

Developing Computational Representations of Disease-Relevant Molecules:

Molecules:

3 Cases Studies for
AI in Biomedicine



Mark Gerstein, Yale

Slides freely downloadable from
Lectures.GersteinLab.org

& "tweetable" (via [@MarkGerstein](https://twitter.com/MarkGerstein)),
No Conflicts for this Talk. See last slide for more info.



Backgrounds

Learning meaningful representations
from large, complex biological data

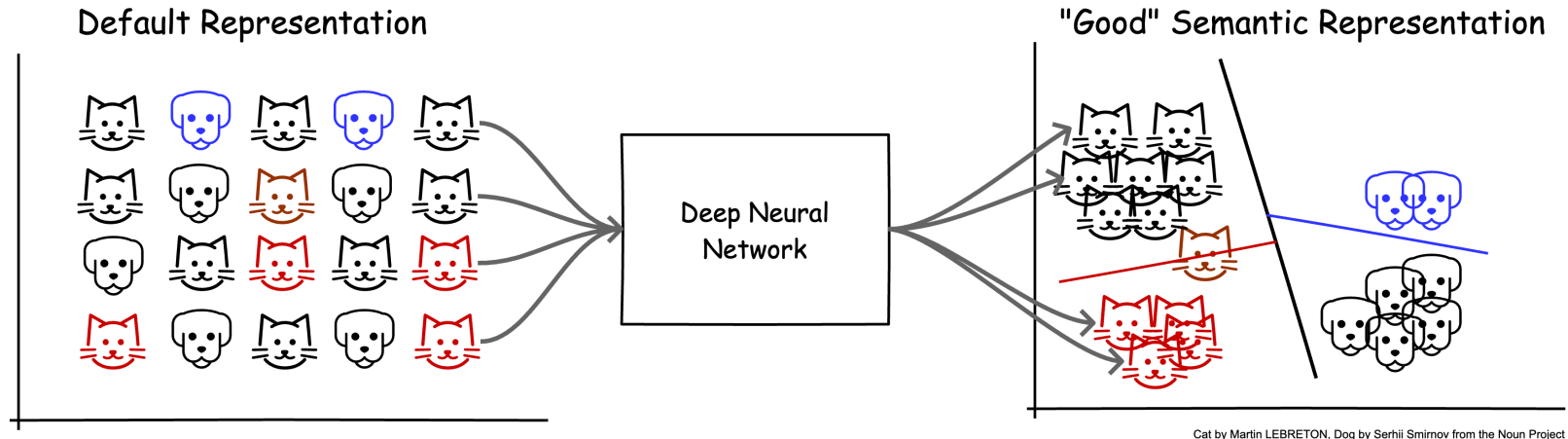


Representation learning

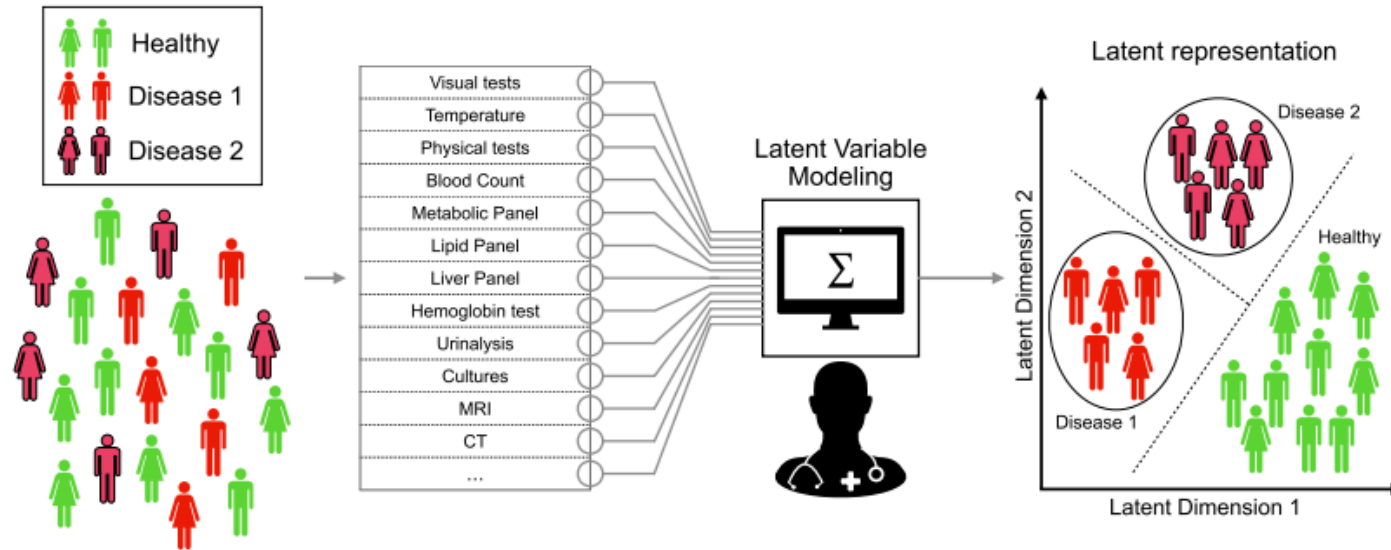
- *“An AI must fundamentally understand the world around us, and we argue that this can only be achieved if it can learn to identify and disentangle the **underlying explanatory factors** hidden in the observed milieu of low-level sensory data.”*

Bengio, Yoshua, Aaron Courville, and Pascal Vincent. "Representation learning: A review and new perspectives." *IEEE transactions on pattern analysis and machine intelligence* 35.8 (2013): 1798-1828.

Representation learning



Representation learning



Representation learning

- Particularly well suited to molecular biomedical data because of
 - its scale and high-dimensionality
 - its difficulty for easy interpretation
- When it comes to scientific data:
 - We only have partial knowledge about their internal structures
 - Thus, if the learned representations could re-discover some known patterns in the data, they could help us discover more potentially meaningful ones.

Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

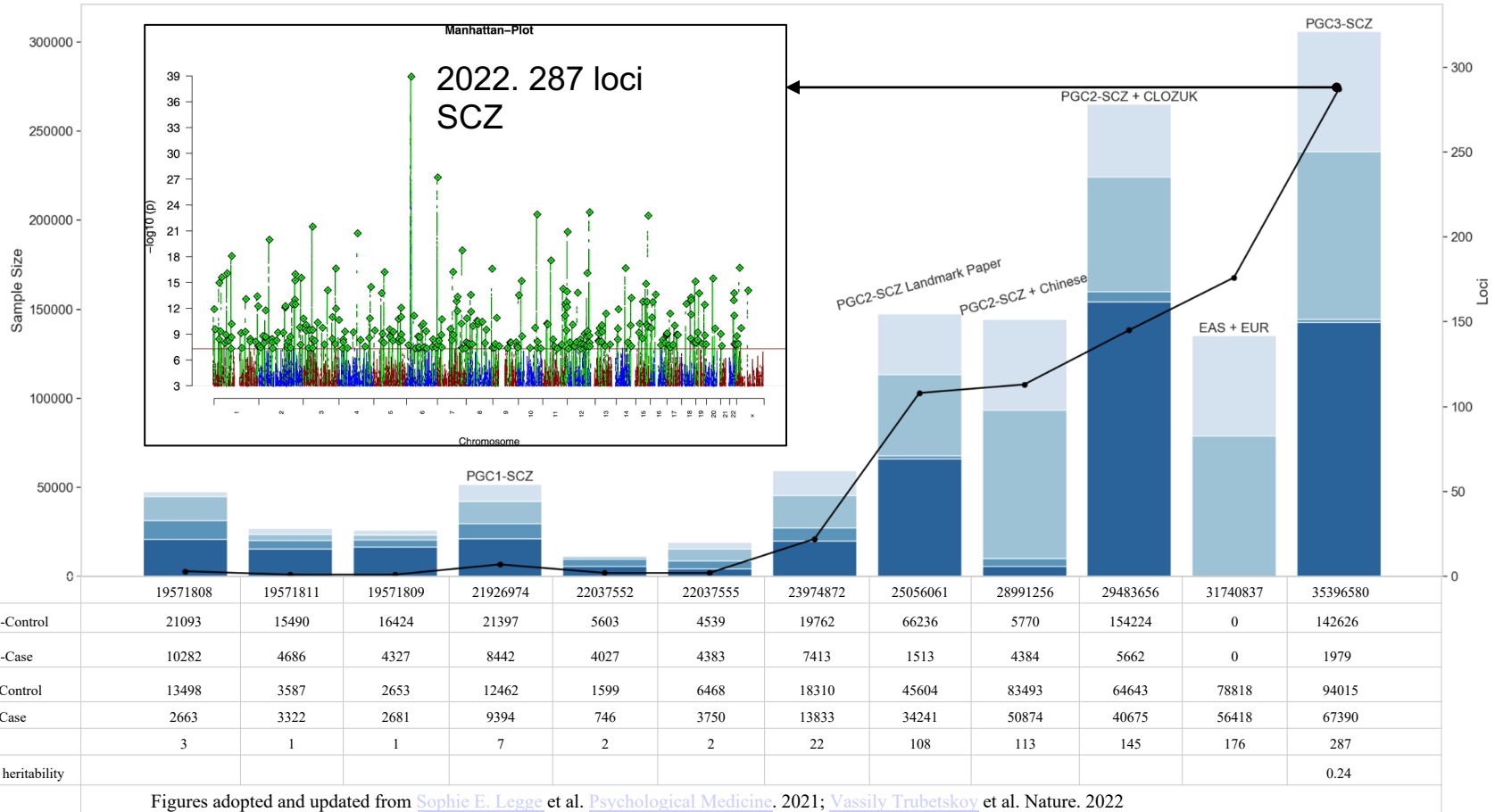
Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Many psychiatric conditions are highly heritable in comparison to other disorders, but their mechanisms are unknown

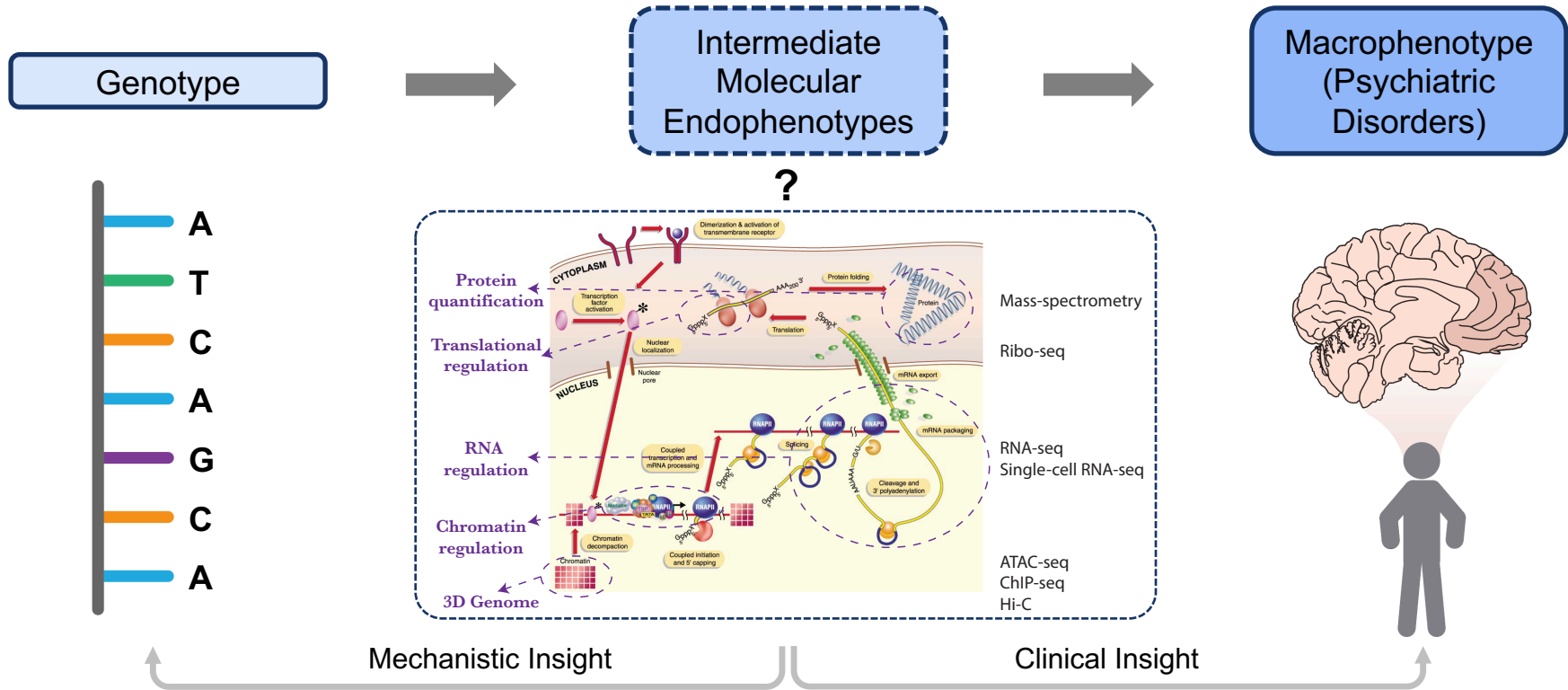
Disease	Heritability	SNP-based Heritability	PMID	Sample Size (1000s)	Molecular Mechanisms
Schizophrenia	79%	24%	35396580	320	(C4A)
Bipolar disorder	60-85%	16-19%	34002096	414	-
Alzheimer's disease	60-80%	3%	34493870	1,127	APOE, Tau
Hypertension	30%	18%	38689001	1,029	Renin–angiotensin–aldosterone
Heart disease	50-60%	6%	26343387	184	Atherosclerosis, VCAM-1
Stroke	32%	17%-21%	33773637	262	Reactive oxygen species (ROS), Ischemia
Type-2 diabetes	26%	20%	30054458	659	Insulin resistance
Breast Cancer	31%	22%	38741014	40	BRCA, PTEN

Great Progress in Finding Variants Related to Brain Diseases: The history of reported schizophrenia GWAS



Assessing gene regulation to understand psychiatric disorders

Addressing the fact that molecular mechanisms are not known for most psychiatric disorders

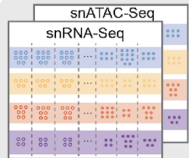
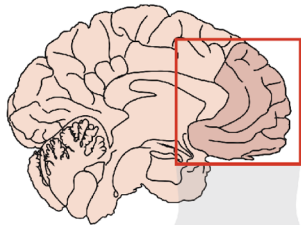


The PsychENCODE Consortium: Focusing on the PFC



PsychENCODE 200 researchers at 40 institutions

Main goal: *Understand the genetic, genomic and epigenomic etiologies of schizophrenia, bipolar disorder, autism spectrum disorder, and other neuropsychiatric disorders*

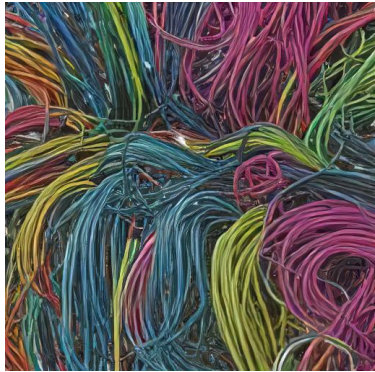


The **prefrontal cortex (PFC)** not only governs executive functions, but is also responsible for:

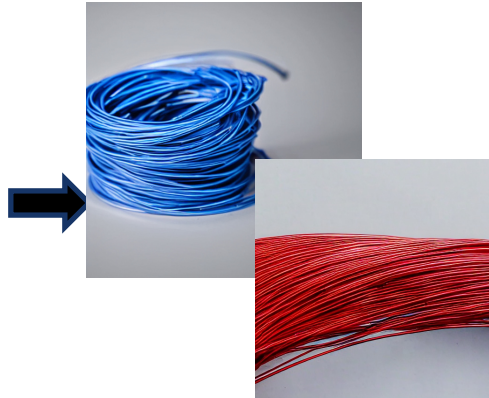
- behavioral regulation and mental health
- development and plasticity
- interplay with neurotransmitter systems

Advantages of single-cell resolution in the brain

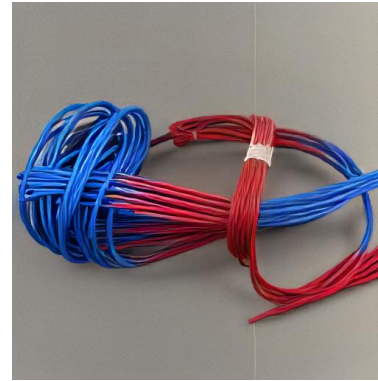
Bulk Datasets



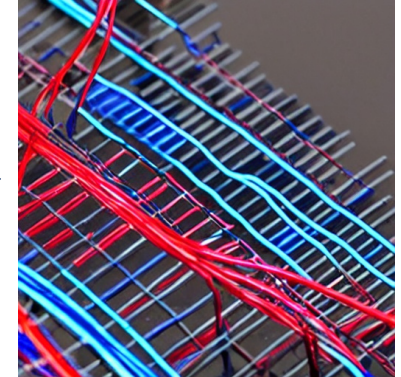
Single-cell Datasets



Spatial Analysis



Cell-to-cell Connectivity



Images generated using the DeepAI Image Generator tool

How PsychENCODE fits into the history of genome annotation



Worm Genome



modENCODE



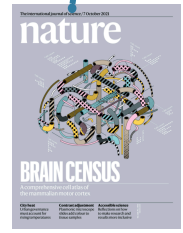
1000 Genomes Pilot



GTEx



1000 Genomes Production



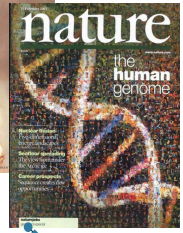
BICCN



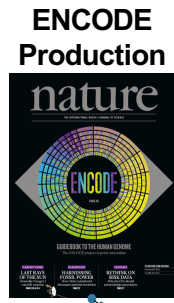
HuBMAP



The Human Genome Project



ENCODE Pilot



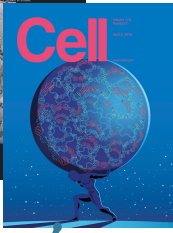
ENCODE Production



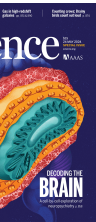
Comparative ENCODE



TCGA



PsychENCODE



2000

2005

2010

2015

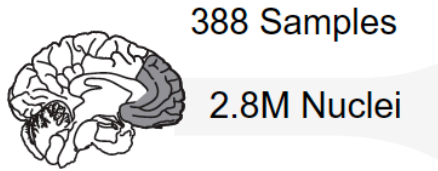
2020

2025

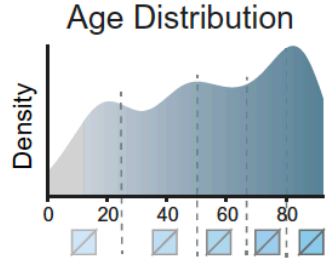
Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Integrating multi-omics data for 388 adult brains



Single-cell data for 388 individuals (snRNA-seq, genotypes, ~60 snATAC-seq) – **one of the largest single-cell collections in the human brain**

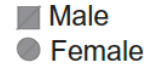


Integrated study derived from 12 cohorts (PEC, AMP-AD, & other studies) for population and cross-disorder comparisons

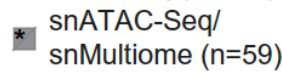
Ancestry



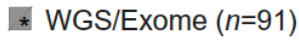
Biol. Sex



snRNA-Seq (n=388)



Genotypes (n=388)



MDD/ PTSD

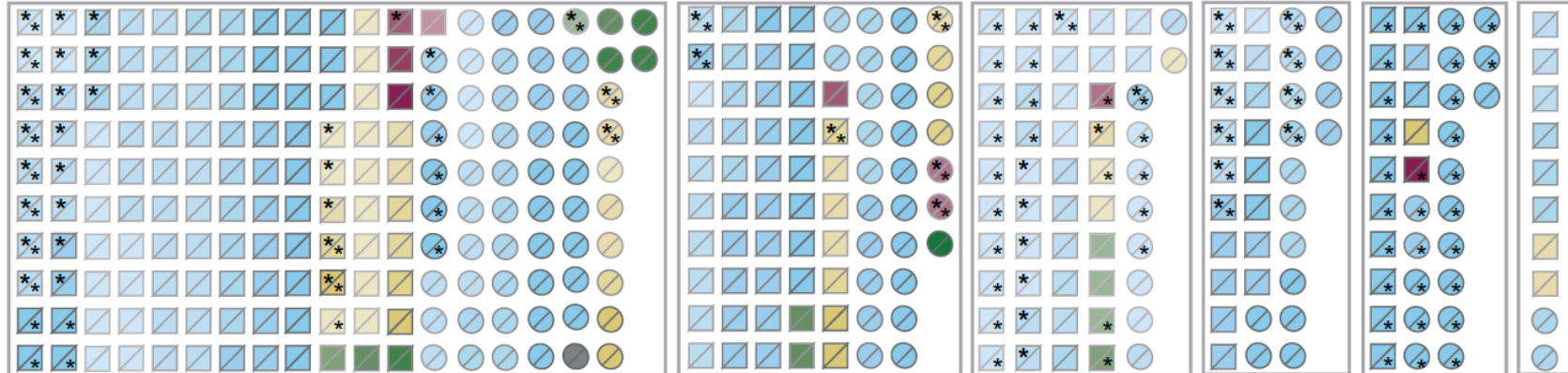
Control (182)

SCZ (77)

ASD (52)

BPD (34)

AD (33) (10)

