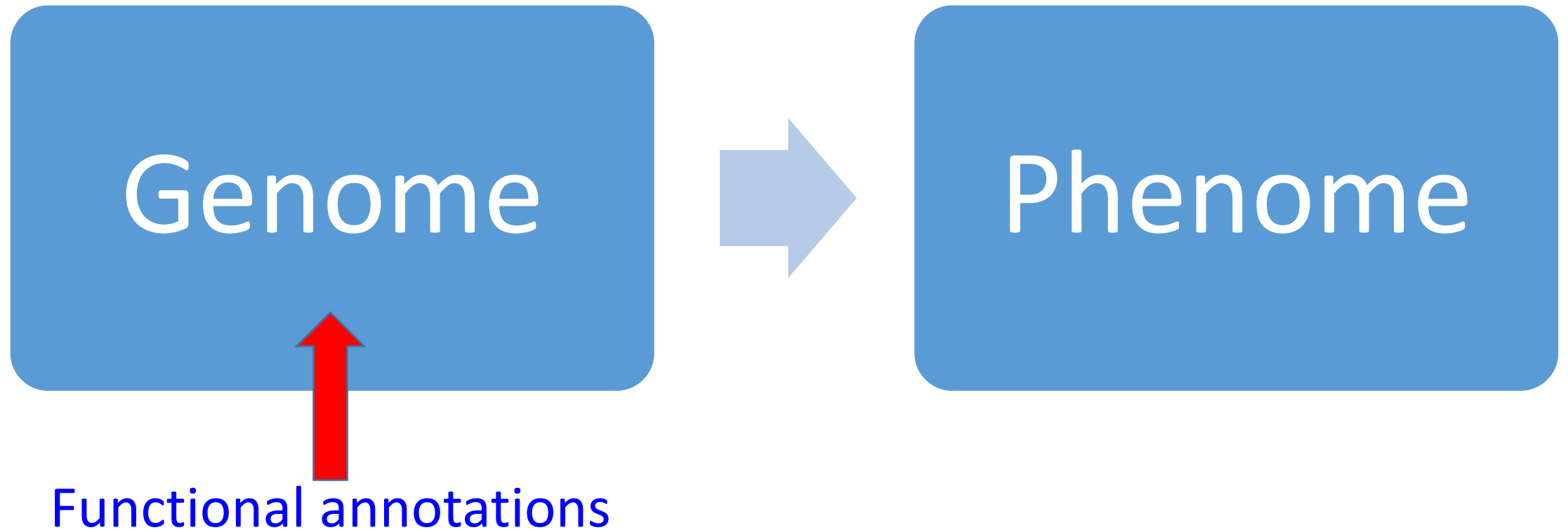


# Functional Annotation in Genomics

Christian Maltecca

# Motivations



# Introduction to Functional Annotation

- **What is Functional Annotation?**

- Functional annotation involves assigning biological information to genomic features, such as genes, SNPs, or other genomic regions.
- Helps in understanding the biological role of genetic variants and their impact on phenotypic traits.

- **Why is Functional Annotation Important?**

- Provides insights into the molecular mechanisms underlying complex traits and diseases.
- Facilitates the interpretation of results from genome-wide association studies (GWAS) and other genomic analyses.
- Aids in identifying potential therapeutic targets and biomarkers for disease.

# Functional annotations

- External information about a variant **in addition to** its **genotypes**
  - <https://useast.ensembl.org/info/genome/index.html>
- Location
  - Whether SNP is within a specific **genomic element**
  - There are many resources to define **genomic elements**.
- Quantification
  - Pathogenicity scores
  - [https://useast.ensembl.org/info/genome/variation/prediction/protein\\_function.html](https://useast.ensembl.org/info/genome/variation/prediction/protein_function.html)

# Types of Functional Annotations

- **Gene-Centric Annotations:**

- **Gene Ontology (GO):** Provides a controlled vocabulary to describe gene products in terms of their associated biological processes, cellular components, and molecular functions.
- **Pathway Annotations:** Describes genes in the context of biochemical pathways, such as KEGG or Reactome pathways.

- **Regulatory Annotations:**

- **Promoters and Enhancers:** Regions that regulate gene expression; identified through methods like chromatin immunoprecipitation (ChIP-seq).
- **Transcription Factor Binding Sites (TFBS):** Locations where transcription factors bind to regulate gene expression.

- **Epigenomic Annotations:**

- **Histone Modifications:** Modifications to histone proteins that affect chromatin structure and gene expression (e.g., H3K27ac marks active enhancers).
- **DNA Methylation:** Addition of methyl groups to DNA, affecting gene expression without changing the DNA sequence.

# Methods for Functional Annotation

- **Experimental Methods:**

- **ChIP-seq (Chromatin Immunoprecipitation Sequencing):** Used to identify DNA-binding sites of proteins, such as transcription factors or histones.
- **RNA-seq (RNA Sequencing):** Provides information on gene expression levels and splicing variants.

- **Computational Methods:**

- **Sequence Homology-Based Annotation:** Uses similarity to known sequences (e.g., BLAST) to infer function.
- **Machine Learning Approaches:** Utilizes algorithms to predict functional elements based on genomic features (e.g., DeepSEA).

- **Integrated Approaches:**

- **ENCODE Project:** Provides a comprehensive map of functional elements in the human genome, integrating multiple experimental data types.
- **Roadmap Epigenomics Project:** Focuses on characterizing the epigenomic landscape across different cell types and tissues.

# Tools and Databases for Functional Annotation

- **UCSC Genome Browser:**

- Offers a wide range of annotations, including genes, regulatory elements, and epigenomic data.

- **Ensembl:**

- Provides comprehensive genome annotation, including genes, variants, regulatory regions, and comparative genomics data.

- **GREAT (Genomic Regions Enrichment of Annotations Tool):**

- Associates genomic regions with biological functions by leveraging annotations from GO, pathways, and other sources.

