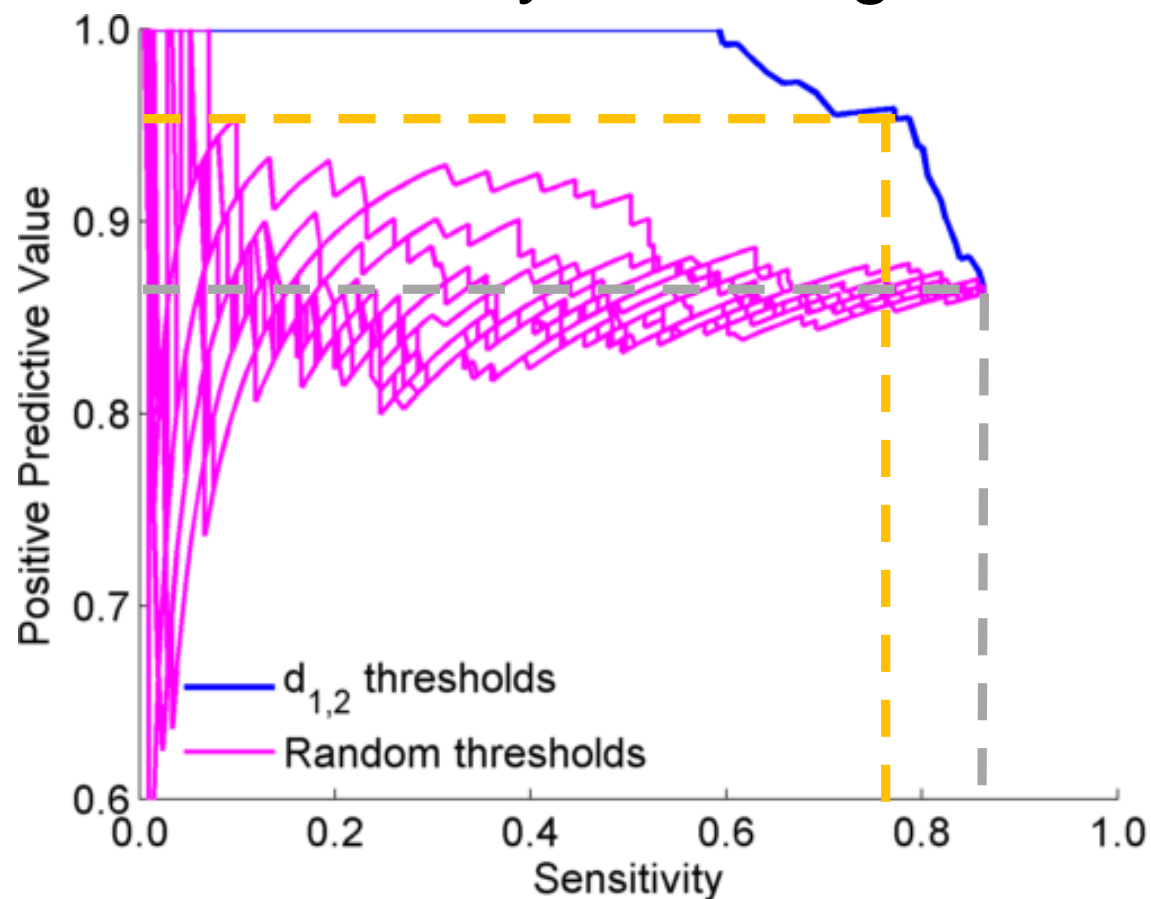
Levels of Expression-Genotype Model Simplifications:



Extremity based linking with homozygous genotypes



Attacker can estimate the reliability of linkings



Sensitivity: Fraction of individuals that are correctly linked PPV: Fraction of selected individuals that are correctly linked

Risk Management Framework