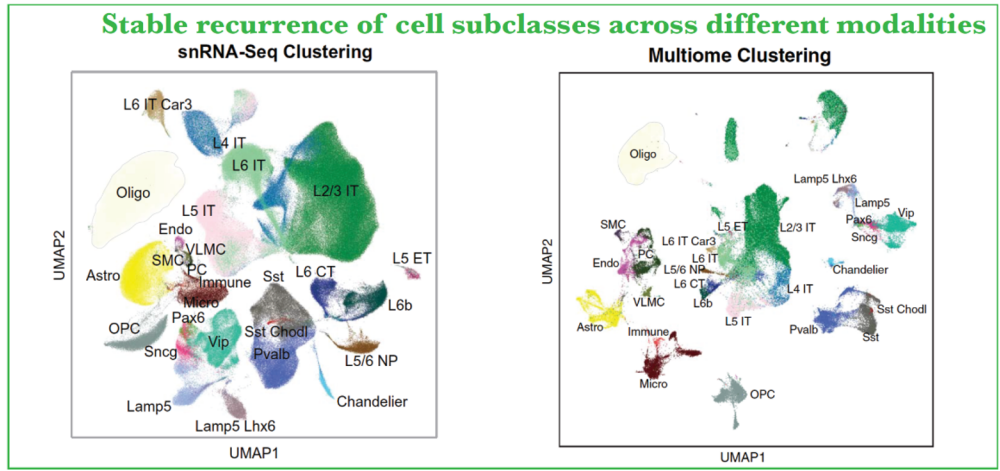
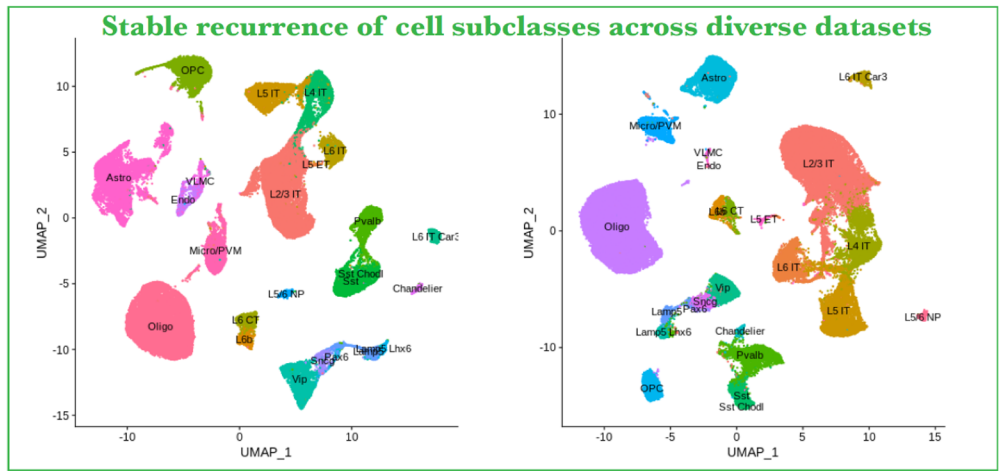
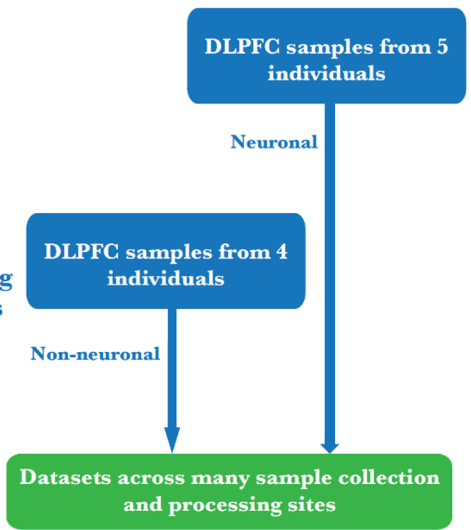


Objective. To synchronize activities across consortia for consistent DLPFC cell type definitions

BICCN:
Multi-modal
Cross-brain-region
Reference Atlas

Ma-Sestan:
Focused PFC study
with deeper sampling
of non-neuronal cells

PsychENCODE:
Validation on psy-
chiatric disorder
snRNAseq datasets

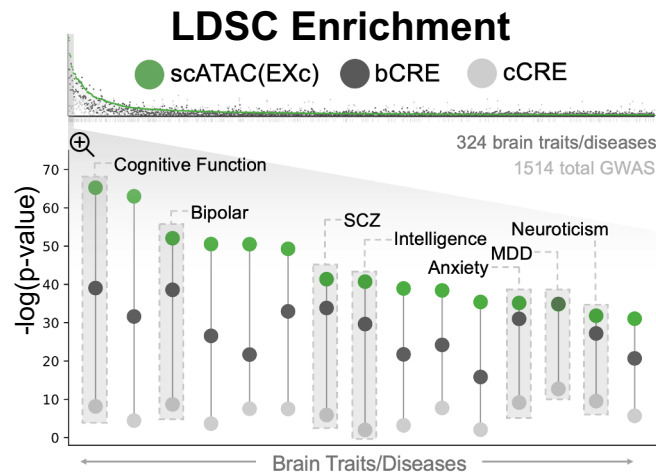
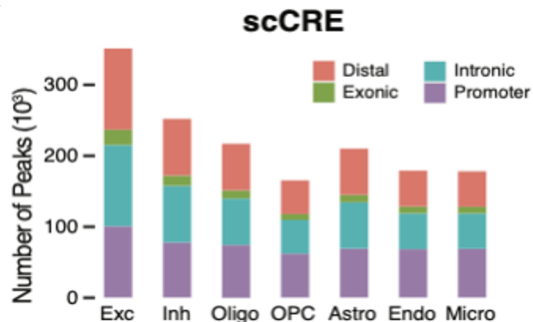


28 Canonical Cell Types

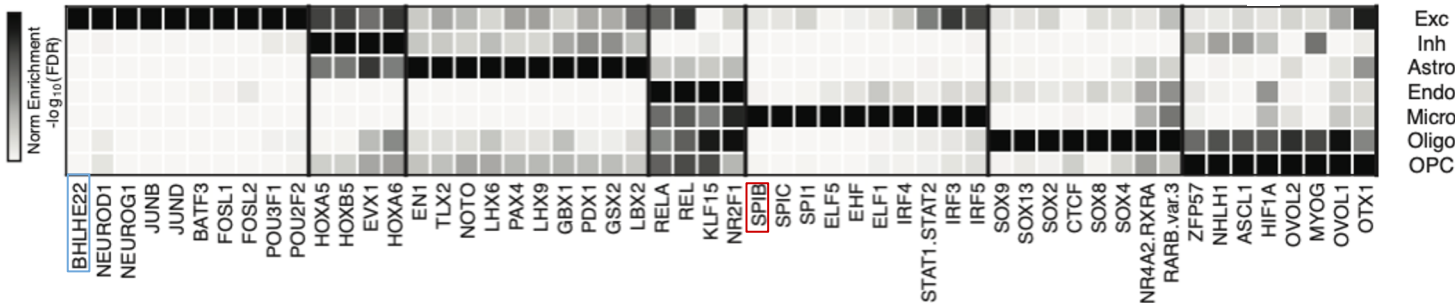
- | Excitatory | | | Inhibitory | | | Non-Neur. | | |
|------------|--------------|--------------|--------------|----------|--------------|--------------|------------|------------|
| — L2/3 IT | — L6 IT Car3 | — Sst | — Lamp5 Lhx6 | — Astro | — VLMC | — Oligo | — PC Δ | — SMC Δ |
| — L4 IT | — L5 ET | — Sst Chodl | — Lamp5 | — Endo | — VLMC | — OPC | — SMC Δ | — Immune Δ |
| — L5 IT | — L5/6 NP | — Pvalb | — Sncg | — SMC | — PC | — Micro | — Immune Δ | — RB Δ |
| — L6 IT | — L6b | — Chandelier | — Vip | — Immune | — Sst | — Micro | — Immune Δ | |
| — L6 CT | | — Pax6 | | — Sncg | — Sst Chodl | — Vip | | |
| | | | | — Pax6 | — Pvalb | — Sst Chodl | | |
| | | | | | — Lamp5 | — Pvalb | | |
| | | | | | — Lamp5 Lhx6 | — L5/6 NP | | |
| | | | | | | — Chandelier | | |
- Δ Non-BICCN cell type

scCREs show specific enrichment for TF motifs and GWAS signals

~560,000 single-cell cis-regulatory elements (scCREs) from ATAC peaks, more enriched for brain traits in GWAS than bulk cCREs



TF Motif Enrichment

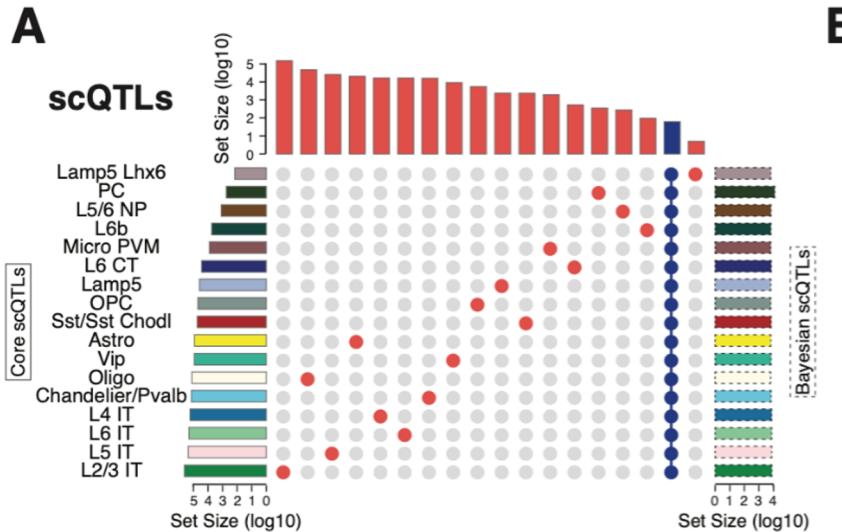


Cell-type-specific TF motif patterns - major brain cell types employ distinct groups of TFs.

Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Cell type-specific eQTLs (scQTLs)

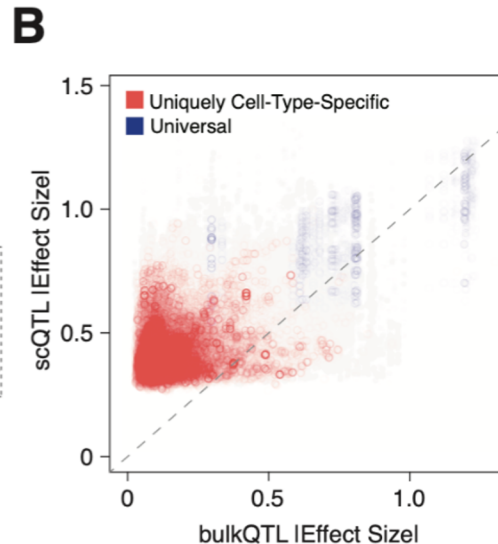


Core set of 1.4M scQTLs
based on standard GTEx QTL scheme

Overall, ~85K scQTLs & ~690 significant eGenes per cell type

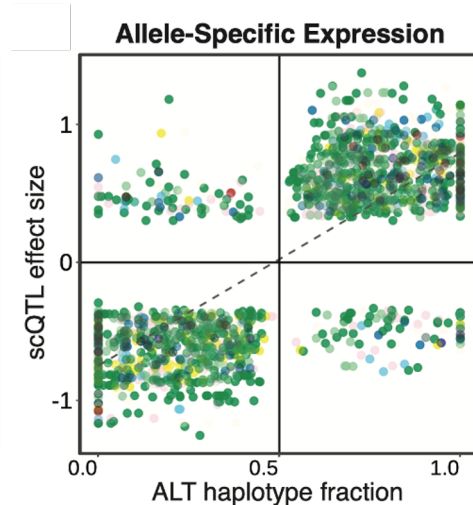
~53% scQTLs cell-type-specific

Bayesian methods sharing information between cell types can identify more scQTLs in rarer cell types



Dilution Effect
implies many
scQTLs will not be
seen in bulk

For scQTLs that overlap with
bulk eQTLs, scQTL effect is
larger & diluted out in bulk



Validation using
single-cell ASE

Consistency betw. scQTL effect
sizes and ALT haplotype fractions

Matrix gymnastics: Cross-study data integration, filtering, and matrix synchronization

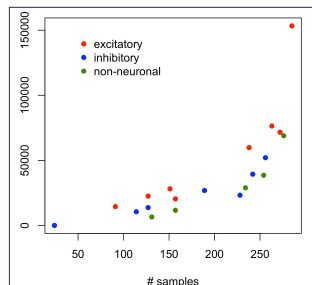
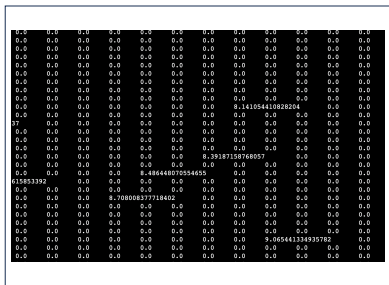
- > Filtering:
 - Genes
 - Nuclei
 - Variants
 - Individuals
- > Matrix synchronization

Challenges in Calling Cell type-specific eQTLs (scQTLs)

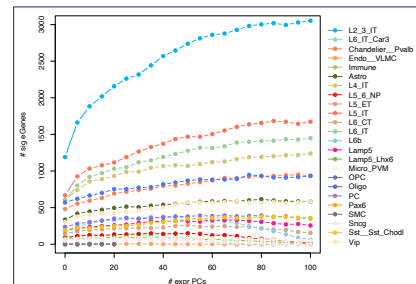
Slowly explore 'decision space' -- details re. pre-processing (e.g., expr normalization, etc)

- log TPM-normalization?
- log CPM-normalization?
- sc-transform normalization?
- TMM-normalization?
- Thresholds for # min nuclei & samples
- Stage to enforce MAF filters

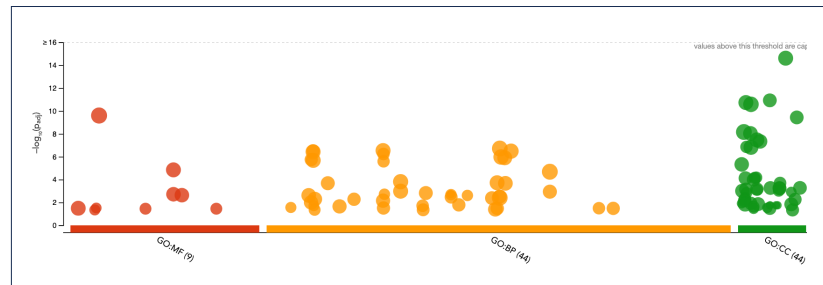
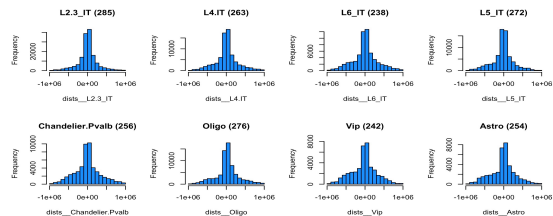
Data sparsity and limited statistical power in snRNA-seq contexts



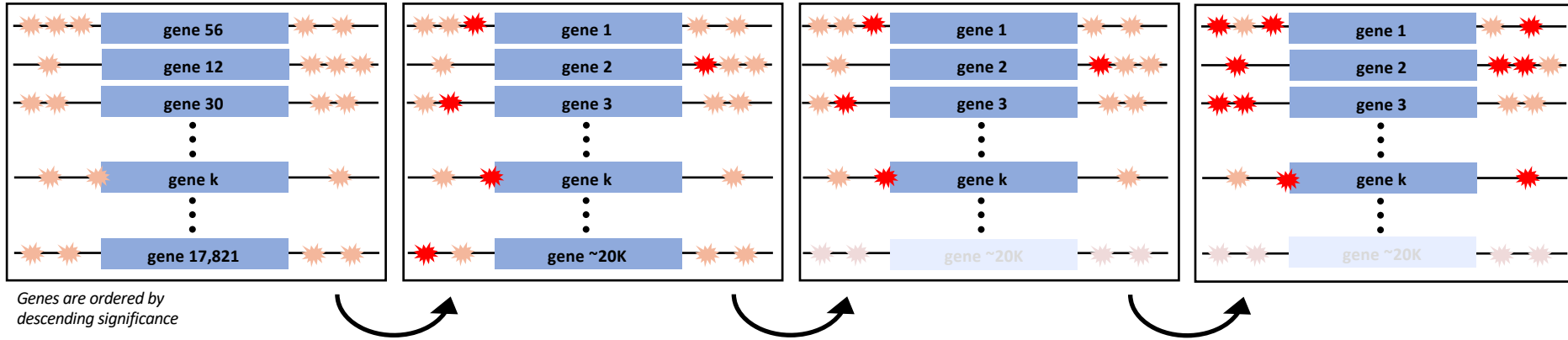
Batch effects (mult. cohorts): Optimizing calling setups (ex: selecting covariates and numbers of PCs to include)



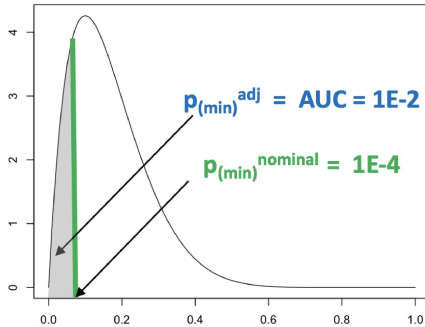
'Semi-blind' validation: Devising and performing quality checks and 'validation' without gold-standard reference dataset for comparisons



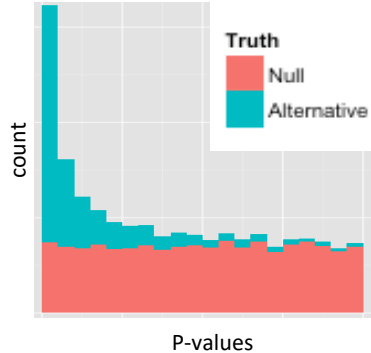
Multi-step (hierarchical) scheme to identify significant eGenes & their associated eSNPs (GTEx compatible approach)



Step 1: Identify the most significant eSNP per gene, and then correct p-values for multiple testing within each gene to derive adjusted gene-level p-values



Step 2: Multiple testing correction (BH to estimate FDR) is applied to the set of all ~20K adjusted gene-level p-values to yield the significant eGenes (FDR 0.05)



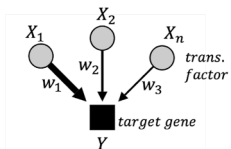
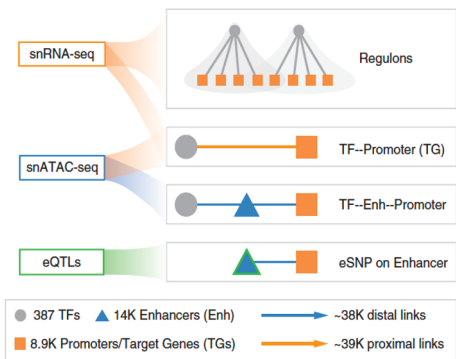
Adopted from D. Robinson (web post)

Step 3: Pull in all significant eSNPs associated with each significant eGene by using the scheme adopted by GTEx

$$p_{\text{nomimanl_thresh}}^{(\text{eGene})} = F^{-1}(p_t)$$

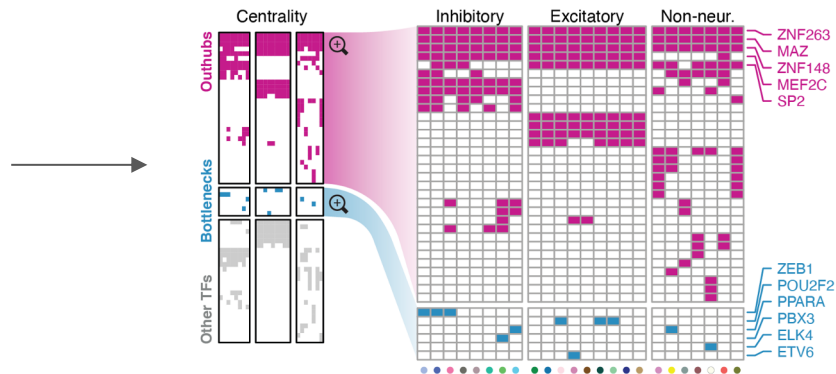
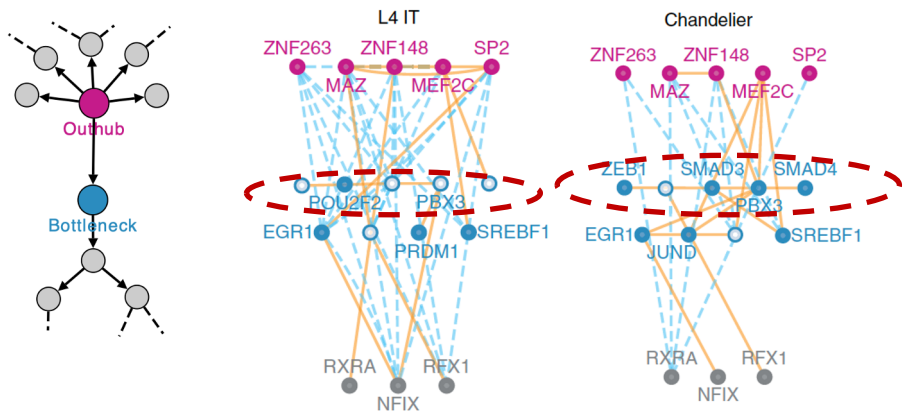
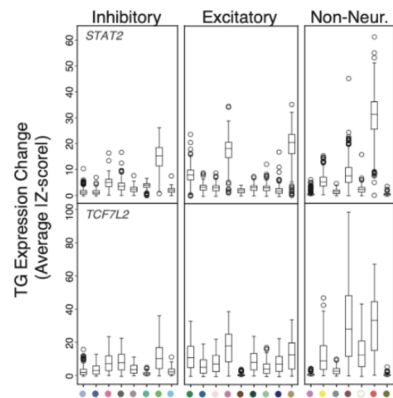
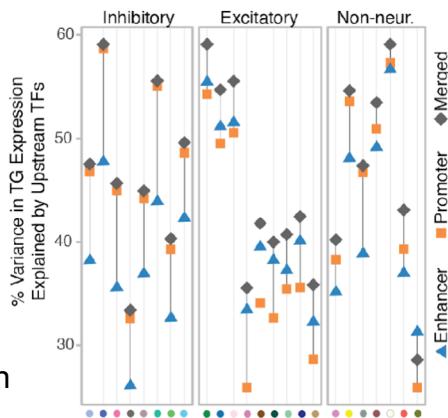
$$F(p_{\text{nomimanl_thresh}}^{(\text{eGene})}) = p_t$$