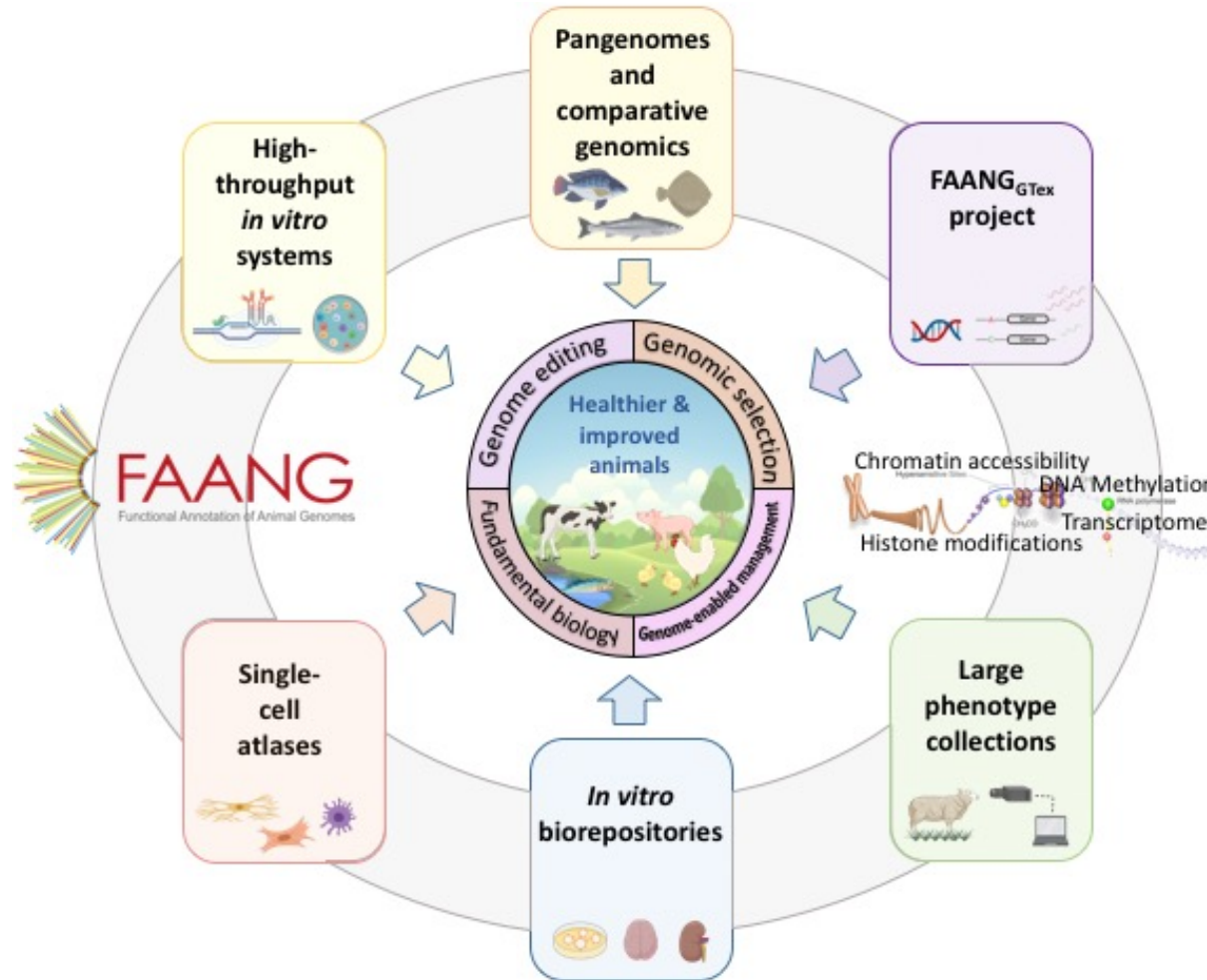
- **RegulomeDB:**

- Integrates various types of functional annotation data to score regulatory potential of non-coding variants.

- **dbSNP and ClinVar:**

- Provide annotations on known genetic variants, including their functional consequences and clinical significance.

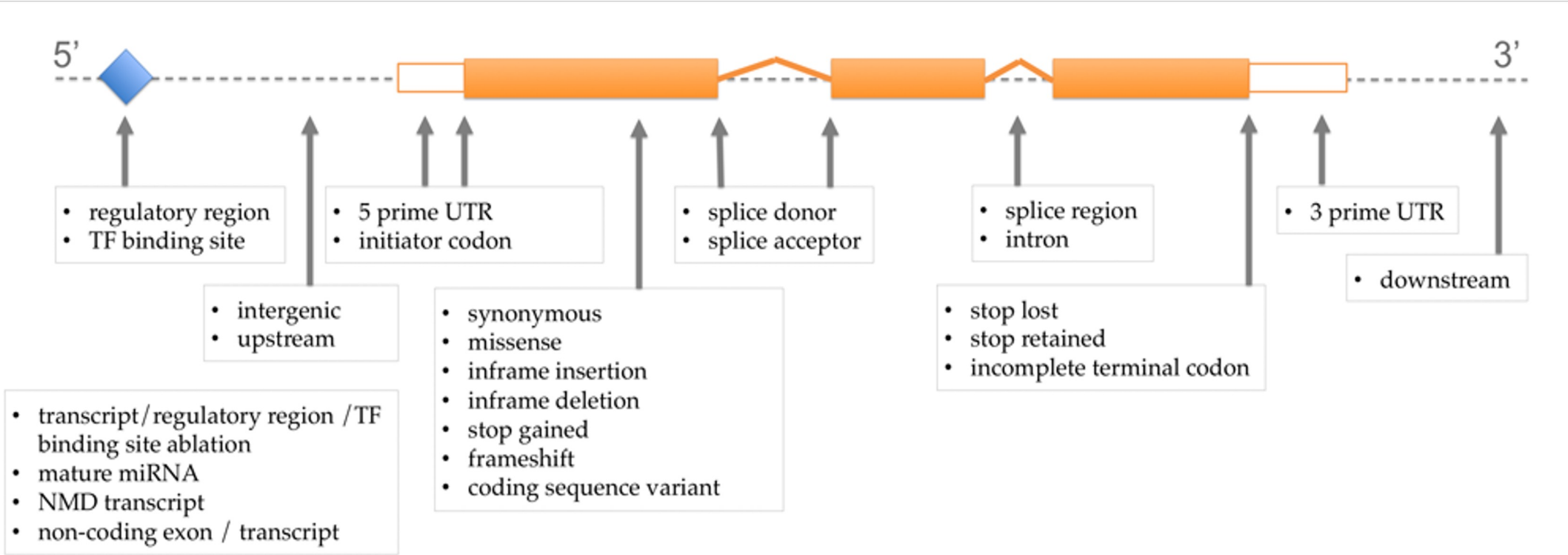
# Functional Annotation of Animal Genomes



<https://www.animalgenome.org/community/FAANG/>



# Location-based



# Quantification-based

Filter SIFT: All PolyPhen: All Consequences: missense variant Source: dbSNP CADD: All Filter Other Columns

Show/hide columns

Variant ID	Chr: bp	Alleles	Evidence	AA	AA coord	SIFT	CADD	IC50	REVEL	Primate	Mutation Assessor
<a href="#">rs1437042753</a>	3:25609282	G/A	gnomAD	R/W	1327	0					0.812
<a href="#">rs1246519283</a>	3:25624786	G/A	gnomAD	R/W	743	0					0.968
<a href="#">rs1173166728</a>	3:25627237	G/A	gnomAD	R/C	651	0	0.988	34	0.52	0.415	0.936
<a href="#">rs1357930313</a>	3:25632778	T/C		D/G	343	0	0.949	34	0.789	0.39	0.837
<a href="#">rs776876327</a>	3:25634002	G/A/T	Ex AC gnomAD	R/C	284	0.01	0.875	34	0.434	0.259	0.673
<a href="#">rs376036396</a>	3:25615236	T/A	Ex AC gnomAD	E/V	1182	0.01	0.786	33	0.362	0.295	0.709
<a href="#">rs1349960311</a>	3:25615269	A/G		L/P	1171	0	0.998	33	0.836	0.801	0.967