Constructing cell-type-specific gene regulatory networks



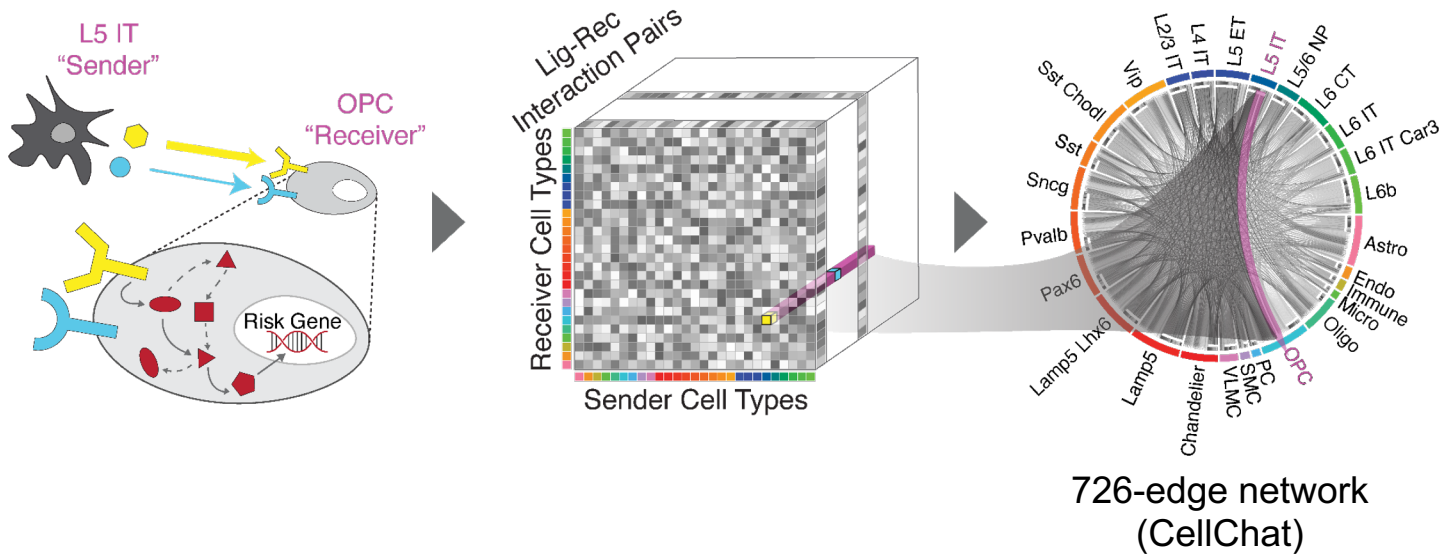
$$Y_i = \beta_0 + \sum_{j=1}^n \beta_j \cdot (X_j \cdot w_j) + \epsilon$$

TF expression explains 52% of variation in target gene expression



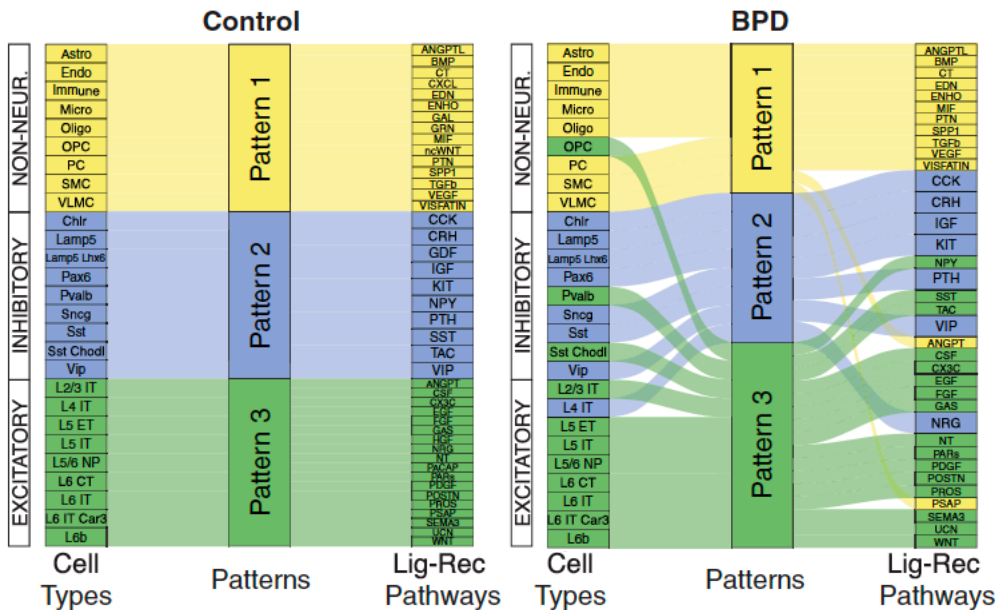
TFs show differential usage (i.e. out-hubs, bottlenecks) across cell-type GRNs

Disease-specific alterations in cell-to-cell communication

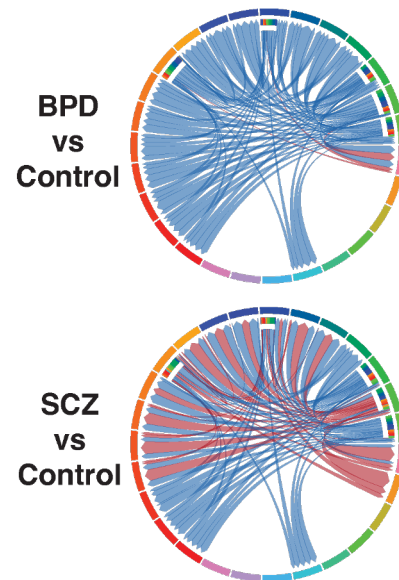


Disease-specific alterations in cell-to-cell communication

Outgoing Communication Pattern Clustering



Differential Analysis of WNT Pathway

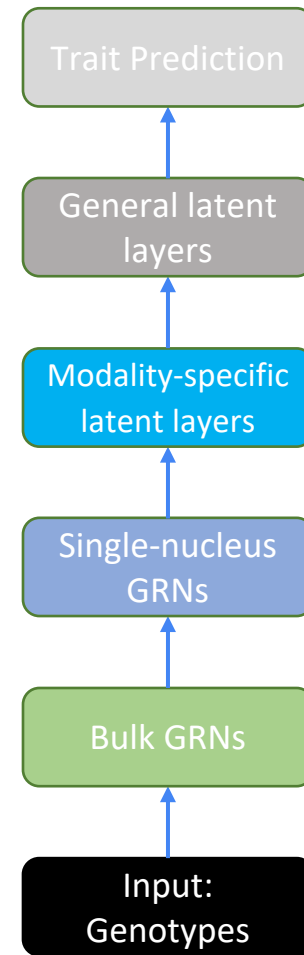
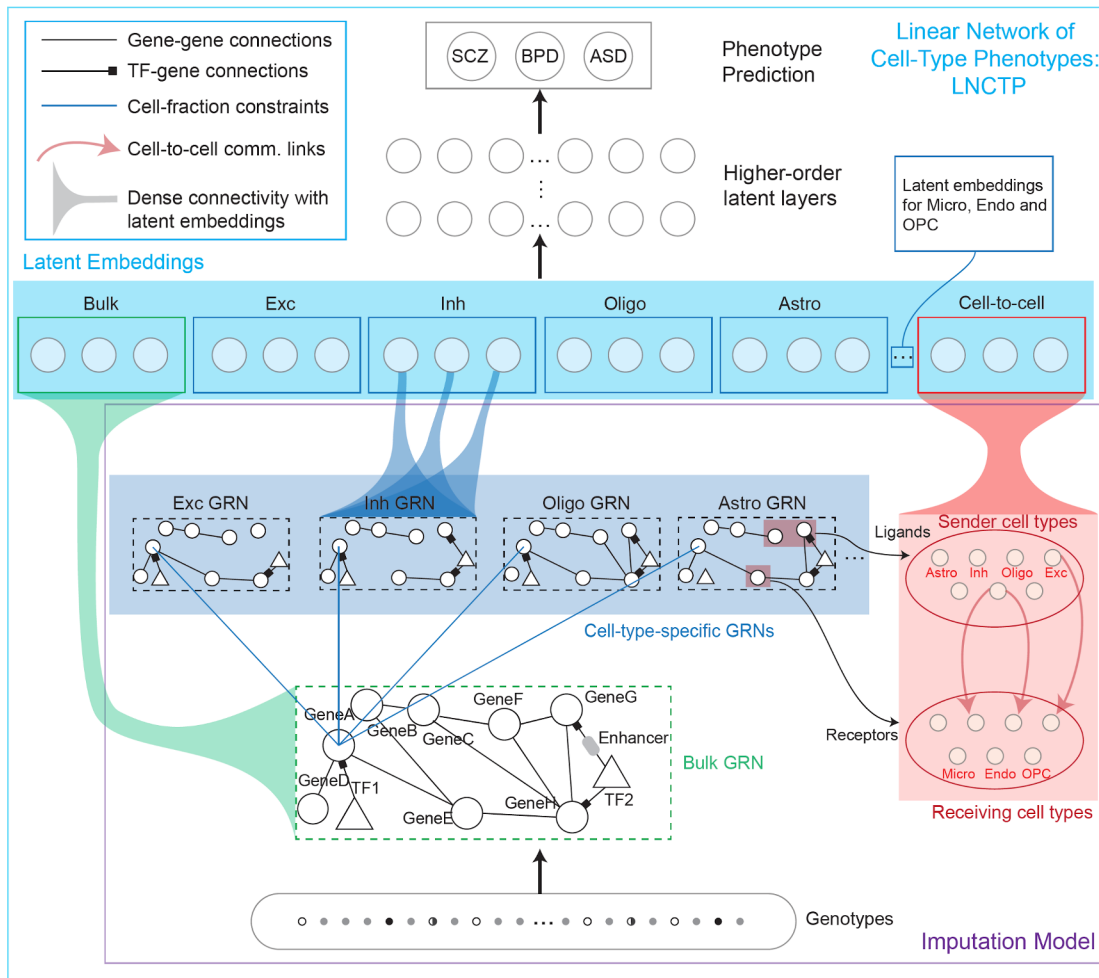


Large-scale changes in cell-cell communication patterns seen in individuals with neuropsychiatric disorders

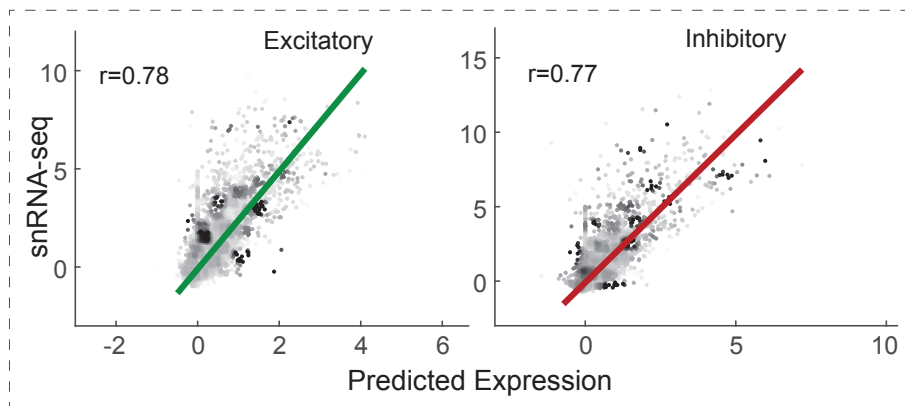
Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Linear Network of Cell-Type Phenotypes (LNCTP) model framework

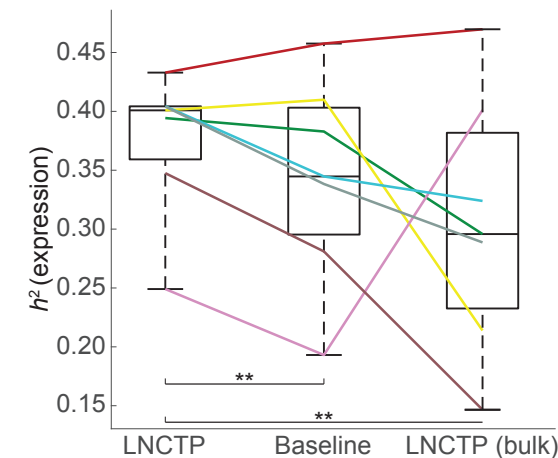


LNCTP shows improved performance for imputing expression



-Prediction of single-cell expression in samples based only on genotypes

-Improves prediction of cell-type expression variance compared with other methods (i.e. baseline or bulk RNA models, PRS)



Cell Types — Exc — Astro — Micro — OPC
— Inh — Endo — Oligo

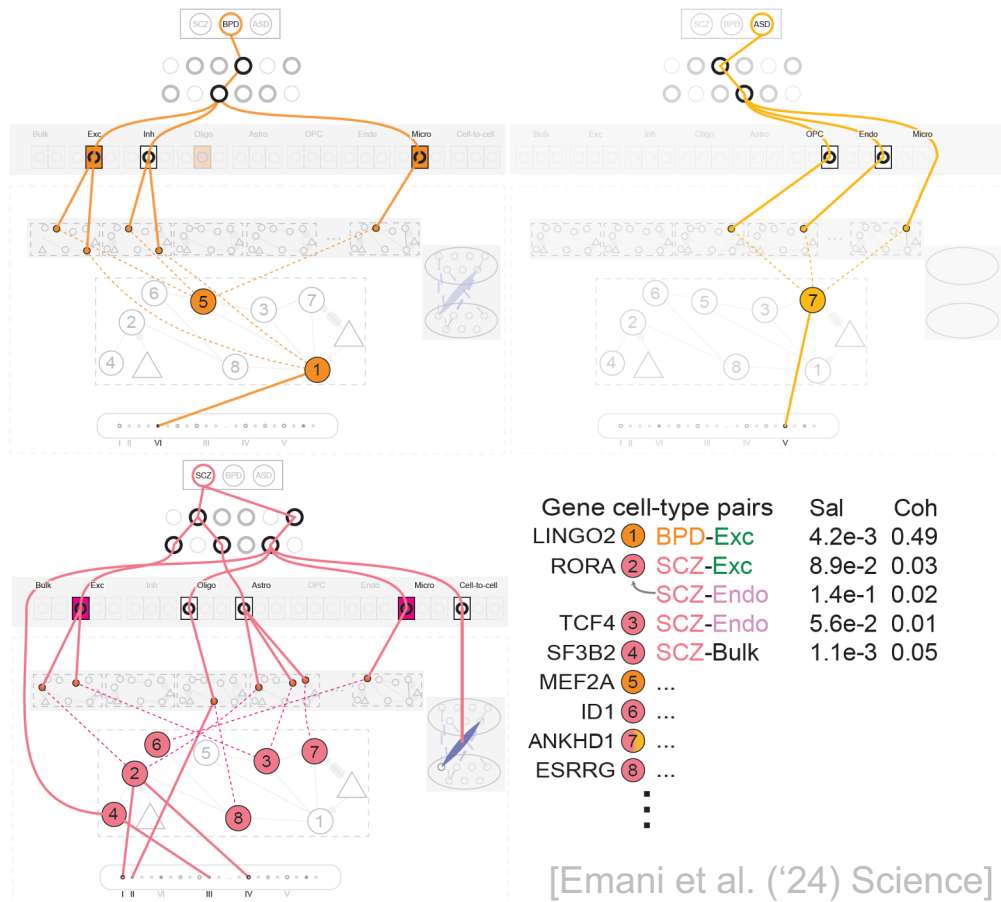
Using LNCTP to link genes, cell types, and phenotypes

Salient pathways from genes through cell types to traits

Results from LNCTP allow the association of traits with genes in a cell-type-specific manner

By tracing the influence of genes through visible and latent layers, **cell-type-specific effects** towards disease can be identified.

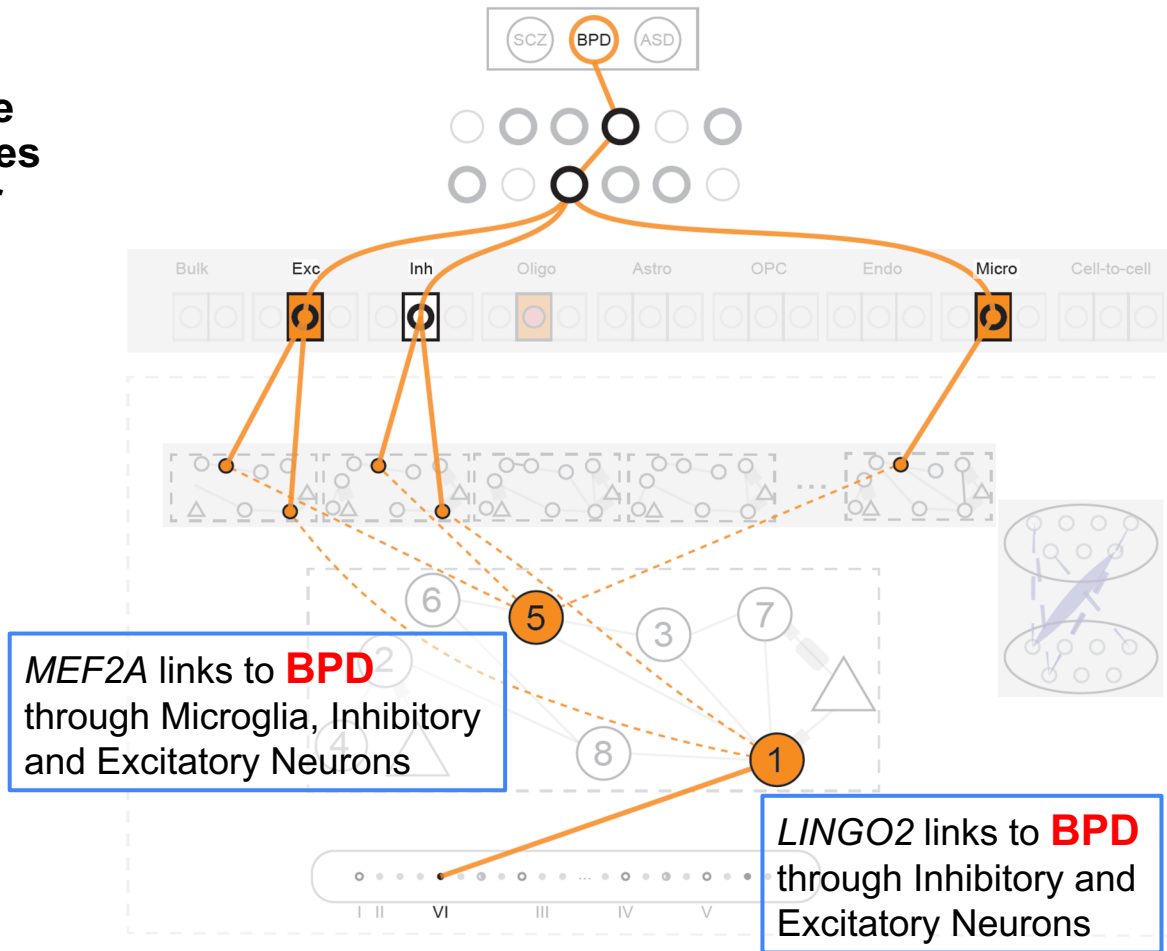
~250 total gene + cell type pairs



[Emani et al. ('24) Science]

LNCTP examples: Prioritized cell types for BPD genes

Results from LNCTP allow the association of traits with genes in a cell-type-specific manner

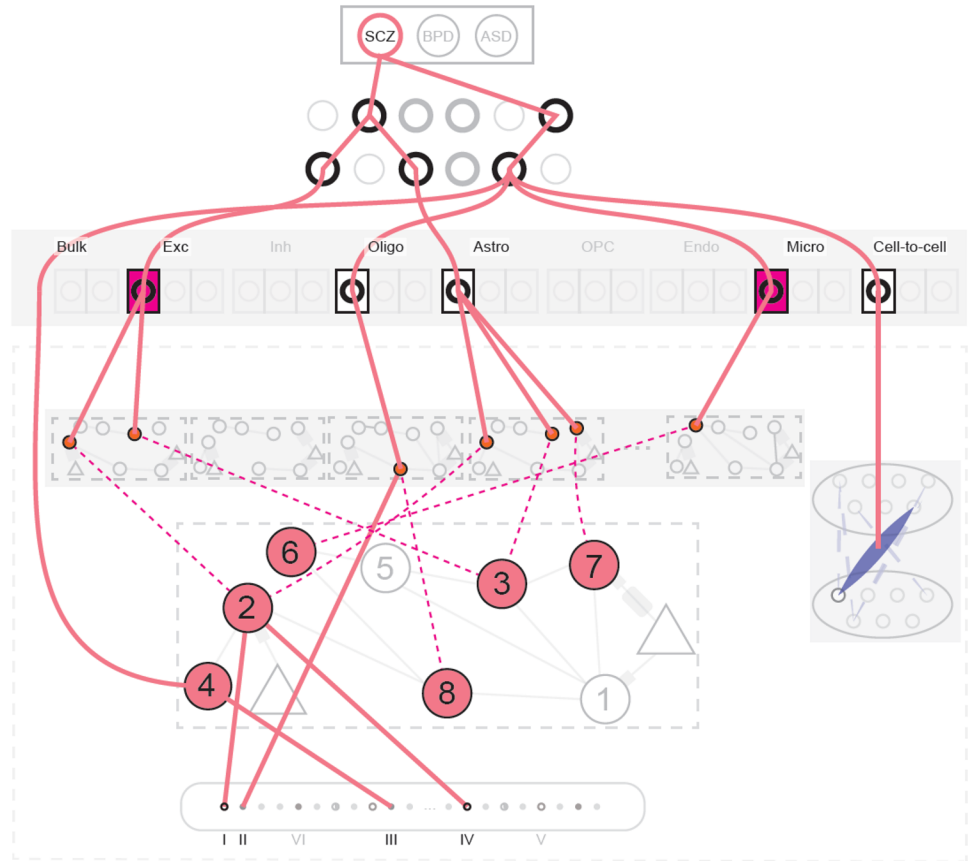


LNCTP examples: Prioritized cell types for SCZ genes

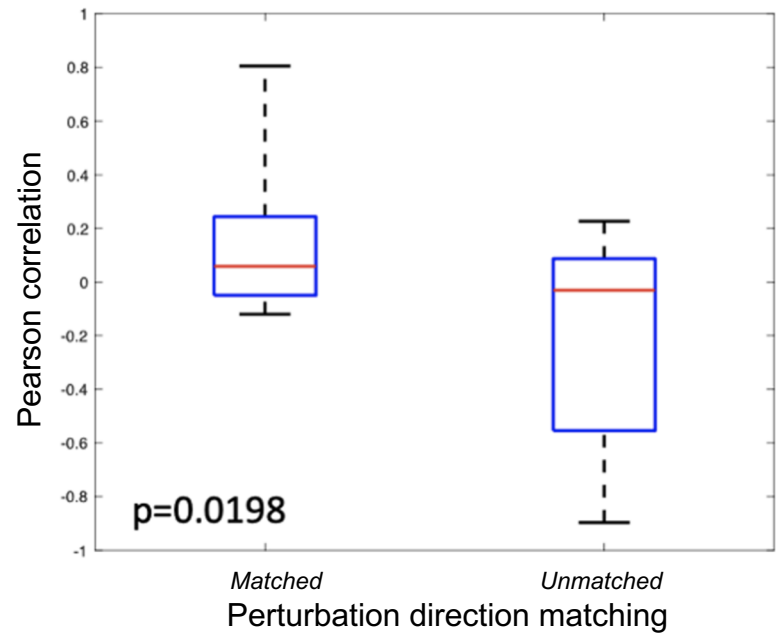
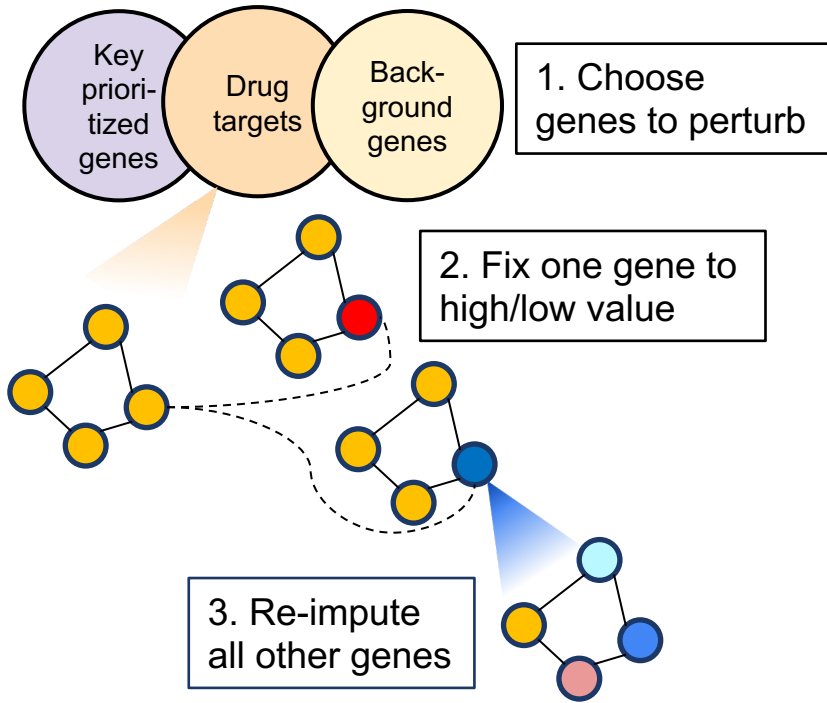
Results from LNCTP allow the association of traits with genes in a cell-type-specific manner

Highlights in **SCZ** include *TCF4*, *RORA*, & Micro-Exc linkage in cell-to-cell network

RORA	②	SCZ-Exc	8.9e-2	0.03
		SCZ-Endo	1.4e-1	0.02
TCF4	③	SCZ-Endo	5.6e-2	0.01
SF3B2	④	SCZ-Bulk	1.1e-3	0.05
MEF2A	⑤	...		
ID1	⑥	...		
ANKHD1	⑦	...		
ESRRG	⑧	...		



Comparison of LNCTP predicted network to Glu neuron CRISPR experiments*



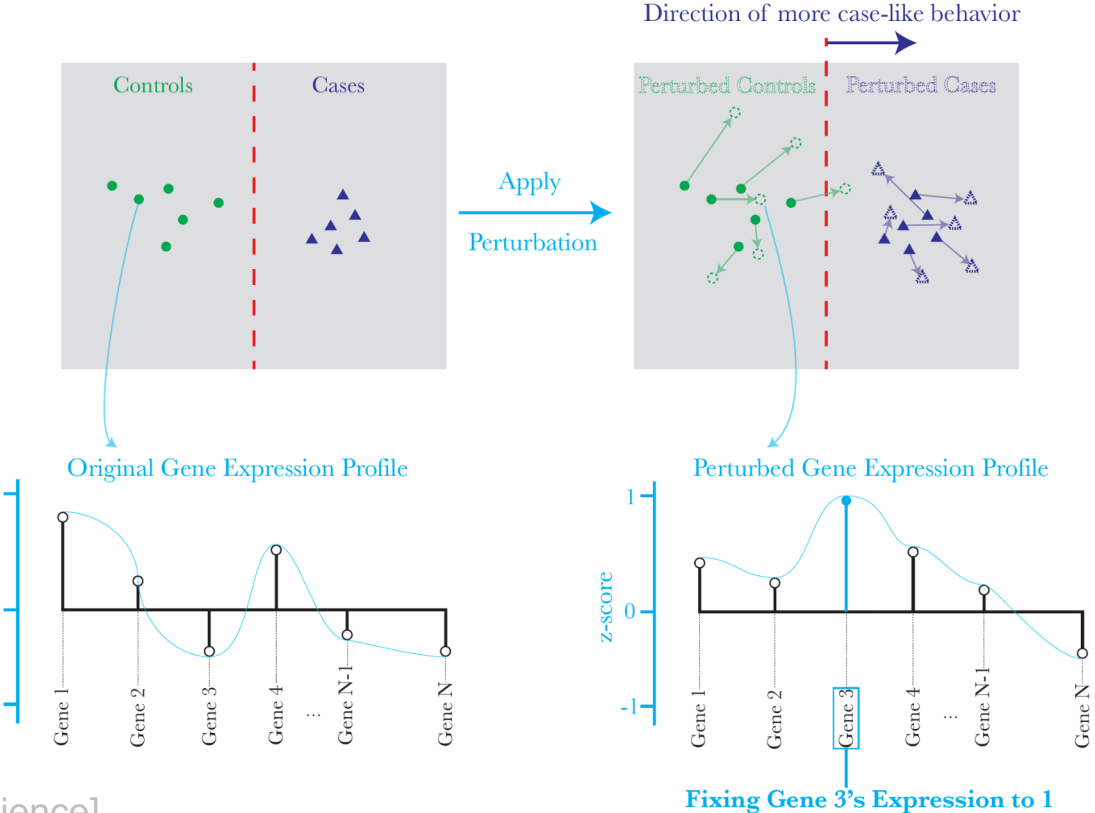
Comparing LNCTP and CRISPR perturbations

- Perturbations in excitatory neuron GRN
- Upper decile of genes according to LNCTP z-score changes
- Perturbation directions are **matched** or **unmatched**
 - **Unmatched** means LNCTP z-score changes correlate with CRISPR fold-change vectors for all genes **except** the perturbed gene

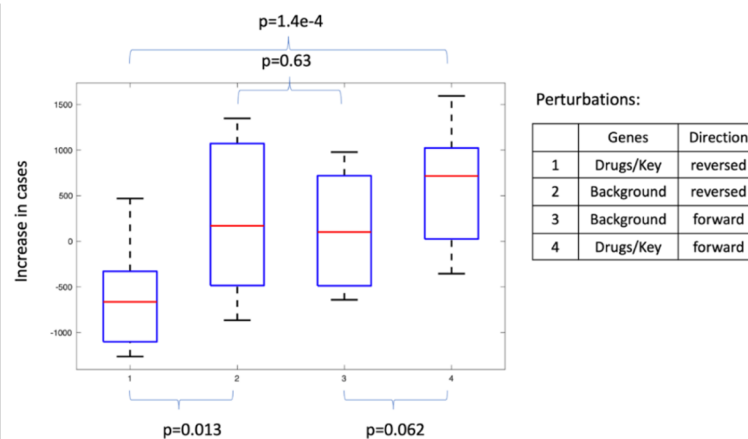
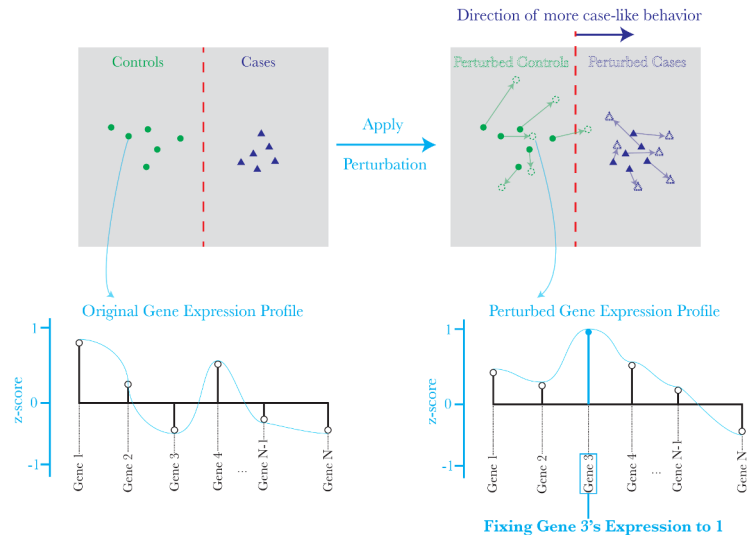
*Tian, R., Abarientos, A., Hong, J. *et al. Nat Neurosci* 24, 1020–1034 (2021).

[Emani et al. ('24) Science]

LNCTP model: Perturbation Analysis



LNCTP model: Perturbation Analysis



LNCTP predicts effects of perturbations on case/control status

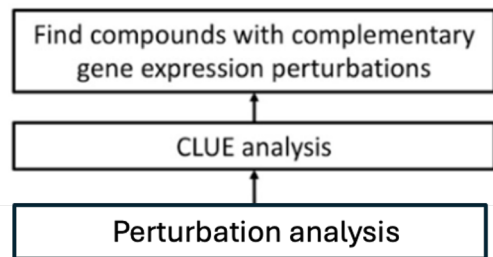
$$p_{LNCTP-perturbed}(y_i, h_i, x_{i,\neg(c^*,g^*)} | z_i, x_{ic^*g^*} = \{1, -1\}) \propto p_{GMRP}(x_{i,\neg(c^*,g^*)} | z_i, x_{ic^*g^*} = \{1, -1\}) \cdot p_{DNN}(y_i, h_i | x_i)$$

$$p_{GMRP}(x_{i,\neg(c^*,g^*)} | z_i, x_{ic^*g^*} = \{1, -1\}) \propto \exp(-E_{GMRP}(x_{i,\neg(c^*,g^*)} | z_i, x_{ic^*g^*} = \{1, -1\}))$$

$$E_{GMRP}(x_{i,\neg(c^*,g^*)} | z_i, x_{ic^*g^*} = \{1, -1\}) = x_{i0}^T J x_{i0} + \sum_g x_{i0g}^T b(z, \beta_g) + \sum_c (x_{ic}^T J_c x_{ic} + x_{ic}^T b_c + x_{ic}^T J_c^{(c2c)} x_{i,LR,c}) +$$

$$\sum_{c_1, c_2} J_{c_1, c_2}^{(c2c)} x_{i, L, c_1}^{(c2c)} x_{i, R, c_2}^{(c2c)} + \lambda \sum_g (x_{i0g} - f(z)^T x_{i, 1 \dots C, g})^2 + \infty \cdot \delta(x_{ic^*g^*} \neq \{1, -1\})$$

LNCTP model: Clue.io Analysis



Well-known drugs used for neuropsychiatric disorders:

dopamine receptor antagonists,
dopamine receptor agonists,
glutamate receptor antagonists,
calcium channel blockers,
GABA receptor agonists,
MAP kinase inhibitors

Compounds with unknown effects:

Cytokine IL-1a: potential in reversing the expression changes of the ID1 gene in microglia
AKT inhibitor 10-DEBC: potential reversing the effects of TCF4, ID1, RORA, SF3B2
Consistent occurrence of bromodomain inhibitors for reversing effects of all eight genes

Number of significant compounds in CLUE database

Gene	Bulk	Astro	Endo	Exc	Inh	Micro	Oligo	OPC
<i>ANKHD1</i>	797	234	15	405	375	13	1808	253
<i>ESRRG</i>	184	52	105	403	89	413	325	206
<i>ID1</i>	400	402	756	389	119	9073	594	1005
<i>LINGO2</i>	431	791	346	25	16	1688	492	1454
<i>MEF2A</i>	335	283	118	8	11	2670	661	349
<i>RORA</i>	242	196	301	156	79	1289	673	4605
<i>SF3B2</i>	176	100	1101	421	1014	4688	774	722
<i>TCF4</i>	489	185	7	253	96	1063	347	385

Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

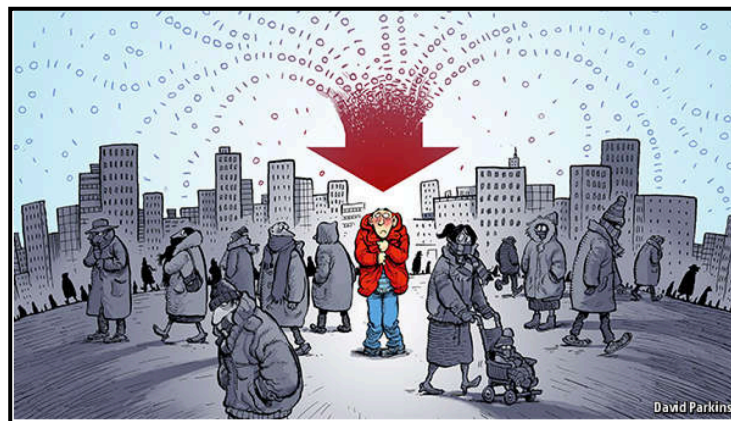


The Other Side of the Coin for Genomics: Why we should share

- Sharing helps **speed research**
 - Large-scale mining of this information is important for medical research
 - Statistical power
 - Privacy is cumbersome, particularly for big data

The Dilemma

- The individual (harmed?) v the collective (benefits)
 - But do sick patients care about their privacy?
- How to balance risks v rewards
 - Quantification



[Economist, 15 Aug '15]

[Yale Law Roundtable ('10). Comp. in Sci. & Eng. 12:8; D Greenbaum & M Gerstein ('09). Am. J. Bioethics; D Greenbaum & M Gerstein ('10). SF Chronicle, May 2, Page E-4; Greenbaum et al. PLOS CB ('11)]

Privacy: Does Genomics has similar "Big Data" Dilemma as in the Rest of Society?

- We confront privacy risks every day we access the internet (e.g., social media, e-commerce).
- Sharing & "peer-production" is central to success of many new ventures, with analogous risks to genomics
 - **EG web search:** Large-scale mining essential



Genetic Exceptionalism :

The Genome is very fundamental data, potentially very revealing about one's identity & characteristics

Personal Genomic info. essentially meaningless currently but will it be in 20 yrs? 50 yrs?

Genomic sequence very revealing about one's children.

Is true consent possible?

Once put on the web it can't be taken back

Ethically challenged history of genetics

Ownership of the data & what consent means (Hela)

Could your genetic data give rise to a product line?

Risk from environmental samples: sparse and noisy genotypes

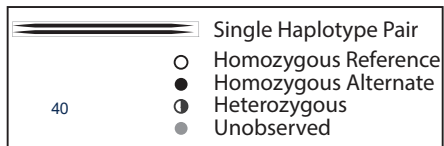
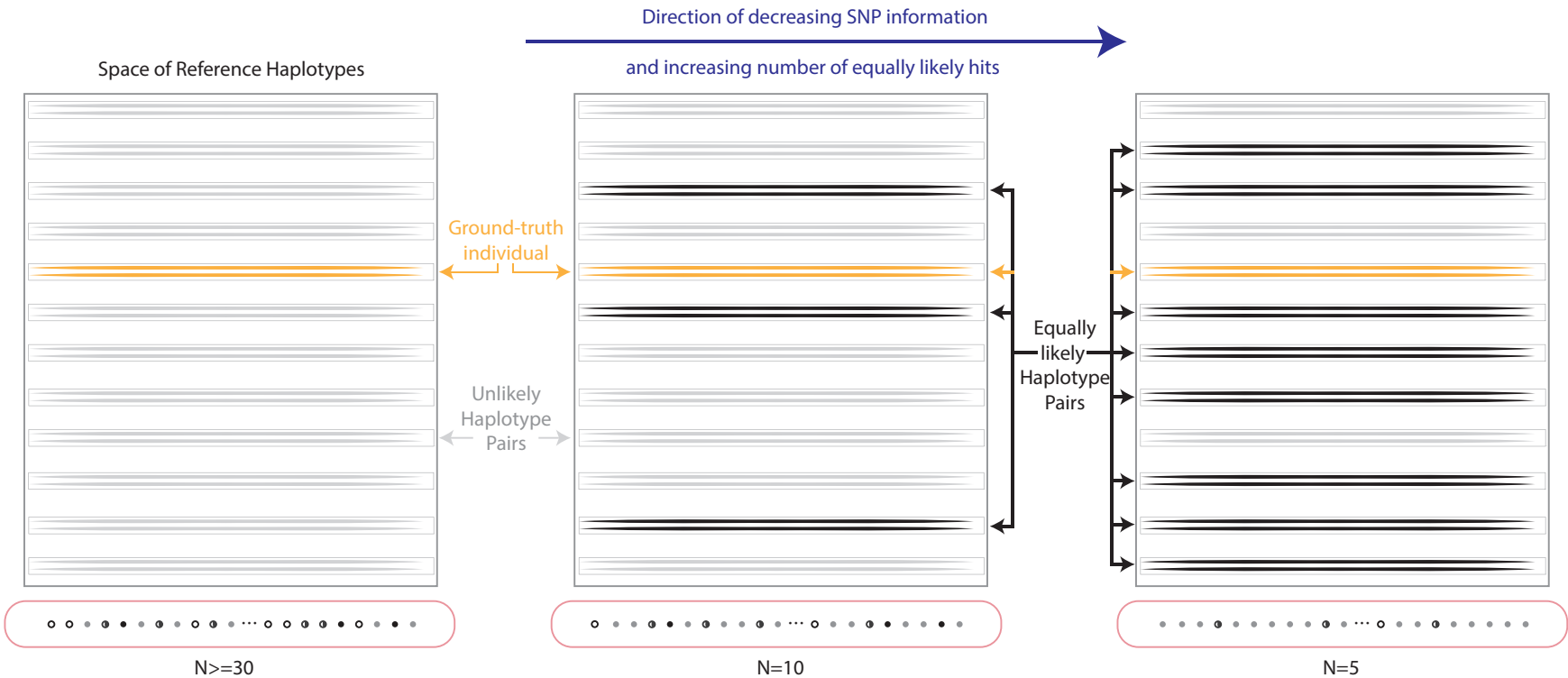
- Swabs of objects are easy to obtain
- How much information is contained in such DNA samples?
- Created a population-genetics-inspired approach* to quantify the risk:

PLIGHT = Privacy Leakage through Inference across Genotypic HMM Trajectories

- Found risk of identification is high even with tens of noisy SNPs
- Identified a way of sanitizing data before publication

*Emani, P.S.; Geradi, M.N.; Gürsoy, G.; Grasty, M.R.; Miranker, A.; Gerstein, M.B. "Assessing and mitigating privacy risks of sparse, noisy genotypes by local alignment to haplotype databases", *Genome Research* (2023), Vol. 33, Iss. 12, Pgs. 2156-2173

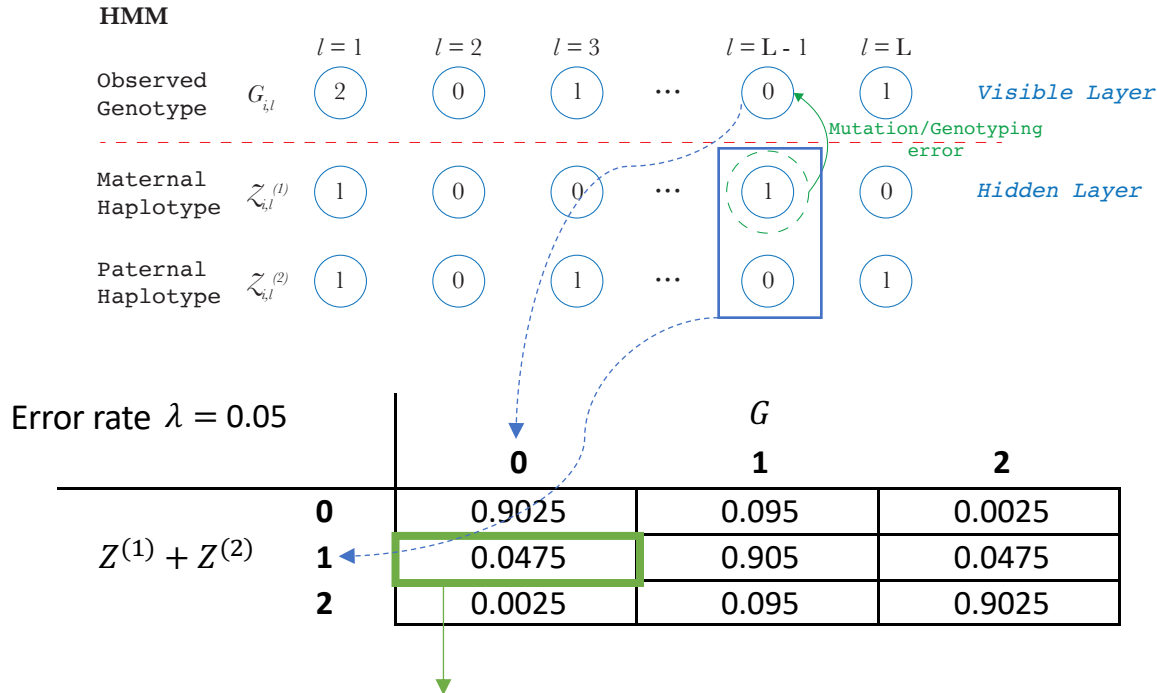
Quantifying identifying information in the limit of small SNP sets



The Li & Stephens HMM: Effect of mutation/genotyping error

(Li, N.; Stephens, M. *Genetics* **2003**, *165*, 2213–2233.)