CADD

0 - 100

0 1

include blank

Apply Cancel

# Categorical annotations

- [https://useast.ensembl.org/info/genome/variation/prediction/predicted\\_data.html](https://useast.ensembl.org/info/genome/variation/prediction/predicted_data.html)

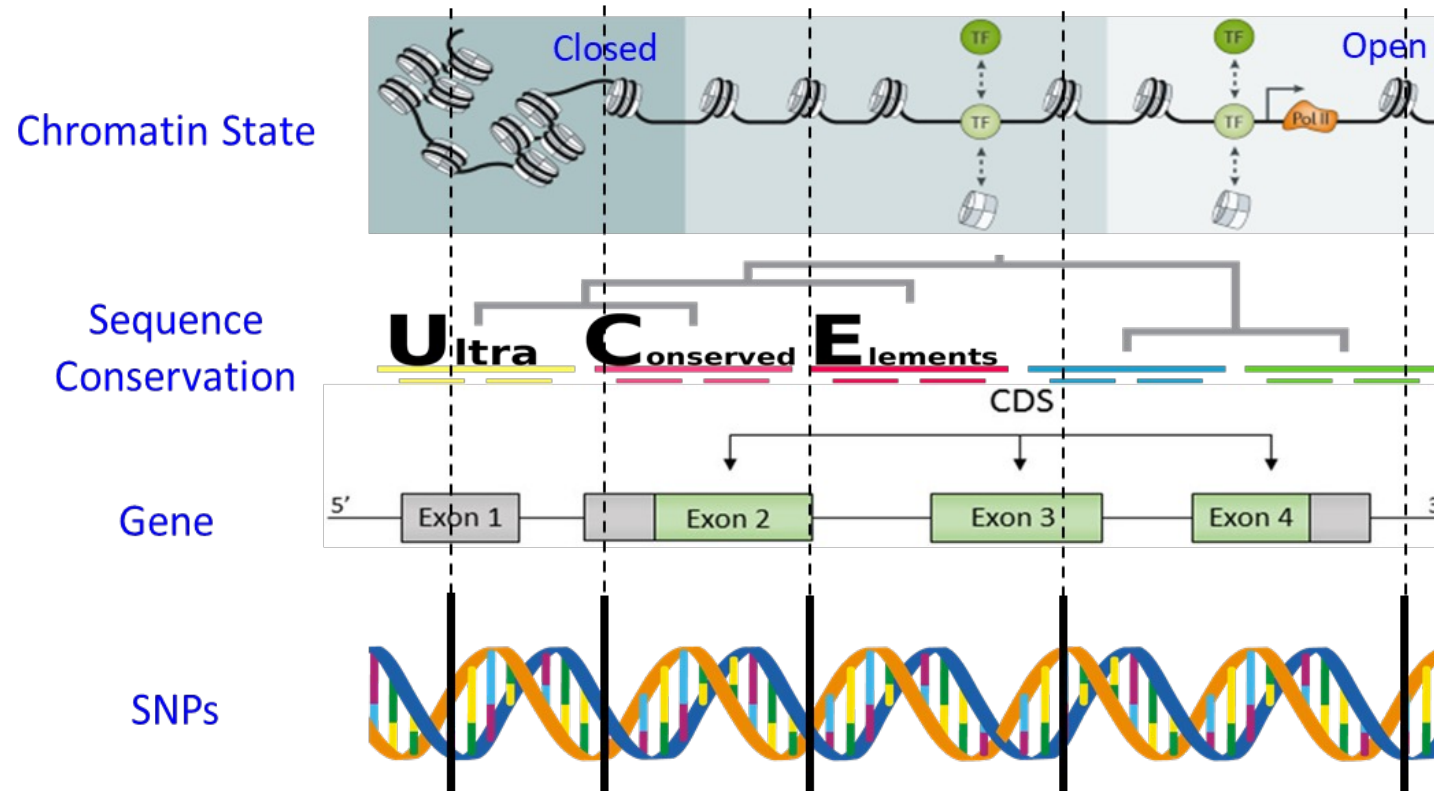
* SO term	SO description	SO accession	Display term	IMPACT
stop_gained	A sequence variant whereby at least one base of a codon is changed, resulting in a premature stop codon, leading to a shortened transcript	<a href="#">SO:0001587</a>	Stop gained	HIGH
frameshift_variant	A sequence variant which causes a disruption of the translational reading frame, because the number of nucleotides inserted or deleted is not a multiple of three	<a href="#">SO:0001589</a>	Frameshift variant	HIGH
stop_lost	A sequence variant where at least one base of the terminator codon (stop) is changed, resulting in an elongated transcript	<a href="#">SO:0001578</a>	Stop lost	HIGH
start_lost	A codon variant that changes at least one base of the canonical start codon	<a href="#">SO:0002012</a>	Start lost	HIGH
inframe_insertion	An inframe non synonymous variant that inserts bases into in the coding sequence	<a href="#">SO:0001821</a>	Inframe insertion	MODERATE
inframe_deletion	An inframe non synonymous variant that deletes bases from the coding sequence	<a href="#">SO:0001822</a>	Inframe deletion	MODERATE
missense_variant	A sequence variant, that changes one or more bases, resulting in a different amino acid sequence but where the length is preserved	<a href="#">SO:0001583</a>	Missense variant	MODERATE
protein_altering_variant	A sequence_variant which is predicted to change the protein encoded in the coding sequence	<a href="#">SO:0001818</a>	Protein altering variant	MODERATE
splice_region_variant	A sequence variant in which a change has occurred within the region of the splice site, either within 1-3 bases of the exon or 3-8 bases of the intron	<a href="#">SO:0001630</a>	Splice region variant	LOW
splice_donor_region_variant	A sequence variant that falls in the region between the 3rd and 6th base after splice junction (5' end of intron)	<a href="#">SO:0002170</a>	Splice donor region variant	LOW

- Often binary
  - eQTL: whether a SNP is an identified eQTL. **Yes or no.**
  - CDS: whether a SNP is within CDS. **Yes or no.**

# Continuous annotations

- Less commonly used than categorical annotations.
- Can be categorized.
- Minor allele frequency
  - <0.01, 0.01-0.05, 0.05-0.10, etc
- Conservation score
  - Constrained element
  - [https://useast.ensembl.org/info/genome/compara/conservation\\_and\\_constrained.html](https://useast.ensembl.org/info/genome/compara/conservation_and_constrained.html)