When identifying individuals, can include the effects of genotyping error/noise using a population genetics approach:

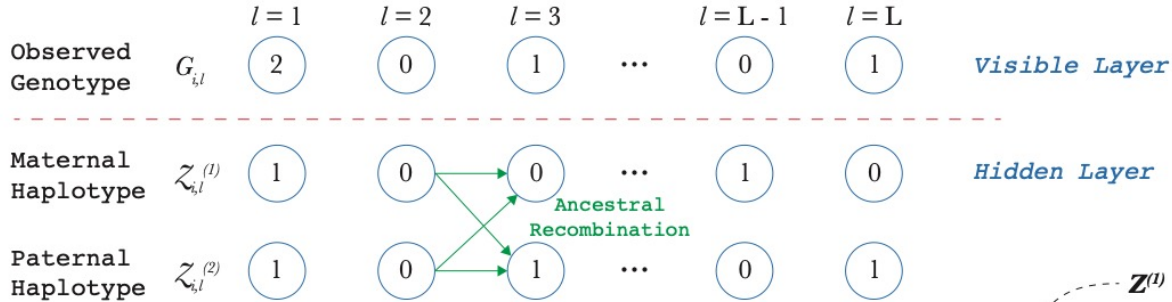


Probability of observing a genotype G of 0, conditional on the sum of the haplotypes being 1

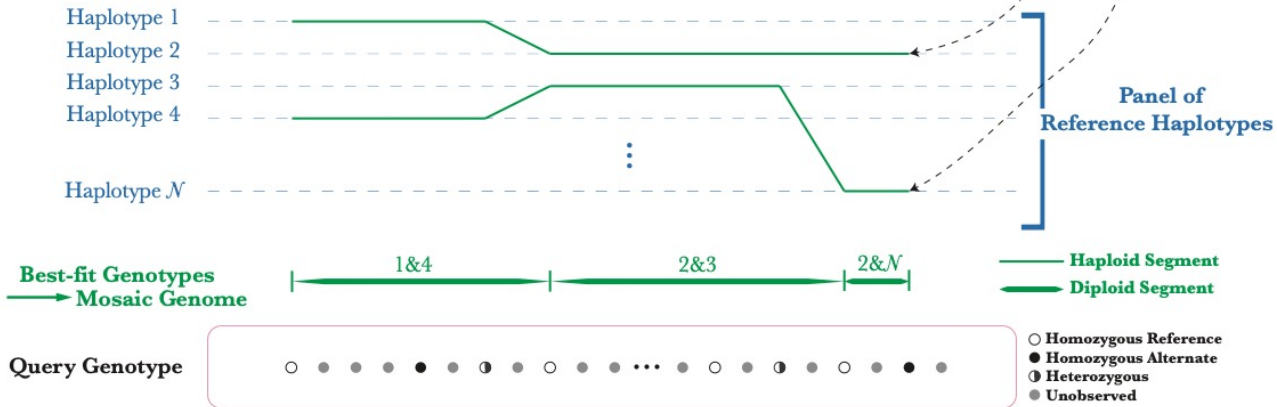
The Li & Stephens HMM: Effect of ancestral recombination

(Li, N.; Stephens, M. *Genetics* **2003**, *165*, 2213–2233.)

HMM with recombination

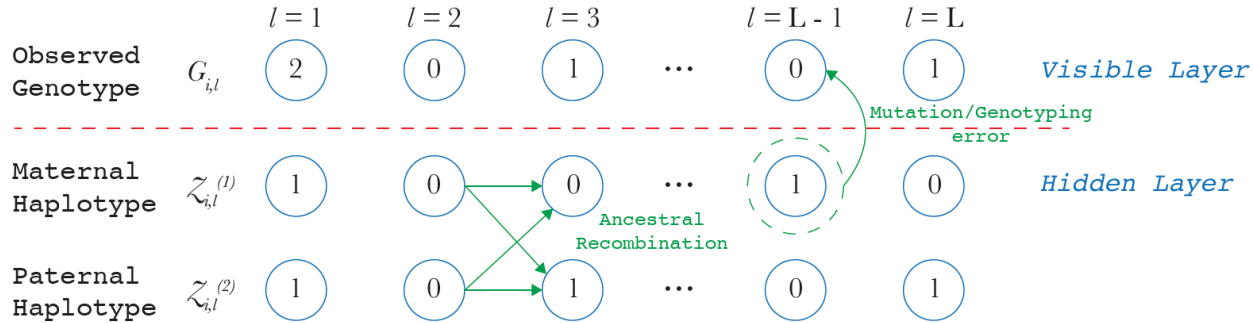


Haplotype Inference: Construction of One Best-fit Piecewise/Mosaic Genome



The Li & Stephens HMM

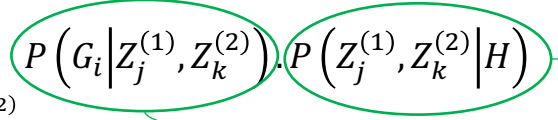
(Li, N.; Stephens, M. *Genetics* **2003**, *165*, 2213–2233.)



Addresses: What is the probability of observing a set of genotypes, based on underlying panel of haplotypes?

$$P(G_i|H)$$

$$= \sum_{Z_j^{(1)}, Z_k^{(2)}} P(G_i|Z_j^{(1)}, Z_k^{(2)}) \cdot P(Z_j^{(1)}, Z_k^{(2)}|H)$$



Encodes recombination

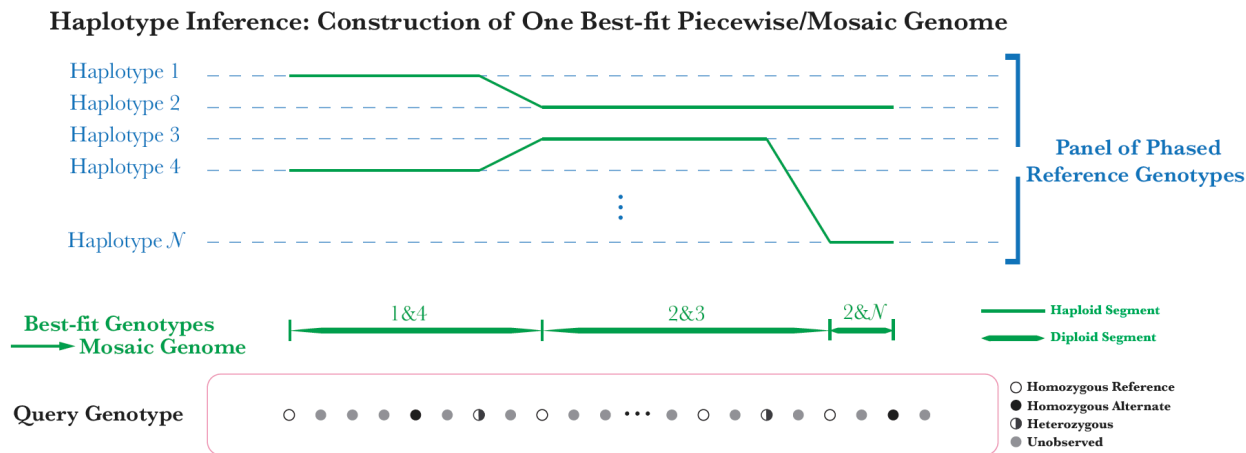
Encodes mutations or genotyping error

H = Set of all reference phased genotypes

$Z_j^{(\alpha)} = \{Z_{j(l),l}^{(\alpha)}\}_{l=1}^L$ = Set of all possible haplotypes at the observed loci l

PLIGHT

Viterbi algorithm for most-likely path



1. Search through the space of haplotype pairs (total dimension = $N \times N$) to match the diploid genotypes
2. Each of the two haplotypes can independently recombine with other haplotypes
3. Allow for mutations/genotyping error
4. **Result:** Piecewise matches of reference haplotypes to observed genotype

PLIGHT

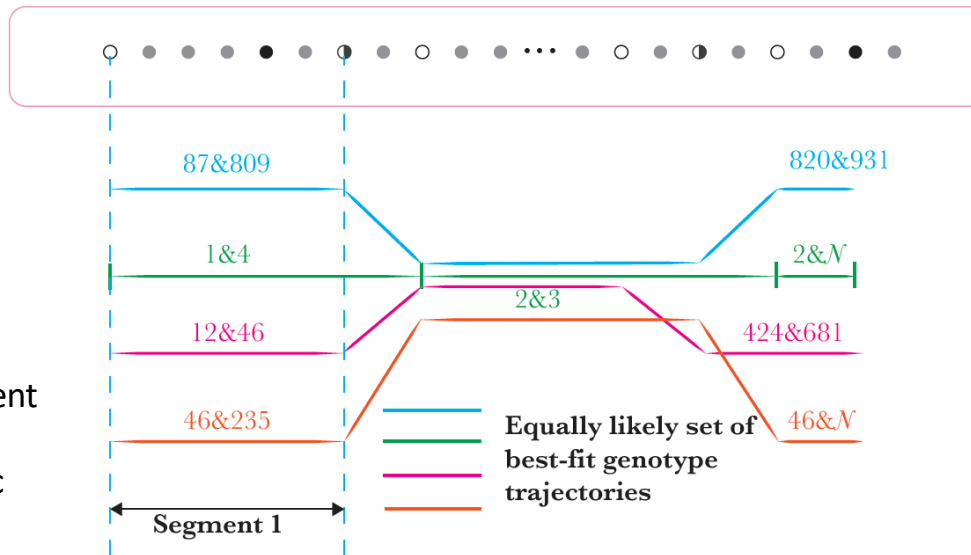
Viterbi algorithm for all equally likely “best-fit” paths

$$\operatorname{argmax}_{Z_j^{(1)}, Z_k^{(2)} \in \text{Set of all possible haplotypes}} P\left(G_i \mid Z_j^{(1)}, Z_k^{(2)}\right) \cdot P\left(Z_j^{(1)}, Z_k^{(2)} \mid H\right)$$

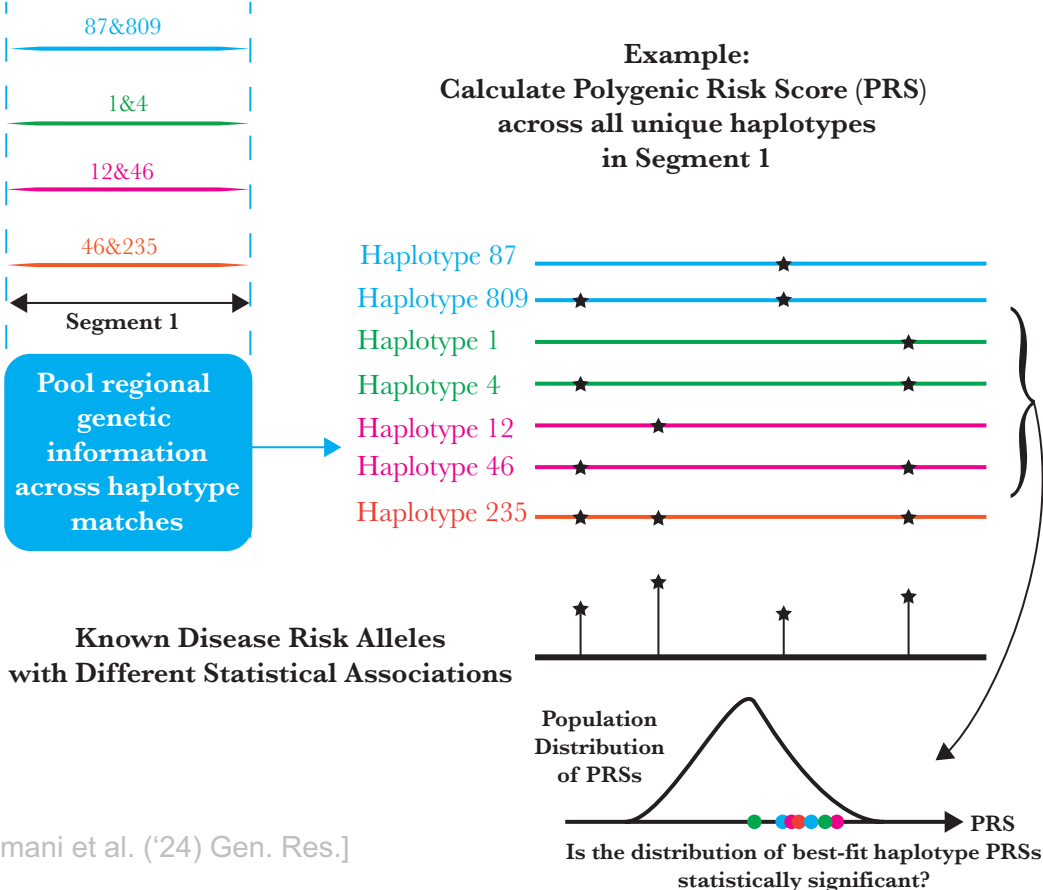
For very sparse data, several paths may be equally likely.

We term these paths ‘**Genotypic Trajectories**’:
to signify the sequential exploration of haplotype space.

The number of independent trajectories in any region gives a sense of genotypic ‘entropy’

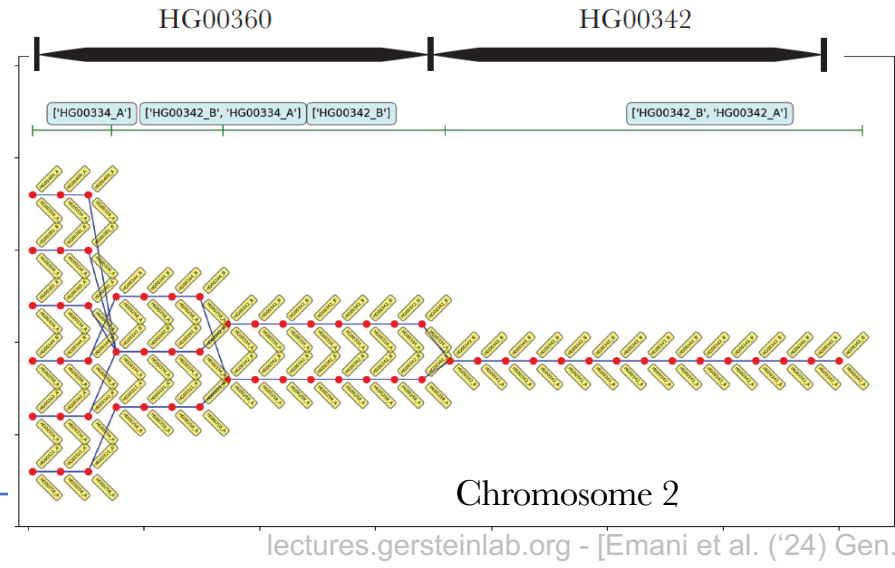
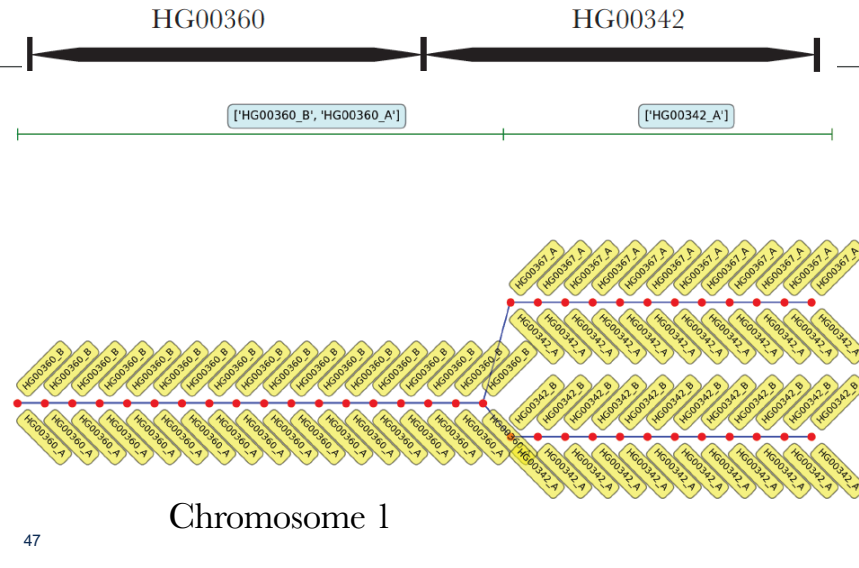


Privacy risk of partial/regional genotype matching



A few noisy SNPs can pose privacy risks

- Identified individuals with as little as 10 SNPs; robust to modest noise
- Found 1st-degree relatives (parents, siblings, children) with ~20-30 SNPs
- Environmental (saliva) swab SNPs still allow identification (with ~30 SNPs), in spite of noise
- Recommended sanitizing SNPs that specify identity with high degree of certainty



Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Overall Problem: Finding Key Variants in Personal Genomes

Millions of variants in a personal genome
Thousands, in a cancer genome
Different **contexts** for prioritization

In **rare disease**, only a few
high-impact variants are associated with disease

In **cancer**, a few positively selected drivers amongst many passengers

In **common disease**, more variants associated & each has weaker effect,
But one wants to find key "functional" variant amongst many in LD



Overall Problem: Finding Key Variants in Personal Genomes

Millions of variants in a personal genome
Thousands, in a cancer genome
Different **contexts** for prioritization

In **rare disease**, only a few
high-impact variants are associated with disease

In **cancer**, a few positively selected drivers amongst many passengers

In **common disease**, more variants associated & each has weaker effect,
But one wants to find key "functional" variant amongst many in LD

**Thus: Need to find & prioritize high impact variants.
Particularly hard for non-coding regions.**

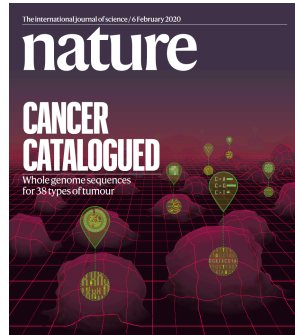
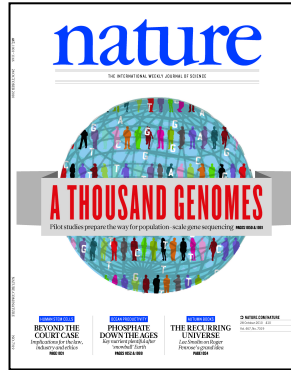


Basic Science → Medicine: Creation of Variant Catalogs of healthy & sick people

Resources for Variant Interpretation

More direct & clear interpretation through **molecular endophenotypes** (gene expression) rather than **"macro phenotype"** (disease diagnosis)