# Link functional annotations to SNPs



Yes or no: whether a SNP is within a functional annotation category

# Theory

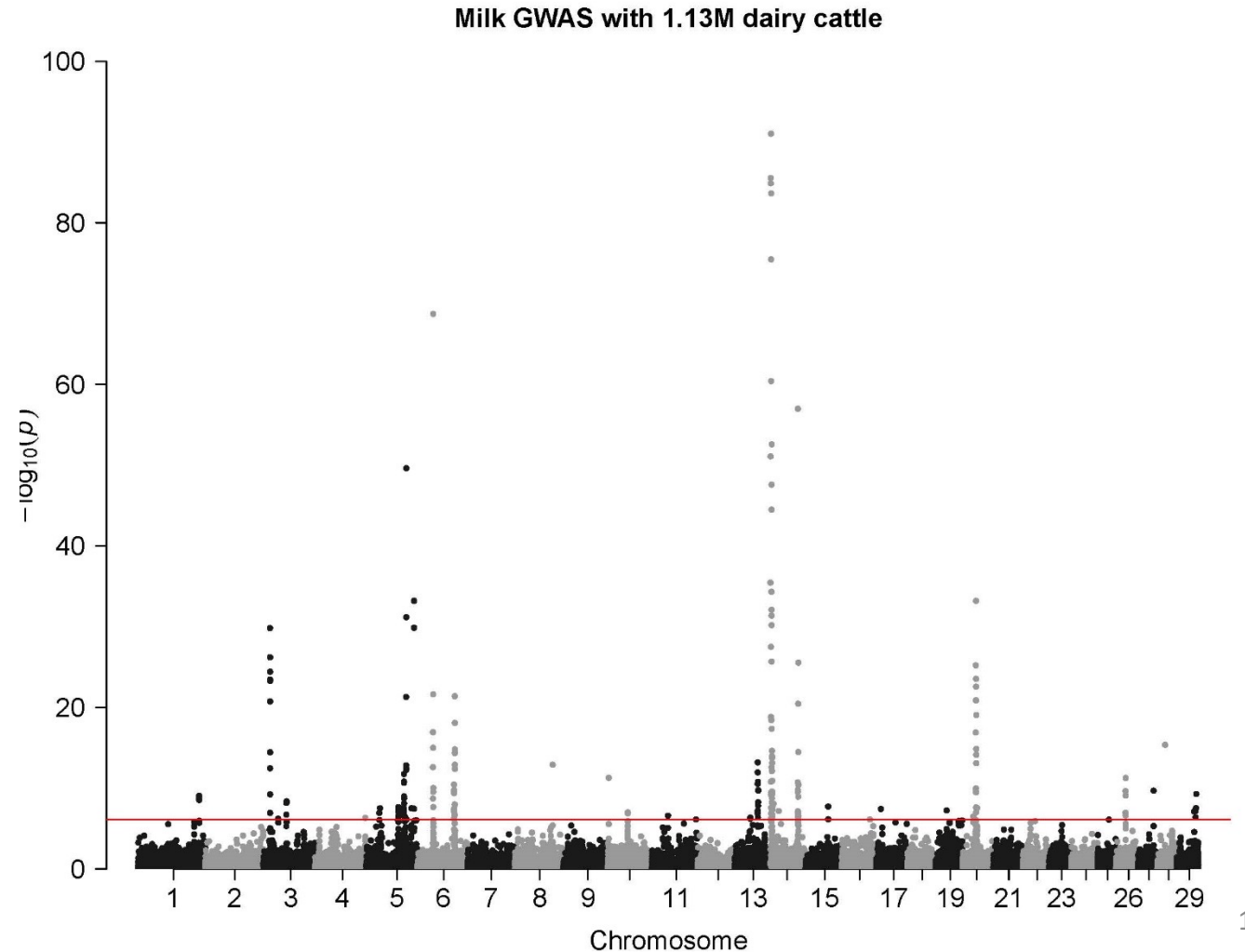
# Genetic Architecture of Complex Traits

- Limit of GWAS

$H_0$ : SNP effect size drawn from a normal distribution  $N(0, \sigma_g^2/M)$ .

$M$ : effective number of independent markers

- Proportion of genetic variance explained by significant SNPs  $\approx$  **18.3%**



# Genetic Architecture of Complex Traits (*cont.*)

- Omnigenic model
  - A small number of “core” genes
    - SNP effect size beyond polygenic effect  $N(0, \sigma_g^2/M)$
  - A lot more “peripheral” genes
    - Contributing much more of heritability
    - SNP effect size **NOT** beyond polygenic effect  $N(0, \sigma_g^2/M)$
- Core genes are easy to deal with in genomic predictions.
  - Bayesian mixture models
    - BayesB, BayesR, BSLMM
  - Fixing big QTL effects
    - LDAK
- Difficult to model peripheral genes **better than GBLUP**



# Modeling Functional Annotations

	Annot1	Annot2	Annot3	Annot4
SNP1	1	1	0	0
SNP2	0	1	1	0
SNP3	1	1	1	1

} **W**

$$\sigma_{\text{SNP1}}^2 = \tau_1 + \tau_2$$

$$\sigma_{\text{SNP2}}^2 = \tau_2 + \tau_3$$

$$\sigma_{\text{SNP3}}^2 = \tau_1 + \tau_2 + \tau_3 + \tau_4$$

$$= \mathbf{W}\boldsymbol{\tau}$$

$\text{Var}(\alpha_j) = \sum_{k:j \in S_k} \tau_k$ , where  $\tau_k$  denotes the per-variant contribution of category  $k$  to  $\text{Var}(\alpha_j)$ .

# Linear Mixed Model with Functional Annotations

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha} + \mathbf{e}$$

$$\boldsymbol{\alpha} \sim N(\mathbf{0}, \mathbf{D})$$

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

$$\mathbf{D} = \begin{bmatrix} \sigma_{\alpha_1}^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_{\alpha_M}^2 \end{bmatrix}$$

$$\begin{bmatrix} \sigma_{\alpha_1}^2 & \sigma_{\alpha_2}^2 & \cdots & \sigma_{\alpha_M}^2 \end{bmatrix}^T = \mathbf{W}\boldsymbol{\tau}$$



$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{g} + \mathbf{e}$$

$$\mathbf{g} \sim N(\mathbf{0}, \sum_{k=1}^K \mathbf{G}_k M_k \tau_k)$$

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

where  $\mathbf{g} = \mathbf{Z}\boldsymbol{\alpha}$ ,  $\mathbf{G}_k = \mathbf{Z}\mathbf{W}^{(k)}\mathbf{Z}^T / M_k$

$\mathbf{W}^{(k)}$  is a **diagonal matrix** whose diagonal elements are  $\mathbf{W}$ 's  $k^{th}$  **column**.

Integrative analysis becomes VC estimation.

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{g} + \mathbf{e}$$

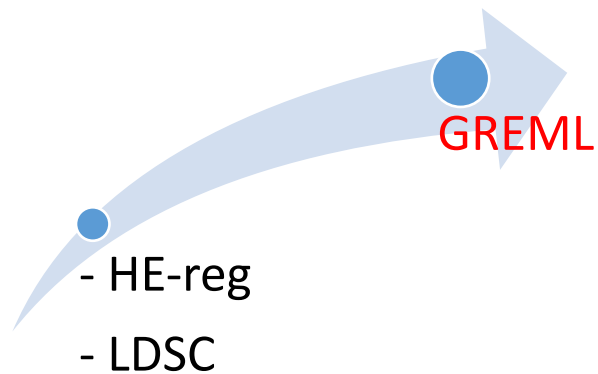
$$\mathbf{g} \sim N(\mathbf{0}, \sum_{k=1}^K \mathbf{G}_k M_k \tau_k)$$

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

- We may need to model many functional annotations, so  $K$  is large.
- VC estimates ( $\tau_k$ ) may be negative.

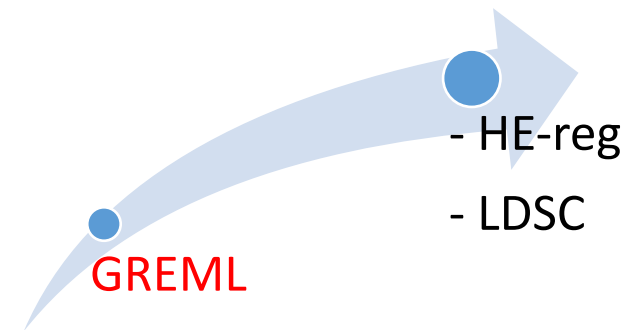
# Existing VC estimation methods

## Statistical efficiency



- HE or LDSC is not statistically efficient.

## Stability and scalability



- GREML may **fail to converge** for many VCs.
- GREML may be **slow** for large samples.

# Linkage Disequilibrium Score Regression

# Introduction to Linkage Disequilibrium (LD) Score Regression

- **What is Linkage Disequilibrium (LD)?**
  - Linkage Disequilibrium (LD) refers to the non-random association of alleles at different loci. It indicates that certain combinations of alleles occur together more frequently than would be expected by chance.
  - Influenced by factors such as genetic drift, selection, mutation, recombination, and population structure.
- **What is LD Score Regression?**
  - A statistical method used to partition heritability and identify confounding biases in genome-wide association studies (GWAS).
  - Utilizes LD patterns to differentiate true polygenic signals from confounding biases in GWAS data.

# Understanding LD Score

- **LD Score Definition:**

- The LD score of a single nucleotide polymorphism (SNP) is the sum of squared correlations ( $r^2$ ) between the SNP and all other SNPs within a certain genomic window.
- Quantifies the amount of genetic variation captured by a SNP.

$$\text{LD Score } (l_j) = \sum_{i=1}^M r_{ij}^2$$

where  $(r_{ij}^2)$  is the squared correlation between SNP ( $j$ ) and SNP ( $i$ ), and  $(M)$  is the number of SNPs within the specified window.

- **Why is LD Score Important?**

- SNPs with high LD scores tag more genetic variation, aiding in differentiating between true signal and noise in association studies.

# How LD Score Regression Works

- **Core Idea:**

- In a polygenic trait, SNPs with higher LD scores should, on average, have higher GWAS test statistics (chi-squared statistics) because they tag more causal variants.

- **Key Components:**

- **Dependent Variable:** GWAS test statistics (chi-squared values) for each SNP.
- **Independent Variable:** LD scores for each SNP.

- **Regression Model:**

$$E[\chi_j^2] = 1 + \frac{N \cdot h^2 \cdot l_j}{M} + \alpha$$

where:

- $(E[\chi_j^2])$  is the expected chi-squared statistic for SNP ( j ),
- $(N)$  is the sample size,
- $(h^2)$  is the heritability explained by the SNPs,
- $(l_j)$  is the LD score of SNP ( j ),
- $(M)$  is the number of SNPs, and
- $(\alpha)$  is an intercept term representing confounding bias.



$$\mathbf{y} = \mathbf{Z}\boldsymbol{\alpha} + \mathbf{e}, \quad \boldsymbol{\alpha} \sim N(0, \mathbf{I}\sigma_{\alpha}^2) \quad \mathbf{e} \sim N(0, \mathbf{I}\sigma_e^2)$$

$$\frac{\mathbf{Z}^T \mathbf{y}}{\sqrt{N\sigma_p^2}} \sim N\left(0, \frac{\mathbf{Z}^T \mathbf{Z} \mathbf{Z}^T \mathbf{Z} \sigma_{\alpha}^2}{N\sigma_p^2} + \frac{\mathbf{Z}^T \mathbf{Z} \sigma_e^2}{N\sigma_p^2}\right),$$

$$\mathbf{s} \sim N\left(0, N\mathbf{C}^2 b^2 / M + \mathbf{C}(1 - b^2)\right),$$

$$\mathbf{s} \sim N\left(0, \text{diag}\left(N\mathbf{C}^2 b^2 / M + \mathbf{C}(1 - b^2)\right)\right);$$

$$s_j \sim N\left(0, N\ell_j b^2 / M + (1 - b^2)\right),$$

# Applications of LD Score Regression

- **Estimating Heritability:**
  - Used to estimate the heritability of complex traits by partitioning genetic variance from GWAS summary statistics.
- **Detecting Confounding Bias:**
  - Intercept ( $\alpha$ ) in the regression can identify confounding biases such as population stratification or cryptic relatedness.
- **Improving GWAS Results:**
  - Differentiates true polygenic signals from spurious associations caused by confounding factors.
- **Understanding Genetic Architecture:**
  - Provides insights into the genetic architecture of complex traits by assessing the genetic variance explained by SNPs across the genome.

# Pros and Cons of LD Score Regression

- **Advantages:**

- **Robust to Confounding:** Effectively identifies and corrects for population stratification and other confounders in GWAS data.
- **Simple and Scalable:** Uses summary statistics and is computationally efficient for large-scale GWAS.
- **Broad Applicability:** Applicable to both quantitative traits and case-control studies.

- **Limitations:**

- **Requires Large Sample Sizes:** Accurate estimates of heritability and bias detection need large GWAS sample sizes.
- **Assumptions:** Assumes GWAS results are polygenic and LD patterns are well-captured, which may not always hold.
- **Sensitive to LD Score Calculation:** Accuracy depends on the quality of LD score calculations, which can vary across populations and reference panels.

# References

1. Bulik-Sullivan, B., Loh, P.-R., Finucane, H. K., et al. (2015). *LD Score Regression Distinguishes Confounding from Polygenicity in Genome-Wide Association Studies*. *Nature Genetics*, 47, 291–295.
2. Finucane, H. K., Bulik-Sullivan, B., Gusev, A., et al. (2015). *Partitioning Heritability by Functional Annotation Using Genome-Wide Association Summary Statistics*. *Nature Genetics*, 47, 1228–1235.
3. Lee, J. J., Wedow, R., Okbay, A., et al. (2018). *Gene Discovery and Polygenic Prediction from a Genome-wide Association Study of Educational Attainment in 1.1 Million Individuals*. *Nature Genetics*, 50, 1112–1121.