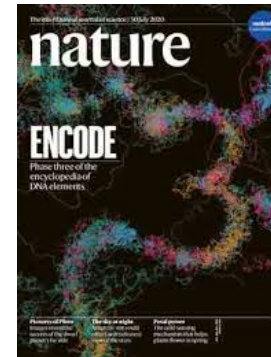
CLASSICS



INITIATIVES



STARTUPS



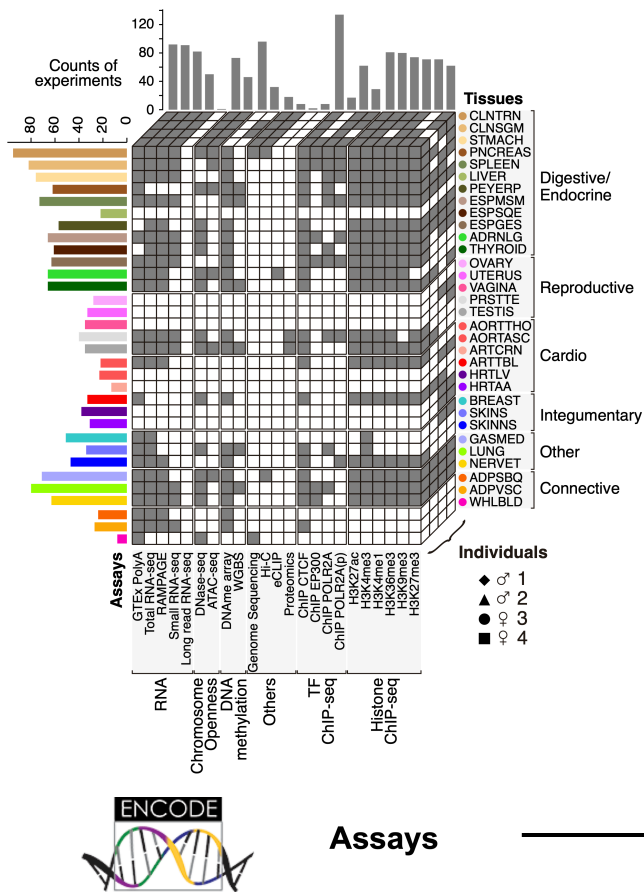
Learning the DNA regulatory grammar

CCAAT Binding Factor		
TACCCGTCAGTCTTCTAATC	TCAGCAAATCAC TTCCATGAGTTGGTGTAGAAGGTGTCAGGGCTGACTAG	490
ATGGGCAGTCAGAAGATTAGAGTCGTTTAGTGAAGGTACTCAACCACATCTTCCACAGTCCCGACTGATC		490
Selenocysteine tRNA Activating Factor		
GGTCCCTGCTGTGAGTGTGGGCACAACTCCAGG	CAC TTCCCCATCTCCAGCCGGCCAGCGAGGTGCCT	560
CCAGGGACGACACTCA CACCCCGTGTGAGTCCGTGAAGGGGTAGAGGGTCGGCCGGTTCGCTCCACGGA		560
ZIC-family		
Ubiquitous GLI- Krueppel Like Zinc Finger		
GGAGAGCCAGGCCAGATGTTCTGGCCAGGGCTGACGTACGTTTGCATGACGT CACCAG CGCTGACGAA		630
CCTCTCGGGT CCGGTCTACAAGACC GGTCCCGACTGC AGTCAAAC GTACTGCAGTGGTCCGGACTGCTT		630
RP58 Zinc Finger Protein cAMP-Responsive Element Binding Protein		
Octamer binding protein		
CCACACAGAGAGGGGCTGTCACTCACCTCAGAGCGCAGGCAACTTCTTAGAGCGGGAACTGGCCGCC		700
GGTGTCTCTCCCCGGACAGTCGAGTGGAGTCTCGCGTCCGTTGAAGAATCTCGCCCTTGACCGGCGG		700
Heat Shock Factor		
TCCTATTCTACCTTG ACGGTCCGAA TAGTCAGAGCTACAGGAGCCGAAGCGCCGAAGGAGCTCAGCAAC		770
AGGATAAGATGGAAC TGCCAGGCTT ATCAGTCTCGAT GTCTCCGGCTT CGCGGCTTCCTCGAGTCGTTGG		770
Vertebrate SMAD Family TF		
Krueppel Like Factor		
GCTACCTCCCTCTTCCCTTG GCCCTTATCCGCGTCTTCTCT TAAAGCCCAAGGAGCGGTCCCGGCAGCGGG		840
CGATGGAGGGAGAAAGGAAC CGGGAATAGGCGC AGAAGAGATTCGGGTCTTCGCCAGGGCCGTCGCCGC		840
GATA Binding Factor		
Vertebrate Homologues of Enhancer Split Complex 1 Nuclear Respiratory Factor 1 Human ETS1 Factor		
CGCTCACGAGTCCGCCGGTGGAGGCGGT CACGCCACGCACGC TTGTGCGCAAGCCCGTTCCCT CAGAAC		910
GCGAGTGTCTCAGGCGGCCACCTCCGCCAGT GCGCGT CGCTGCGA ACACGCGTTCGGGCAAGGGAGTCTTG		910
EGR/nerve growth factor induced protein C		
Myc Associated Zinc Fingers Vertebrate Homologue of Enhancer Split Complex 2		
CCGGAAGTCCGAGGC CGAGG A GAGCCACGTGGGT CGAGCT GGGGCGCAGCTCGC ATG GGGGAGTCTATCC		980
GGCCTTCACGCTCCGGCTCCTCTCGGTGCACCCAGCTCGACCCCGCTCGAGCGTACCCCTCAGATAGG		980

Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

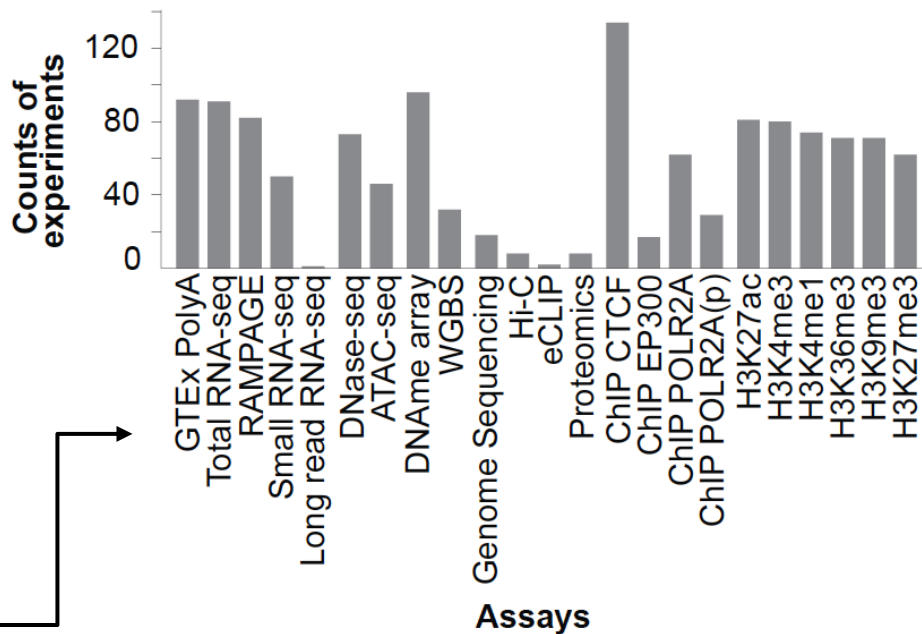
- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

The EN-TECH resource of multi-tissue experimental functional genomics data

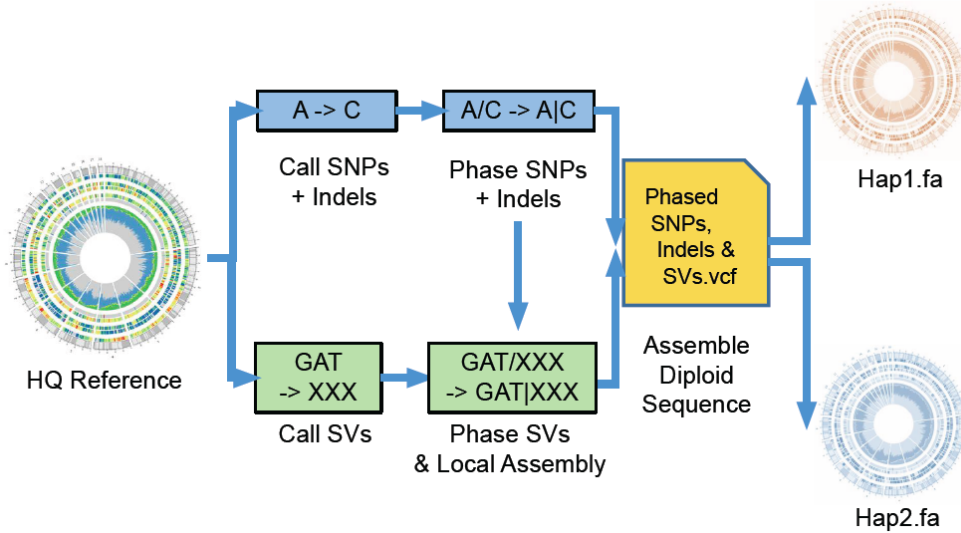


Tissues
&
Individuals

~1500 experiments

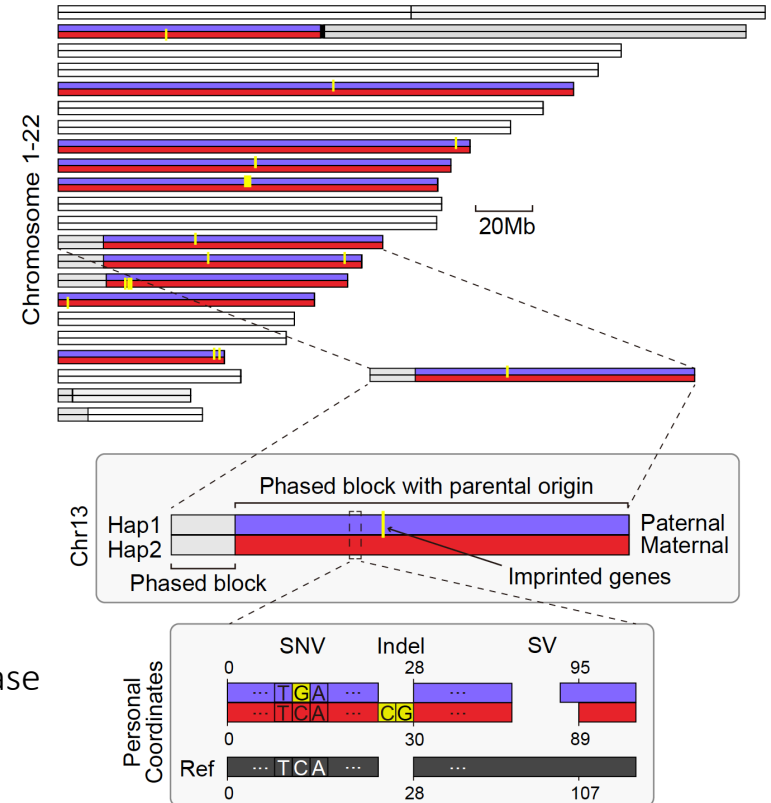


Construction of personal diploid genomes

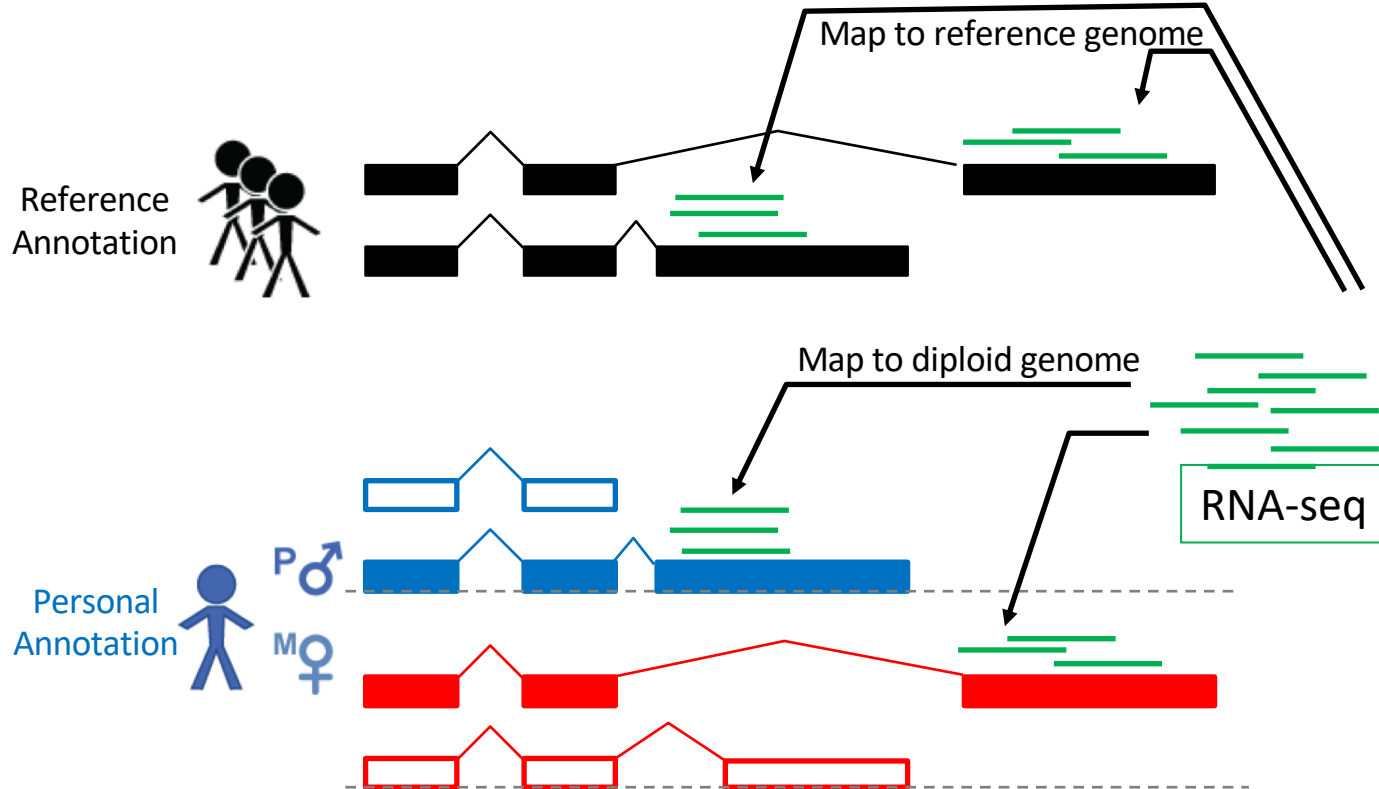


- long-read (PacBio, 10x) and short-read (Illumina) DNA-seq
- >20K long structural variants (SVs) per person
- We integrate WGBS data to call imprinted regions and phase maternal and paternal haplotypes

Diploid genome of individual 3



How to map functional genomic data

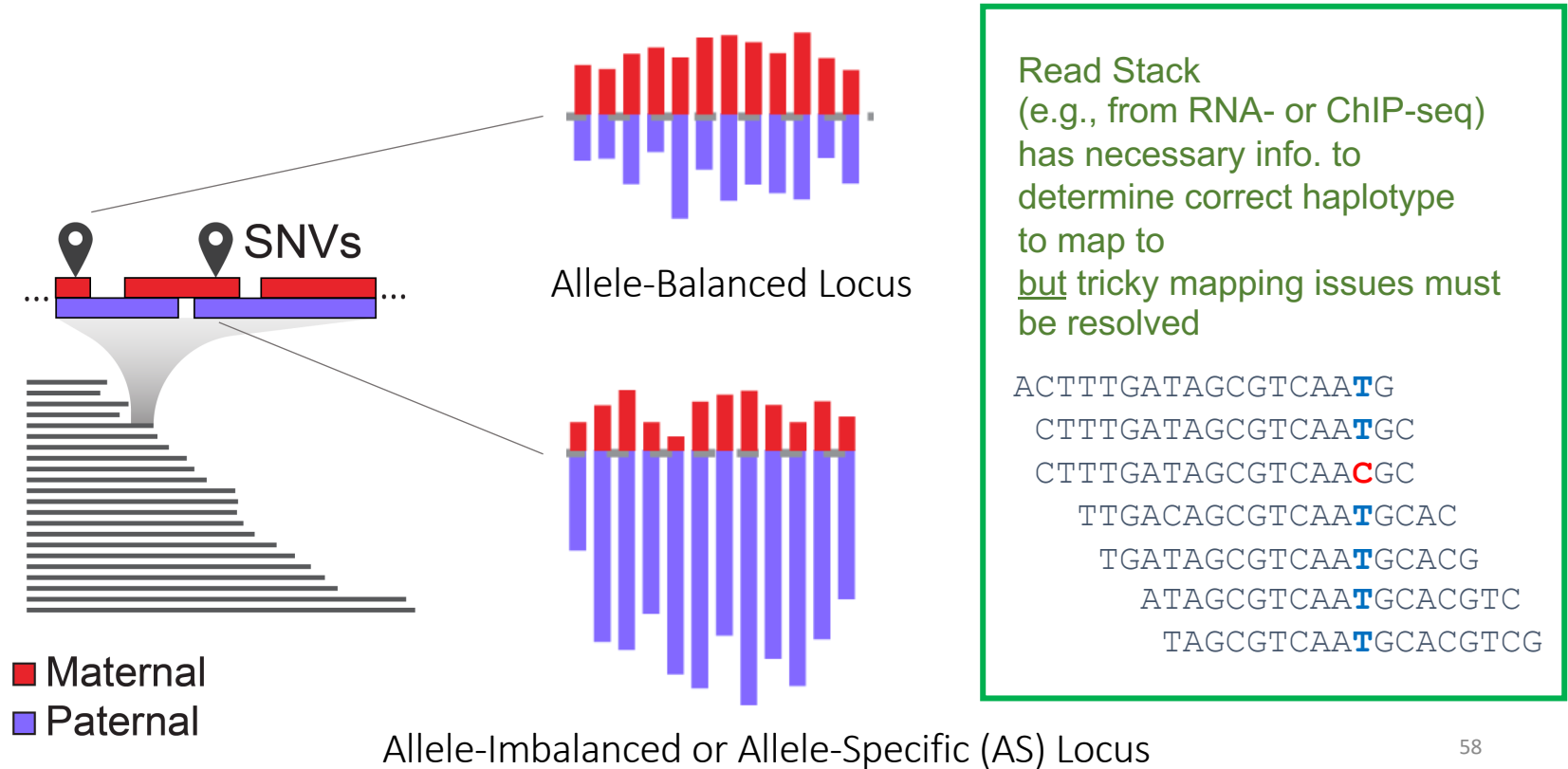


Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

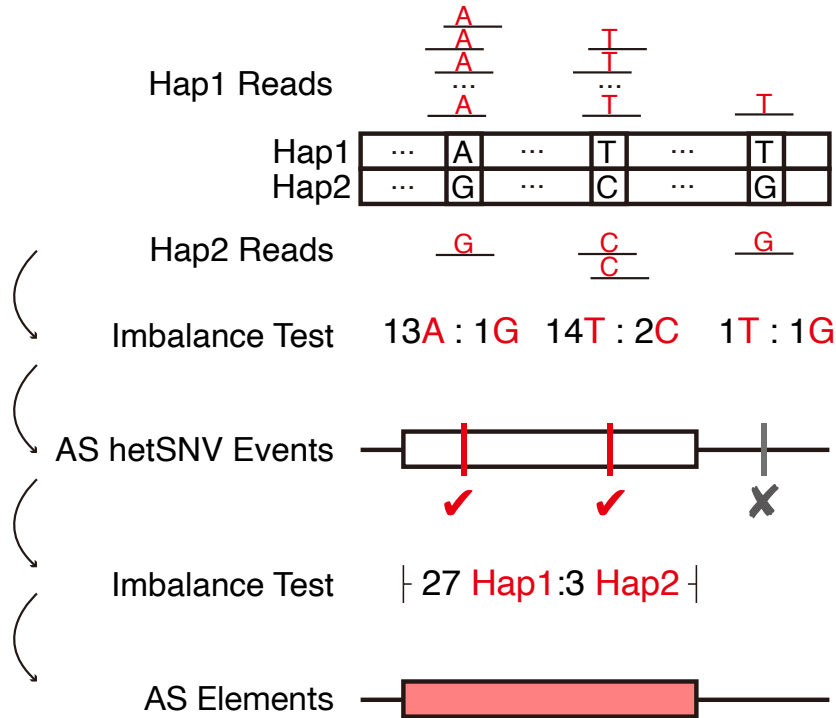
Personal genomes and the detection of allele-specific events

Personal genomes allow to study a broader set of SNVs, including individual-specific SNVs



Personal genomes and the detection of allele-specific events

AS event calling



AlleleSeq pipeline

Identify accessible SNVs

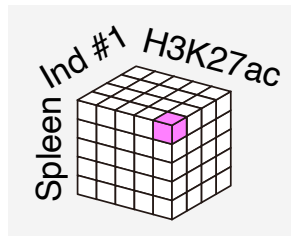
Perform beta binomial test

Identify heterozygous AS SNVs

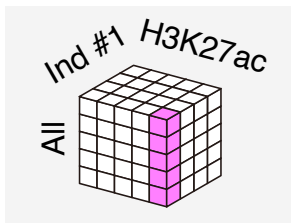
Merge reads within an element (genes, cCREs)
Perform beta binomial test

Identify AS elements (genes, cCREs)

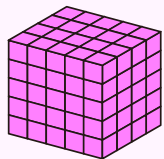
Aggregating AS Events into Catalog & their GWAS/eQTL Enrichment



EX:
H3K27ac, Ind. #1, Spleen
has ~2,600 AS SNVs



AS SNVs Across Tissues:
Union: 5,500 SNVs
Pooling: 27,000 SNVs
(Allows for Joint Calling)