# Partitioning Heritability using GREML

# Introduction to GREML

- **What is GREML?**

- Genomic-Relatedness-Based Restricted Maximum Likelihood (GREML) is a statistical method used to estimate the proportion of phenotypic variance that is attributable to genetic variance (heritability).
- GREML utilizes information from genomic data to partition heritability into components explained by different sets of genetic variants, often grouped by their minor allele frequency (MAF), functional annotation, or chromosomal location.

- **Applications of GREML:**

- Understanding the genetic architecture of complex traits.
- Estimating heritability attributable to specific genetic components.
- Identifying the contribution of different genomic regions or functional annotations to the overall genetic variance of a trait.

# The Basics of Heritability and GREML

- **Heritability  $h^2$ :**

- The proportion of the total phenotypic variance of a trait that can be attributed to genetic variance.
- Broad-sense heritability ( $H^2$ ) includes all genetic variance components, while narrow-sense heritability ( $h^2$ ) focuses only on additive genetic variance.

- **GREML Framework:**

- Uses a linear mixed model to partition phenotypic variance into genetic and environmental components.
- The model estimates the genetic variance explained by SNPs captured in a genomic relationship matrix (GRM).

$$y = X\beta + g + e$$

where:

- ( $y$ ) is the vector of phenotypic values,
- ( $X$ ) is the matrix of fixed effects,
- ( $\beta$ ) is the vector of fixed effect coefficients,
- ( $g$ ) is the vector of genetic effects (random effects),
- ( $e$ ) is the vector of residual (environmental) effects.

# Constructing the Genomic Relationship Matrix (GRM)

- **Genomic Relationship Matrix (GRM):**

- Represents the genetic similarity between all pairs of individuals based on their SNP genotypes.
- The GRM is constructed using genotype data to estimate the proportion of the genome shared identical-by-state (IBS) between individuals.

$$\text{GRM}_{ij} = \frac{1}{M} \sum_{k=1}^M \frac{(x_{ik} - 2p_k)(x_{jk} - 2p_k)}{2p_k(1 - p_k)}$$

where:

- $(M)$  is the number of SNPs,
- $(x_{ik})$  and  $(x_{jk})$  are the genotypes for individuals  $i$  and  $j$  at SNP  $k$ ,
- $(p_k)$  is the allele frequency of SNP  $k$ .



# Partitioning Heritability with GREML

- **Objective:**

- Partition the heritability of a complex trait into components attributable to different sets of SNPs (e.g., by MAF, functional category, or genomic region).

- **GREML Model for Partitioning:**

$$V(y) = \sigma_g^2 K + \sigma_e^2 I$$

where:

- $V(y)$  is the phenotypic variance-covariance matrix,
- $\sigma_g^2$  is the additive genetic variance explained by SNPs,
- $K$  is the GRM, capturing genetic relationships,
- $\sigma_e^2$  is the residual variance,
- $I$  is the identity matrix.

# Steps in GREML Analysis

## 1. Prepare Data:

- Obtain SNP genotype data and phenotype data.
- Quality control (QC) steps to filter SNPs and individuals (e.g., removing SNPs with low MAF, high missingness, or poor quality).

## 2. Construct the GRM:

- Calculate the GRM using quality-controlled SNPs.

## 3. Fit the GREML Model:

- Use software like GCTA, BOLT-LMM, or LDAK to fit the linear mixed model and estimate variance components.

## 4. Partition Heritability:

- Decompose total genetic variance into components explained by different sets of SNPs (e.g., different MAF bins or functional categories).

# Interpretation of GREML Results

- **Total Heritability Estimate:**
  - The overall proportion of phenotypic variance explained by all SNPs in the model.
- **Partitioned Heritability Estimates:**
  - Heritability explained by specific subsets of SNPs, allowing for interpretation of the contribution of various genetic factors.
- **Implications for Genetic Architecture:**
  - Helps in understanding the distribution of genetic effects across the genome.
  - Provides insights into whether certain genomic regions or functional annotations contribute disproportionately to the trait's heritability.

# Advantages and Limitations of GREML

- **Advantages:**

- Allows estimation of SNP-based heritability using only genotype data.
- Can partition heritability based on functional annotations or other criteria.
- Robust to population stratification if properly accounted for in the GRM.

- **Limitations:**

- Requires large sample sizes for accurate heritability estimates.
- Sensitive to the quality and representativeness of the reference population used to calculate the GRM.
- Assumes a polygenic model, where many variants of small effect contribute to the trait.

# References

1. Yang, J., Lee, S. H., Goddard, M. E., & Visscher, P. M. (2011). GCTA: A Tool for Genome-wide Complex Trait Analysis. *American Journal of Human Genetics*, 88(1), 76-82.
2. Lee, S. H., Wray, N. R., Goddard, M. E., & Visscher, P. M. (2011). Estimating Missing Heritability for Disease from Genome-wide Association Studies. *American Journal of Human Genetics*, 88(3), 294-305.
3. Yang, J., Zaitlen, N. A., Goddard, M. E., Visscher, P. M., & Price, A. L. (2014). Advantages and Pitfalls in the Application of Mixed-Model Association Methods. *Nature Genetics*, 46(2), 100-106.

# Jicai Jiang NC State University

<https://jiang18.github.io/mph/>

MPH Home Install Documentation Examples Performance Scripts Search Previous Next

MPH

Citation

License

Author and contact

## MPH

**MPH** (MINQUE for Partitioning Heritability) is a C++ program for fast REML estimation of genetic (co)variance components.

The method features:

- Effective implementation of Fisher-scoring REML (equivalent to iterative MINQUE)
- Fast and memory-efficient performance
- Robust convergence
- Support for analyses of dominance, epistasis, and genetic correlation
- Capability to perform complex genome-partitioning of quantitative genetic variation and covariation

## Citation

Jiang J. MPH: fast REML for large-scale genome partitioning of quantitative genetic variation. *Bioinformatics*. 2024 May 2;40(5):btæ298. doi: [10.1093/bioinformatics/btæ298](https://doi.org/10.1093/bioinformatics/btæ298). PMID: [38688661](https://pubmed.ncbi.nlm.nih.gov/38688661/); PMCID: [PMC11093526](https://pubmed.ncbi.nlm.nih.gov/PMC11093526/).

## License

MPH is distributed under the [GPL-3.0 license](https://www.gnu.org/licenses/gpl-3.0.html).

## Author and contact

[Jicai Jiang](#)

If you want to submit an issue concerning the software, please do so using the [MPH GitHub repository](#).

<https://github.com/jiang18/slemm>

jiang18 Add expected		0e3a4d3 · last week	🕒 214 Commits
bin	fixed shebang		last year
docs	Update link to filling_bed.R		last week
examples	Add expected		last week
scripts	Ease examples		last week
src	Update pred.cpp		2 weeks ago
.gitattributes	Initial commit		2 years ago
CHANGELOG.md	Update CHANGELOG.md		2 weeks ago
LICENSE	Create LICENSE		2 years ago
Makefile	oneAPI		last year
README.md	Update README.md		3 months ago

README GPL-3.0 license

## SLEMM (Stochastic-Lanczos-Expedited Mixed Models)

SLEMM is a software tool for large-scale (up to millions) genomic predictions and genome-wide association studies.

## Citation

Cheng J, Maltecca C, VanRaden PM, O'Connell JR, Ma L, Jiang J. SLEMM: million-scale genomic predictions with window-based SNP weighting. *Bioinformatics*. 2023 Mar 1;39(3):btad127. doi: [10.1093/bioinformatics/btad127](https://doi.org/10.1093/bioinformatics/btad127). PMID: [36897019](https://pubmed.ncbi.nlm.nih.gov/36897019/); PMCID: [PMC10039786](https://pubmed.ncbi.nlm.nih.gov/PMC10039786/).

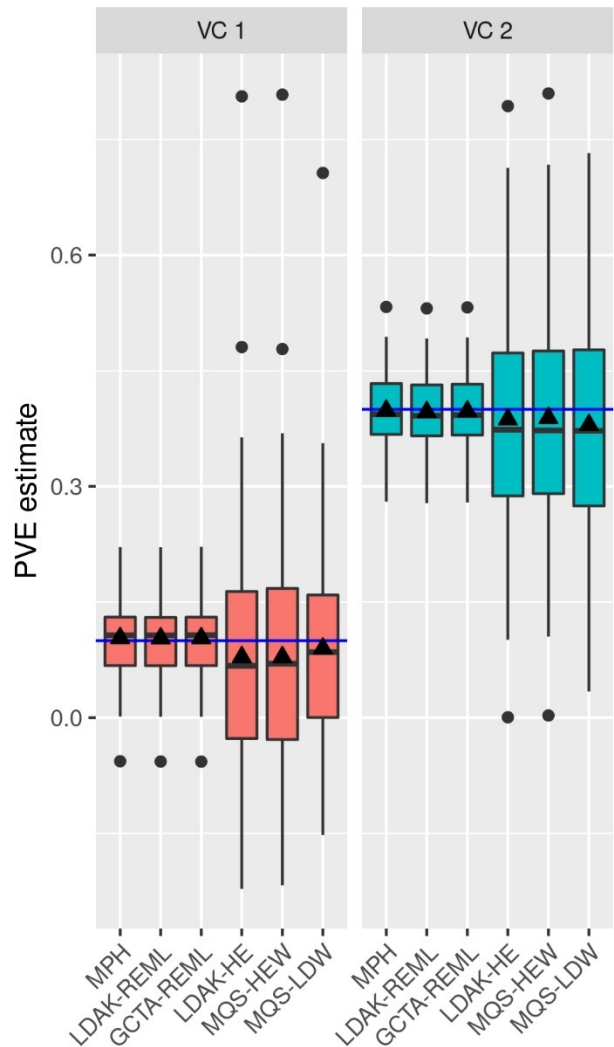
# LDSC and stratified LDSC

$$t_j \sim N\left(0, \frac{n}{m} \tilde{l}_j h^2 + c\right) \longrightarrow t_j \sim N\left(0, \frac{n}{m} \sum_{k=1}^K \tilde{l}_{jk} \tau_k + c\right)$$

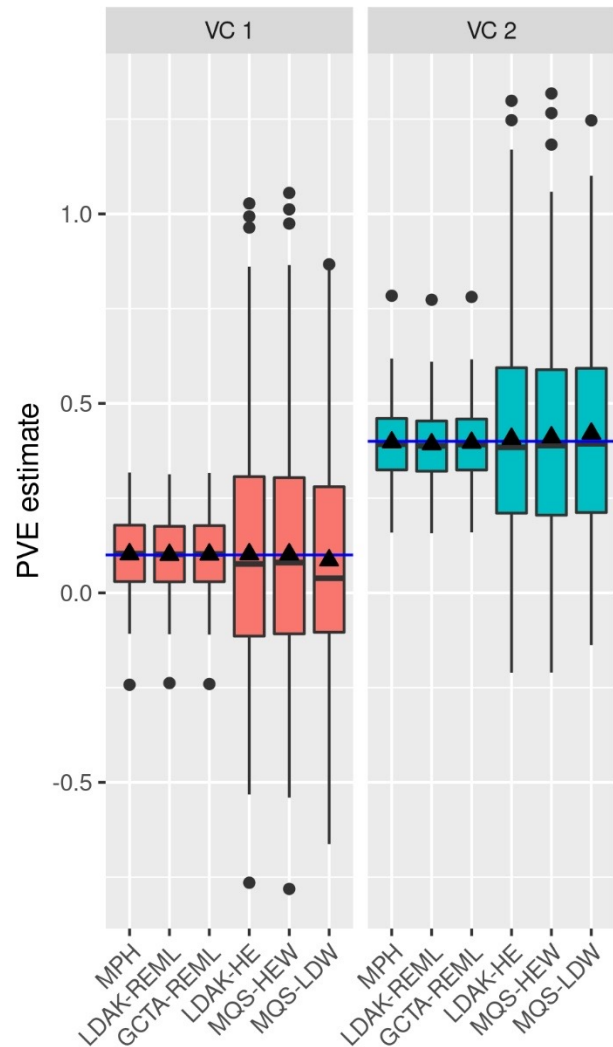
$c$  is estimated in LDSC, which can measure the genomic inflation of GWAS statistics.

Parameter estimation is based on **iteratively reweighted least squares**.

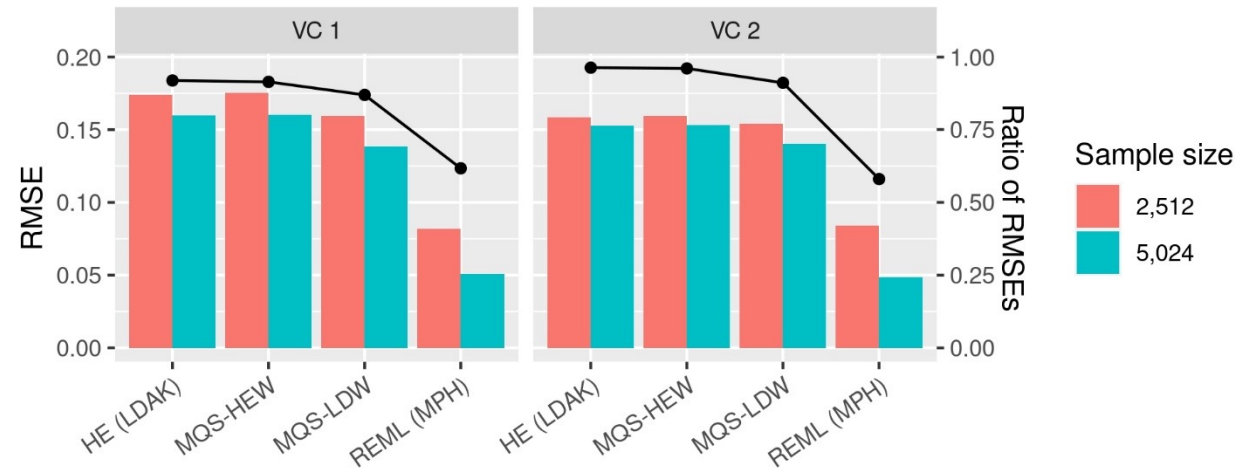
Dairy bull (5,024 individuals)