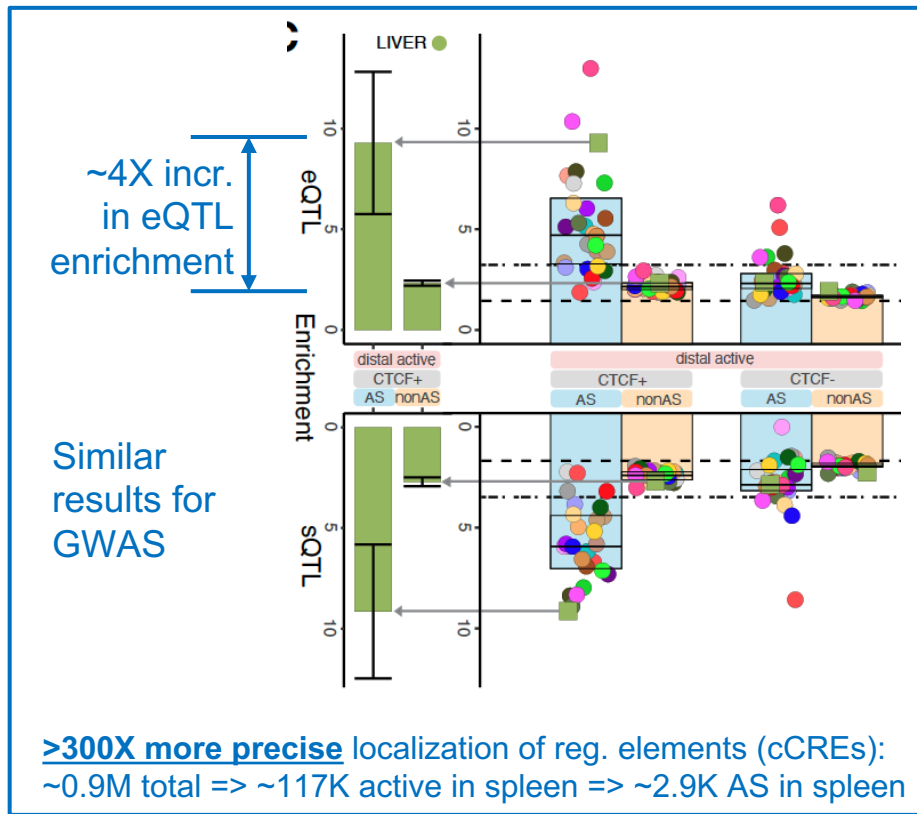


AS catalog

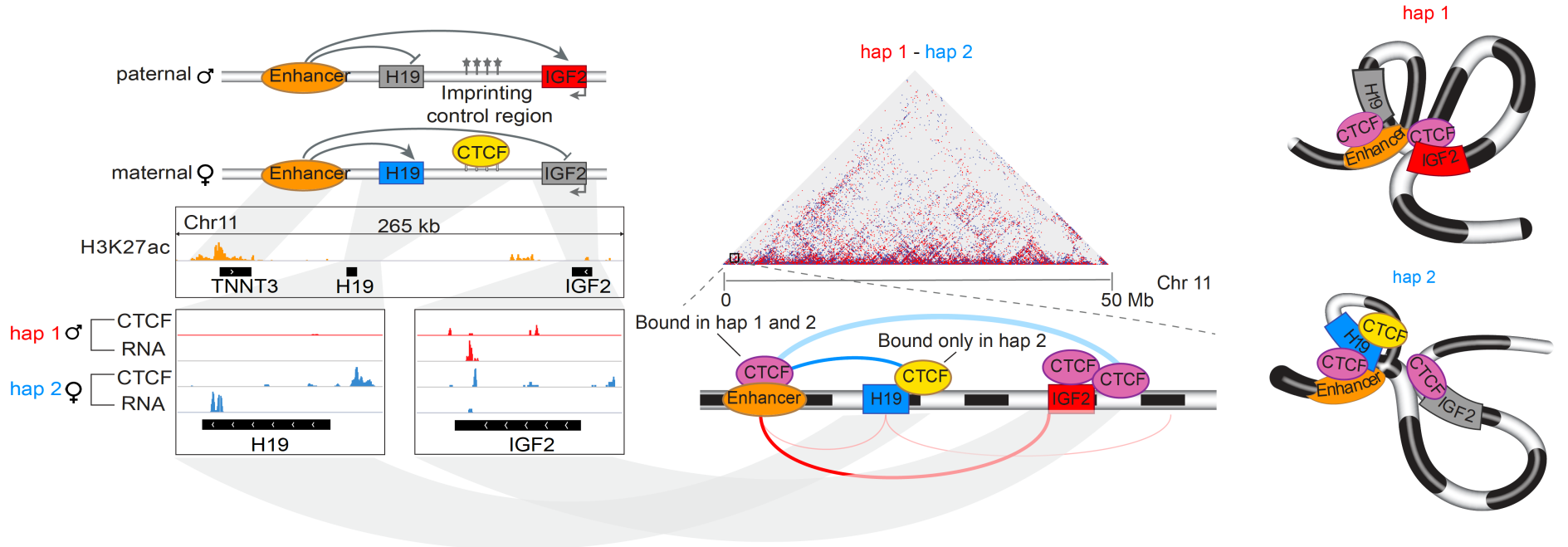


4 Individuals
12 Assays
31 Tissues

~1.3 M
SNVs
total



The classic example of *H19* and *IGF2* allele-specific activity



Recapitulating a Classic Story

AS Hi-C analysis

==>

Different M & P Chromatin

Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Relating AS SNVs and TF binding motifs



▲ AS SNV ▲ non-AS SNV
(for a given assay)

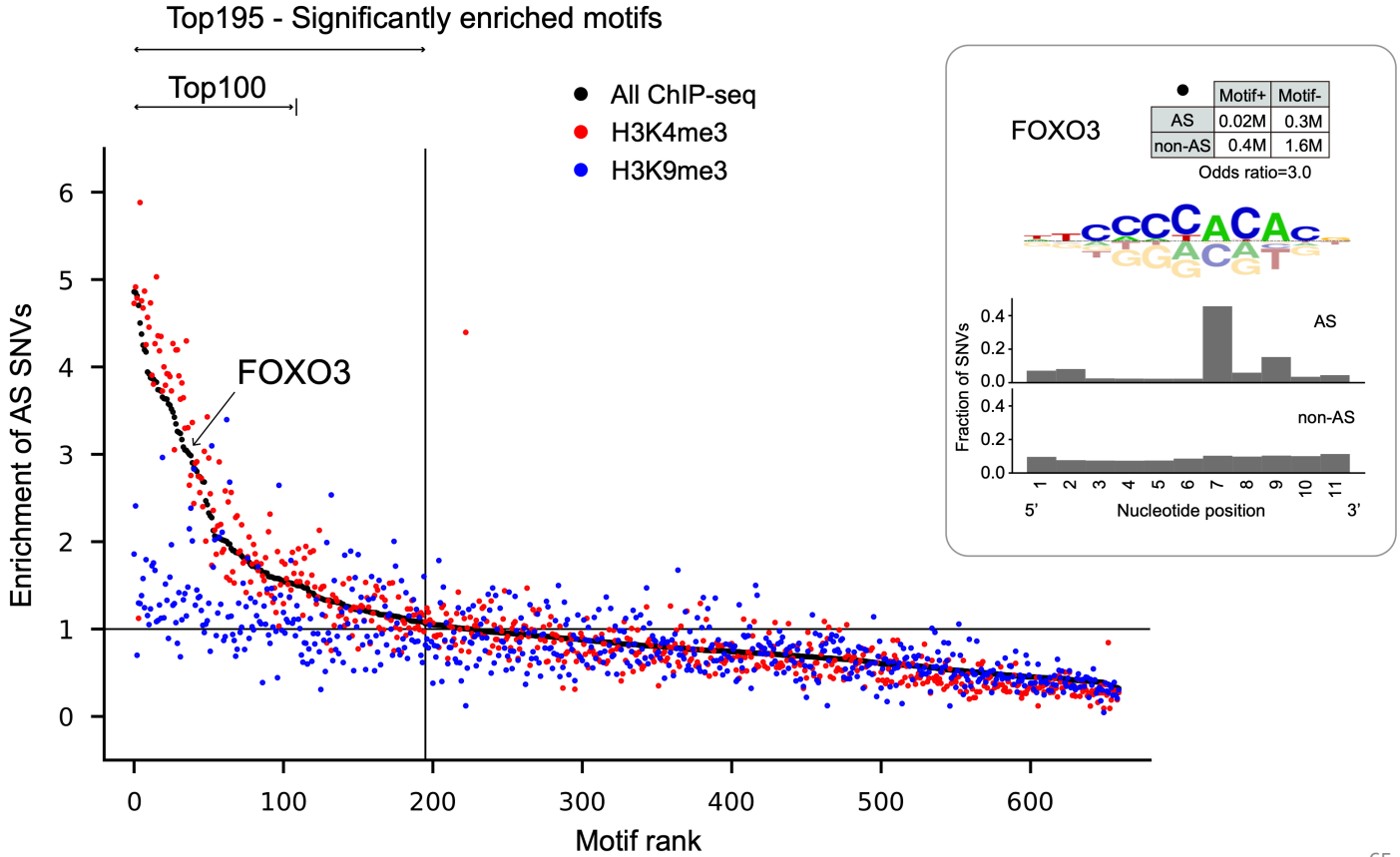
■ TF binding

There are many SNVs in the genome

- Some can impact activity (AS SNVs)
- Some don't have any impact (non-AS SNVs)

We can identify TFs with motifs more sensitive to mutations, showing enrichment in AS events

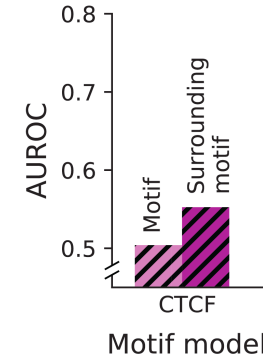
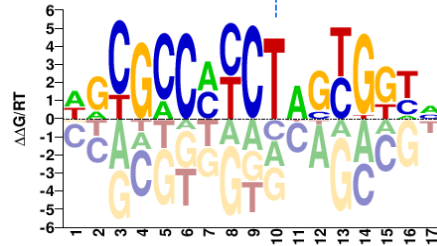
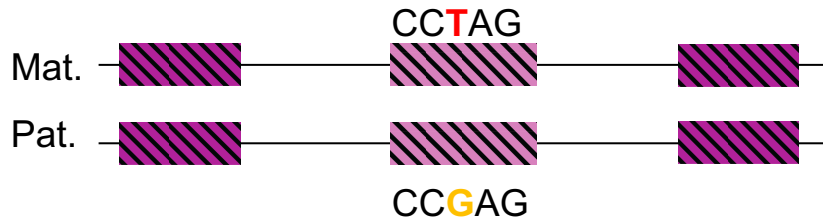
Cross-referencing AS SNVs and TF binding motifs identified AS “sensitive” TFs



Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Predicting AS activity just from nucleotide sequence (Simple)

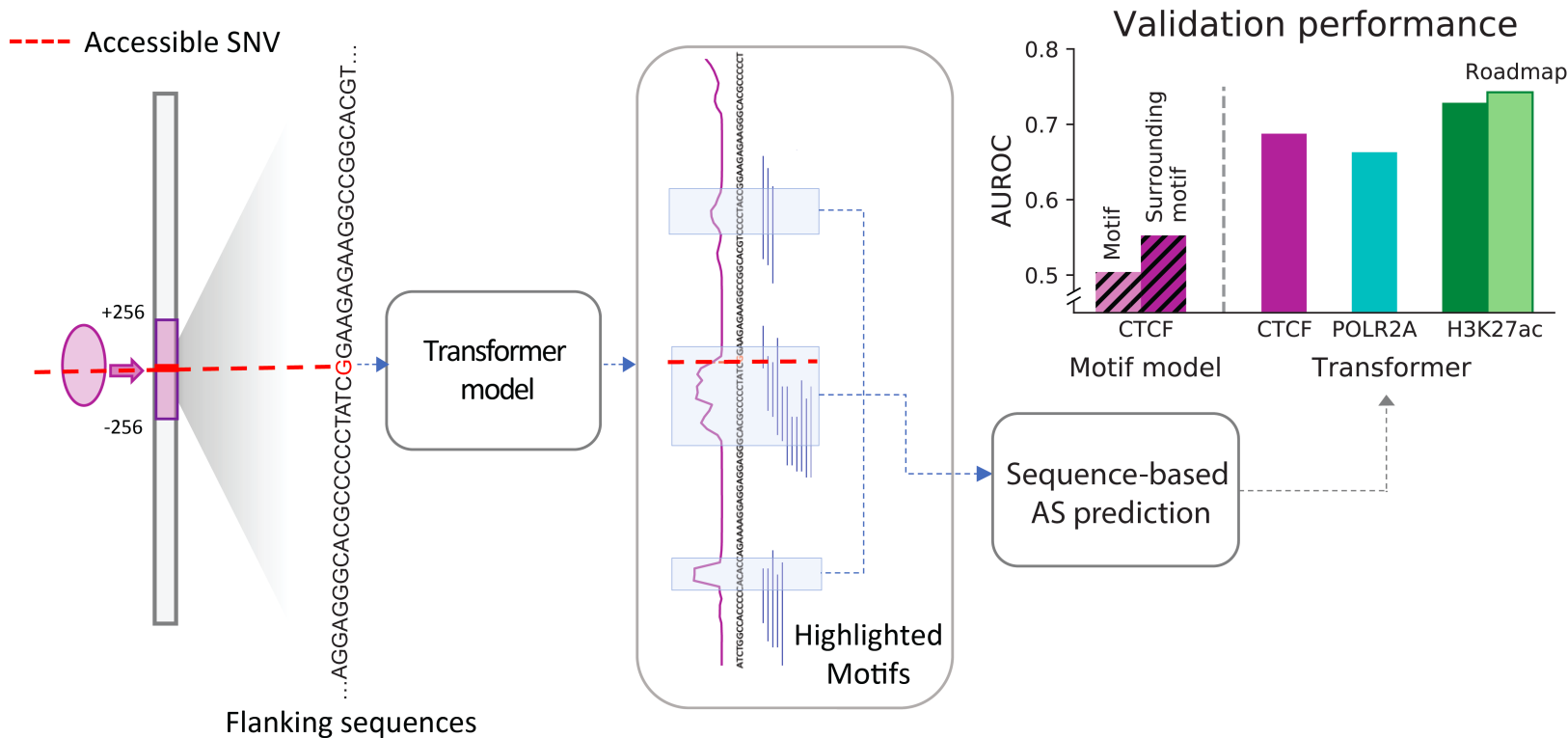


Simple logistic regression models based on overlapping & nearby motifs

- Each TF has a specific sequence that defines its motif. Can we predict AS events from those overlapping TF sequence motifs?
- EX: Predicting CTCF AS activity from its motif
- Intuition: Might think AS variant “knocking-out” a TF motif would give rise to differential AS binding

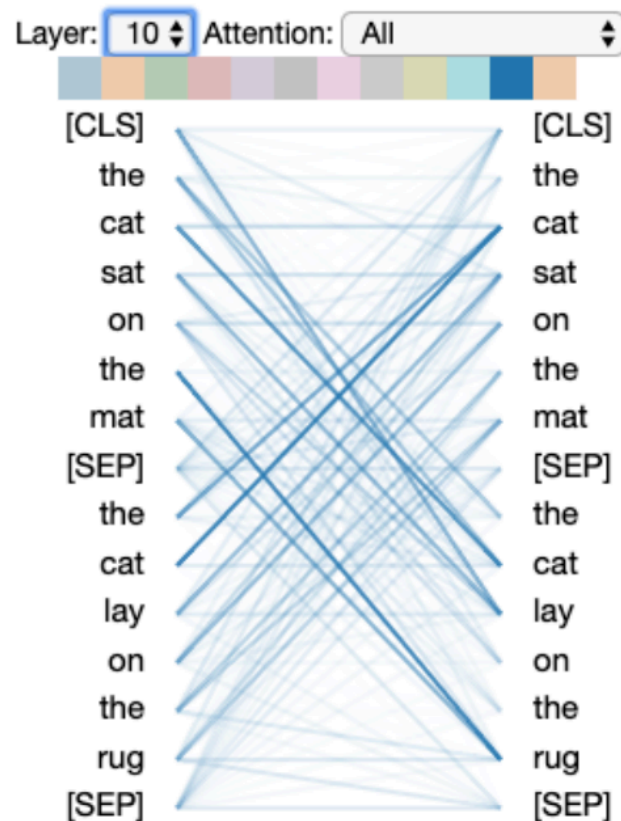
Predicting AS activity just from nucleotide sequence (Deep Learning)

Model taking into account a larger 250-bp window around the SNV yields better predictions



Contextual language model

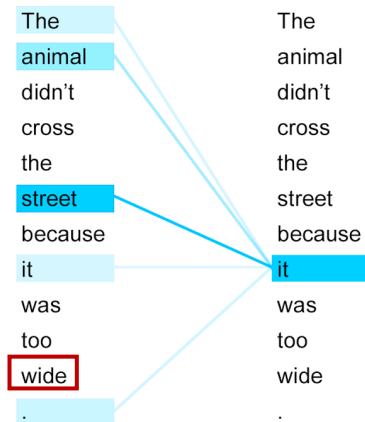
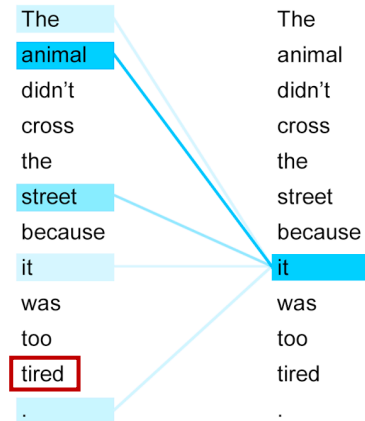
- Attention
 - How “relevant” each part of the sentence is to the rest
 - $L \times L$ attention matrix for **every word** against **every word**
 - Calculated using the whole sentence



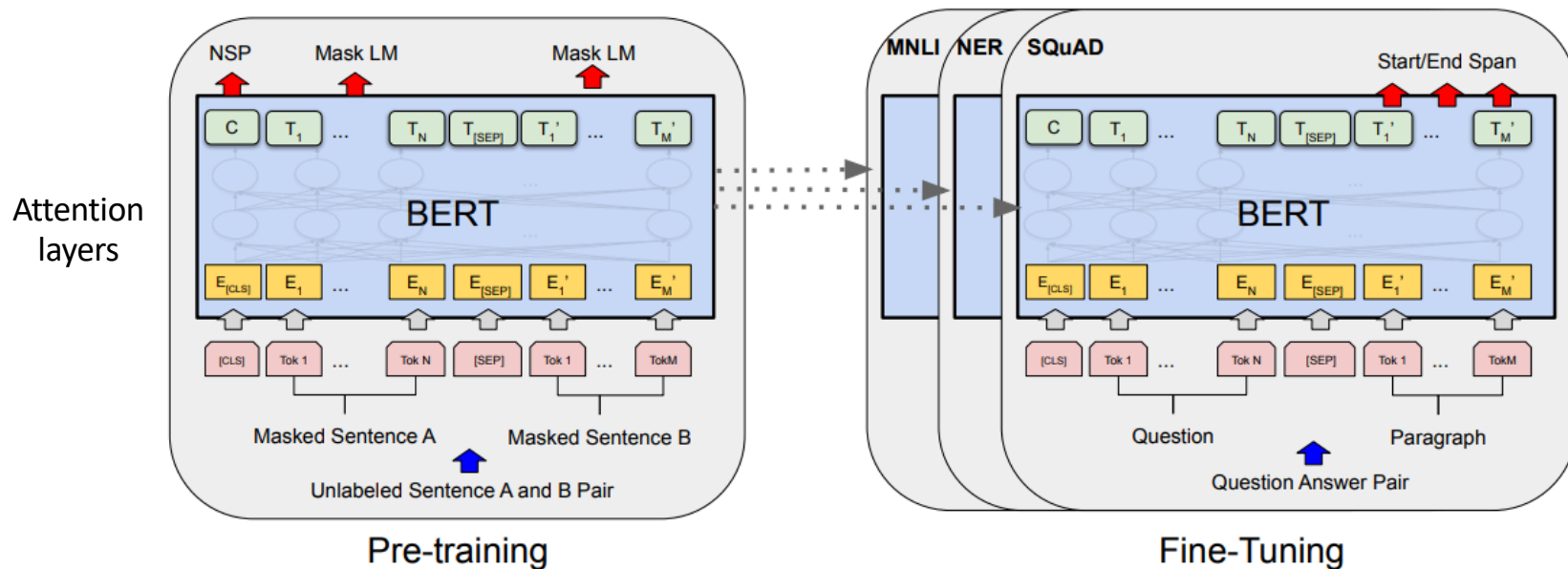
Contextual language model

*The animal didn't cross the street because **it** was too tired.*
*L'animal n'a pas traversé la rue parce qu'**il** était trop fatigué.*

*The animal didn't cross the street because **it** was too wide.*
*L'animal n'a pas traversé la rue parce qu'**elle** était trop large.*



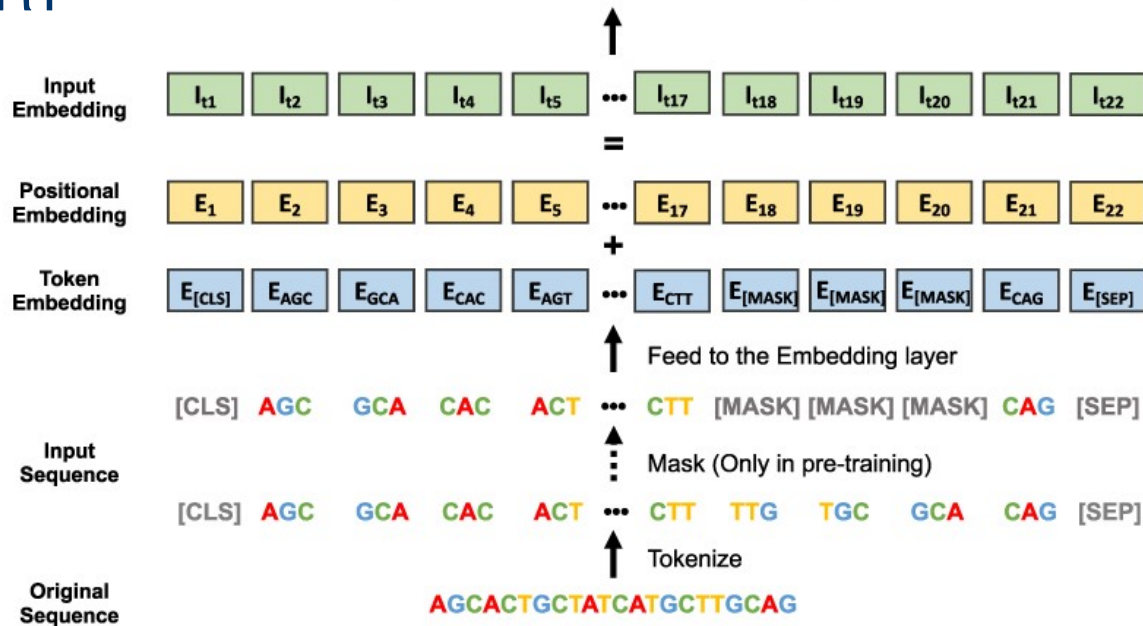
Transformer encoder & BERT



Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).