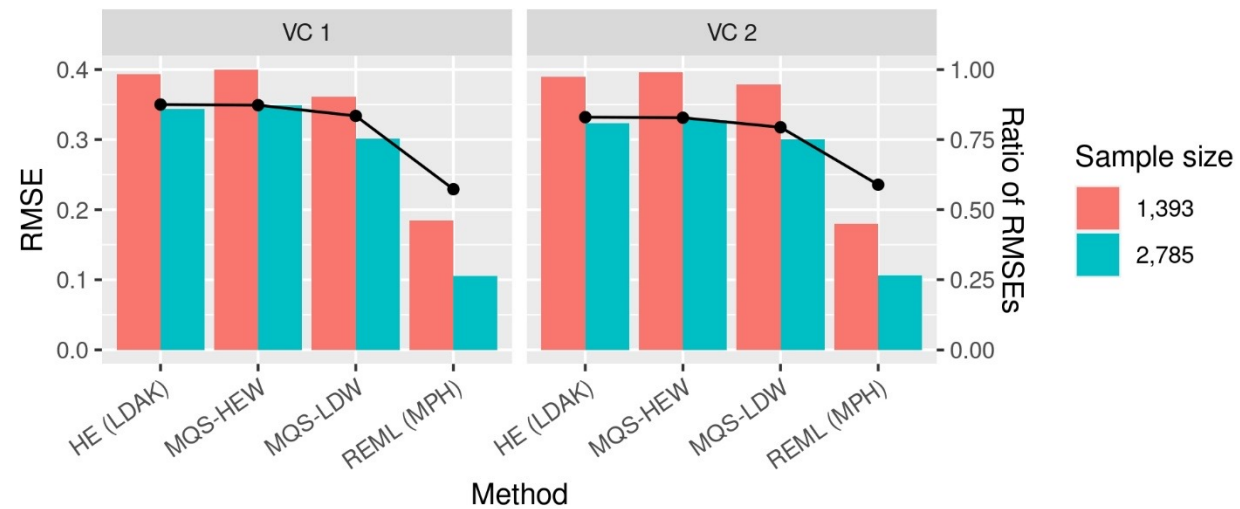
Duroc pig (2,785 individuals)



Dairy bull



Duroc pig



MQS-LDW (like LDSC) lacks statistical power.

MQS-LDW responds poorly to increasing sample size.



# Solution

- We need to use REML for the integration of functional annotations.
- Better REML implementation
  - Much more computationally efficient
  - More robust
- MPH (<https://jiang18.github.io/mph/>)

# Significance Tests and Predictions

- Does a function annotation contribute to heritability?
  - Wald tests on  $\tau_k$
  - LRT tests
- Functional-annotation-informed variance components
  - Genetic variance estimate for each marker
- Genomic predictions
  - BLUPs enhanced by the integration of functional annotations

# Functional enrichment

	Annot1	Annot2	Annot3	Annot4
SNP1	1	1	0	0
SNP2	0	1	1	0
SNP3	1	1	1	1

} **W**

$$\sigma_{\text{SNP1}}^2 = \tau_1 + \tau_2$$

$$\sigma_{\text{SNP2}}^2 = \tau_2 + \tau_3$$

$$\sigma_{\text{SNP3}}^2 = \tau_1 + \tau_2 + \tau_3 + \tau_4$$

$$= \mathbf{W}\boldsymbol{\tau}$$

$\text{Var}(\alpha_j) = \sum_{k:j \in S_k} \tau_k$ , where  $\tau_k$  denotes the per-variant contribution of category  $k$  to  $\text{Var}(\alpha_j)$ .

More generally,  $\text{Var}(\alpha_j) = \sum_k W_{jk} \tau_k$  for a binary ( $W_{jk}=0/1$ ) or continuous functional annotation.

# Functional enrichment (*cont.*)

$$\begin{bmatrix} \sum \text{var}(\alpha_1) \\ \sum \text{var}(\alpha_2) \\ \vdots \\ \sum \text{var}(\alpha_M) \end{bmatrix} = \mathbf{W}\hat{\boldsymbol{\tau}}$$

Genetic variance explained by SNPs in annotation category  $k$

$$v_k = \sum_{j:j \in S_k} \text{Var}(\alpha_j)$$

Per-SNP enrichment in annotation category  $k$

$$\rho_k = \left( \sum_{j:j \in S_k} \text{Var}(\alpha_j) / M_k \right) / \left( \sum_j \text{Var}(\alpha_j) / M \right)$$

# Functional enrichment (*cont.*)

Let  $\mathbf{Q}$  be a variant-to-annotation incidence matrix.

	Annot1	Annot2	Annot3	Annot4	Annot1,3
SNP1	1	1	0	0	1
SNP2	0	1	1	0	1
SNP3	1	1	1	1	1

$\uparrow$                      $\uparrow$                      $\uparrow$                      $\uparrow$                      $\uparrow$   
 $\mathbf{q}_1$                      $\mathbf{q}_2$                      $\mathbf{q}_3$                      $\mathbf{q}_4$                      $\mathbf{q}_5$

Genetic variance explained by SNPs in annotation category  $k$

$$v_k = \mathbf{q}_k^T \mathbf{W} \hat{\boldsymbol{\tau}}$$

Per-SNP enrichment in annotation category  $k$

$$\rho_k = \frac{\mathbf{q}_k^T \mathbf{W} \hat{\boldsymbol{\tau}} / M_k}{\mathbf{1}^T \mathbf{W} \hat{\boldsymbol{\tau}} / M}$$

# Functional enrichment (*cont.*)

Note that  $\hat{\boldsymbol{\tau}}$  and  $\text{var}(\hat{\boldsymbol{\tau}})$  are computed by REML.

Given  $\hat{\boldsymbol{\tau}}$  and  $\text{var}(\hat{\boldsymbol{\tau}})$ , standard errors for enrichment estimates can be readily computed.

Genetic variance explained by SNPs in annotation category  $k$

$$\text{var}(v_k) = \text{var}\left(\mathbf{q}_k^T \mathbf{W} \hat{\boldsymbol{\tau}}\right) = \mathbf{q}_k^T \mathbf{W} \text{var}(\hat{\boldsymbol{\tau}}) \mathbf{W}^T \mathbf{q}_k$$

Per-SNP enrichment in annotation category  $k$

$$\text{var}(\rho_k) = \text{var}\left(\frac{\mathbf{q}_k^T \mathbf{W} \hat{\boldsymbol{\tau}} / M_k}{\mathbf{1}^T \mathbf{W} \hat{\boldsymbol{\tau}} / M}\right)$$

Delta method

# Inverse-variance weighted average of estimates across seven type traits