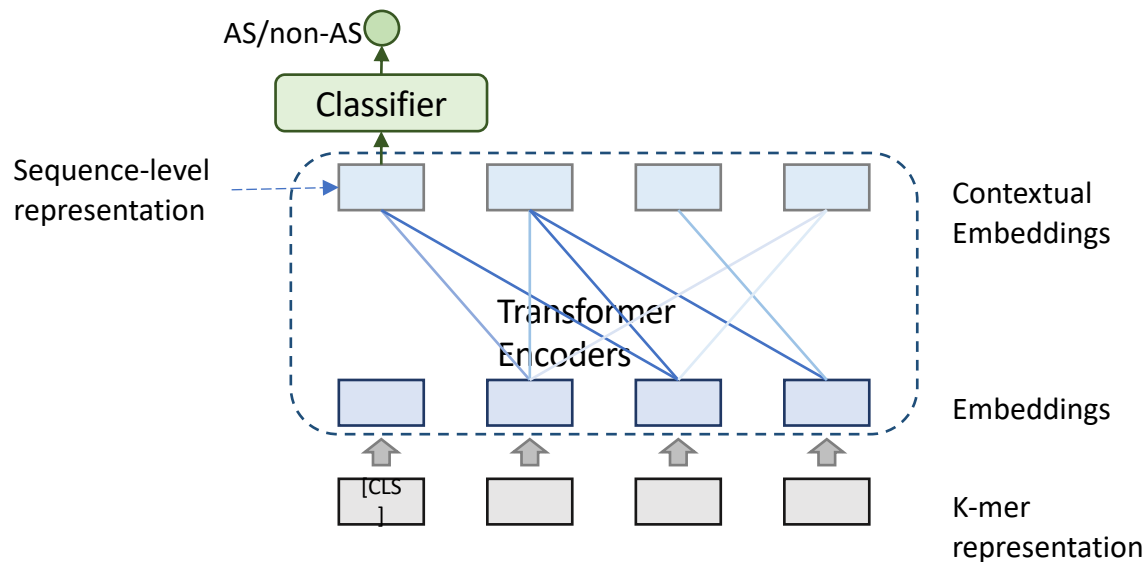
Transformer encoder & BERT

- DNABERT

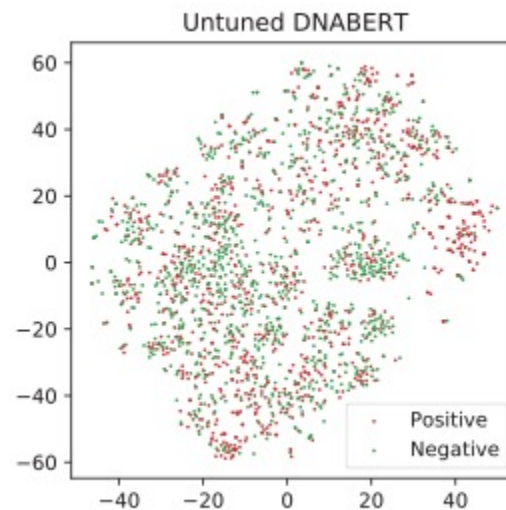
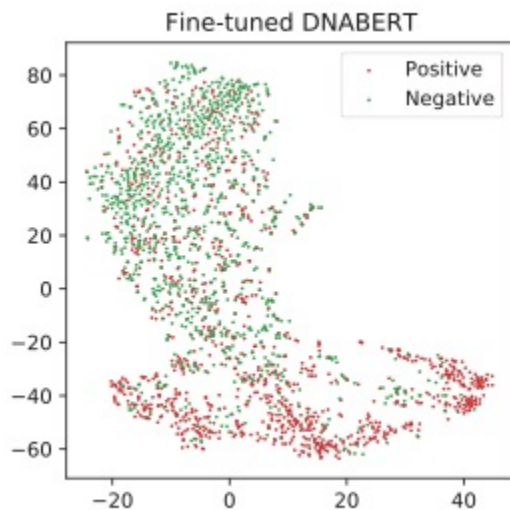


Allele specificity prediction

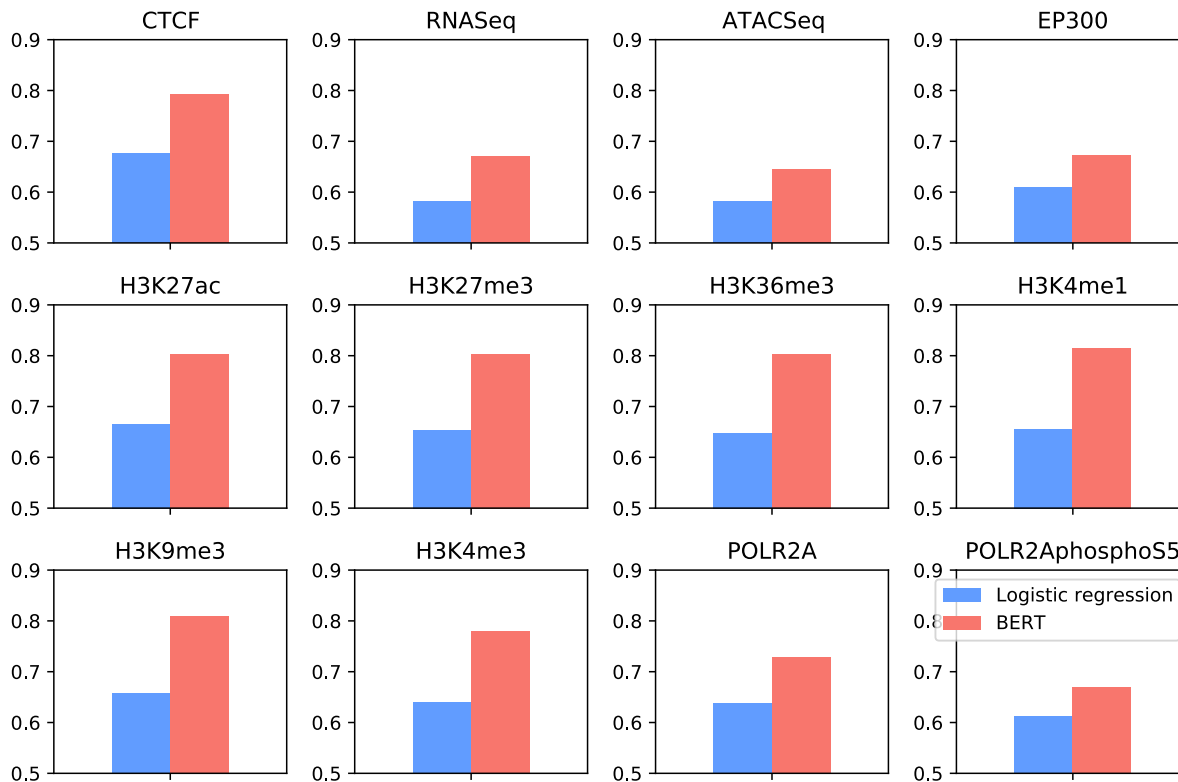
- Trained on the aggregated AS set of donor 3
 - ± 128 bp sequence context
 - Positive vs negative: AS vs non-AS heterozygous variants



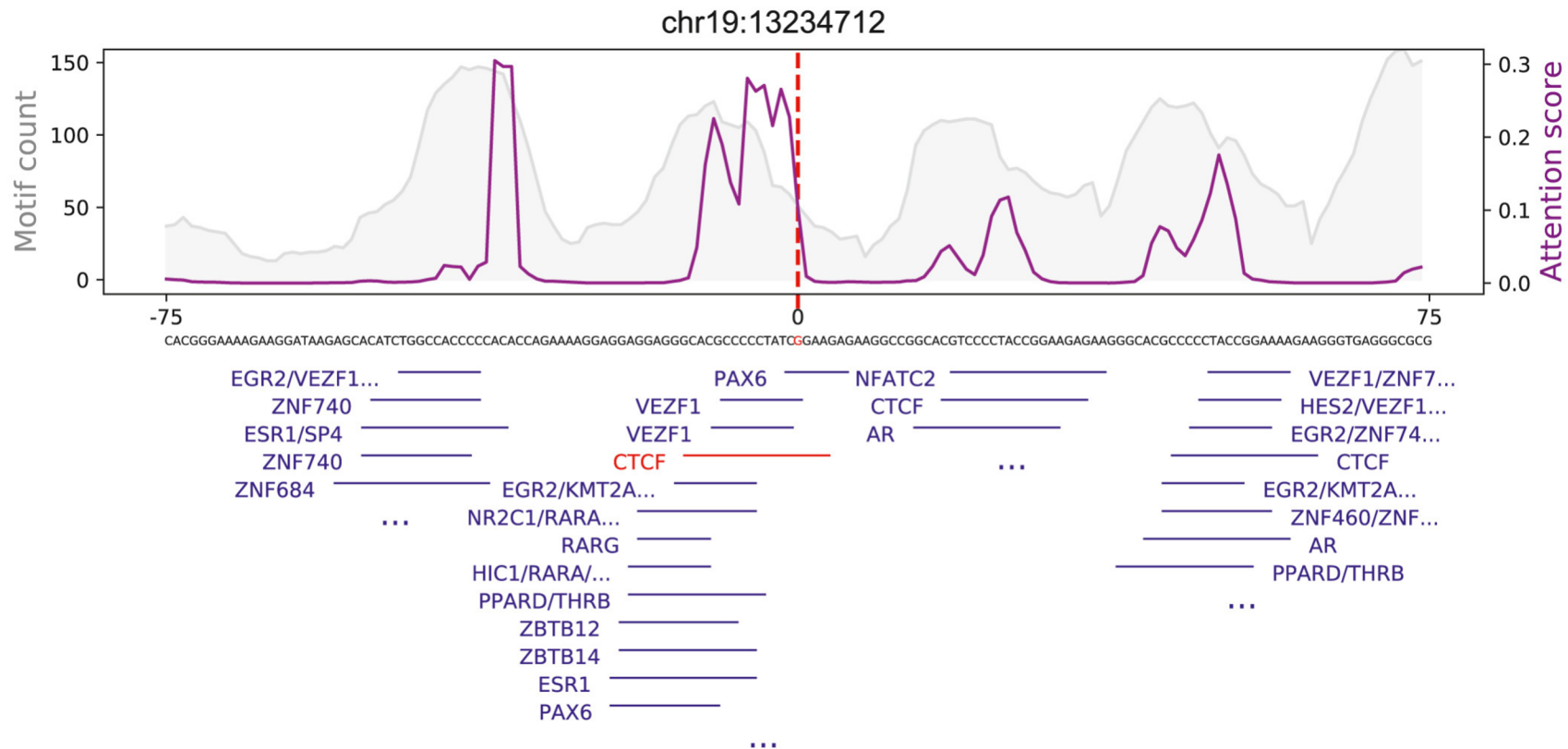
Sequence embedding



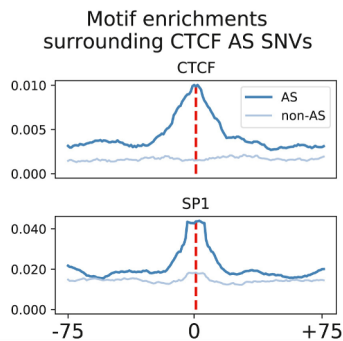
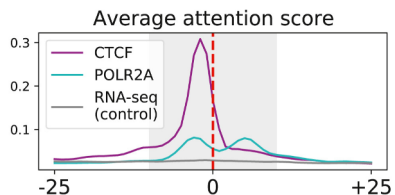
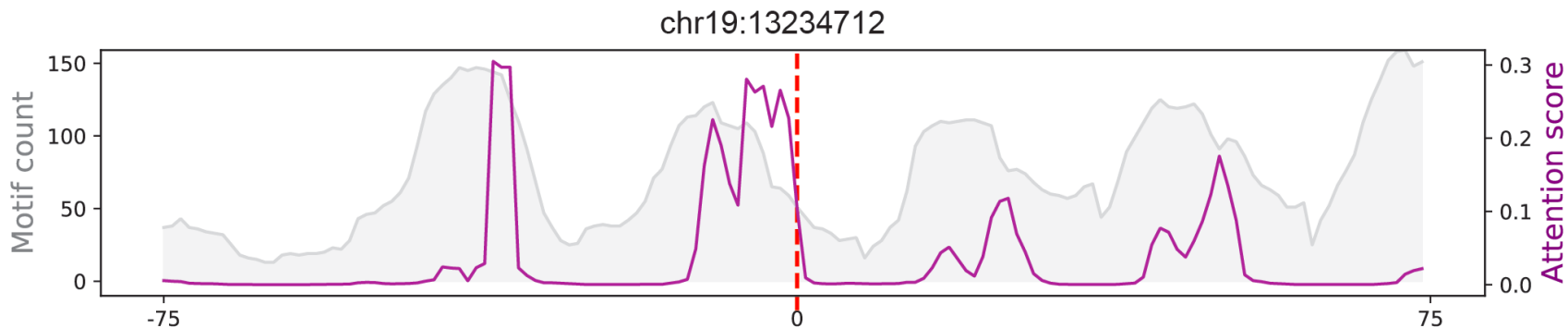
Allele specificity prediction



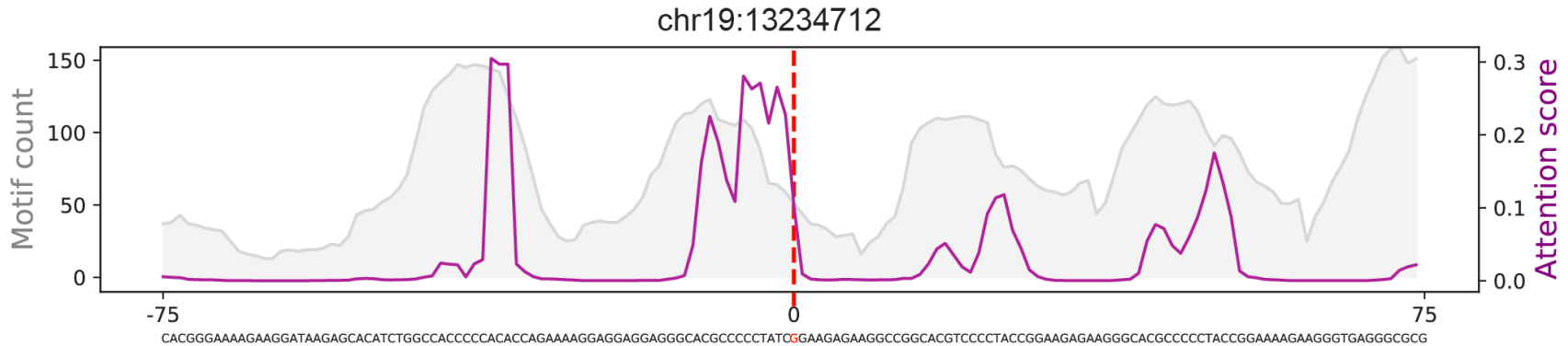
Predicting AS activity just from nucleotide sequence (Attention)



Predicting AS activity from just nucleotide sequence (Attention)

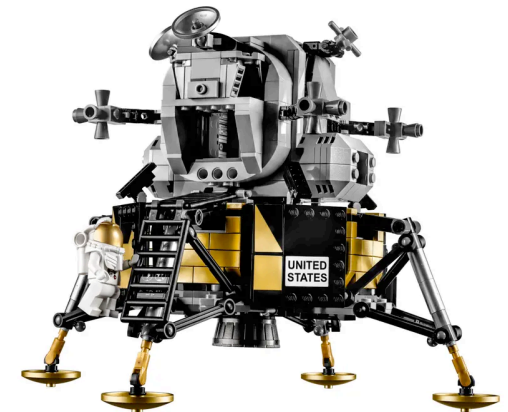


Predicting AS activity from just nucleotide sequence (Attention)



Both the motif centered at the SNV position & the surrounding sequence motifs of other TFs are relevant for AS behavior

- For instance, a mutated TF-binding site could be stabilized by other cofactors & show no AS behavior
- Similar to the legs of the Lunar Module:
If one doesn't work, the three other legs can still anchor properly



Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEx: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Developing Computational Representations of Disease-Relevant Molecules: 3 Cases Studies for AI in Biomedicine

- **Part 1: Multi-scale Modelling for Brain Disorders (Representing Molecular & Cellular Networks in a DL Framework)**
 - PsychENCODE consortium, leveraging high heritability of psychiatric diseases with functional genomics
- **Uniform Processing of Single Cell Data for 388 Brains**
 - 28 cell types, merging BICCN with a PFC-focused study
 - >500K scCREs
- **Creating Cellular Networks**
 - 1.4 M scQTLs from GTEx methods
 - Building cell-type-specific regulatory networks from scQTLs, scCREs & single-cell co-expression
 - Cell-to-cell communication networks, with changes in disease
- **Integrative Models Using These Networks**
 - Embedding regulatory networks & cell-to-cell communication networks in a deep-learning model to predict disease from genotype
 - Using this to prioritize specific pathways & genes.
 - Modelling perturbations & using these to suggest pot. drug targets
- **Part 2: Measuring Genomic Privacy Risk from a Few, Noisy SNPs (Understanding Information Leakage in terms of Constraints on Haplotype Trajectories)**
 - The Dilemma of Genomic Privacy: The genome as fundamental, inherited info that's private v. need for large-scale sharing & mining for med. research
 - 30 SNPs from "environmental" coffee cup sample sufficient for ID
 - Based on finding most likely haplotype "trajectories" in a genome DB
 - Single trajectory for a unique match for an ensemble of equivalent ones for near match
 - Calculating a PRS score over an ensemble
- **Part 3: Variant Impact for Precision Medicine (Learning Distributed Sequence Patterns via Transformer Attention)**
 - Need for variant catalogs & interpretation resources
- **EN-TEX: a Resource for Variant Interpretation**
 - >1500 functional experiments with diploid genomes of 4 individuals
 - Differential mapping to haplotypes
- **Development of AS Catalog**
 - >1M allele-specific events, over all samples from jt. calling
 - Useful biological interpretation: chrX, SVs, Igf2-H19
 - Association of AS events & eQTL/GWAS variants; provides a source of variant interpretation
 - Relating AS events to tissue-specificity & conservation
- **Transformer model to Predict AS Variants**
 - Identification of sensitive TF binding motifs
 - Model can successfully predict if a SNV will be AS purely from sequence; however, it requires extended context (~150 bp) around SNV
 - Interpretation in terms of anchoring co-factors

Acknowledgements

github.com/gersteinlab/**PLIGHT**

P **Emani**, Geradi, M.N.; Gürsoy, G.; Grasty, M.R.; Miranker, A

Also: **JOBS**.gersteinlab.org

Team work with >80 authors

13 Labs

P. S. Emani, J. J. Liu, D. Clarke, M. Jensen, J. Warrell,
C. Gupta, R. Meng, C. Y. Lee, S. Xu, C. Dursun, S. Lou

Y. Chen, Z. Chu, T. Galeev, A. Hwang, Y.
Li, P. Ni, X. Zhou

T. Chatterjee, Y. Dai, Zi. Duan, M.
Gancz, D. Garrido-Martín, E. Henry,
G. E. Hoffman, A. Huang, Y. Jiang,
T. Jin, S. Khullar, J. Liu, S. Liu, J.
Moore, E. Nguyen, N. Phalke, M.
Pjanic, H. Pratt, A.S. Rajagopalan, T.
R. Riesenmy, N. Shedd, G. Wang, Y.
Xia, A. C. Yang, S. Zheng, D. Lee, Z.
Weng, H. Won

T. E. Bakken, J. Bendl, L. Bicks, L. Cheng, Y.
Cheng, M. Flaherty, J. F. Fullard, S. Gaynor-
Gillett, J. Grundman, N. Hawken, N. L.
Jorstad, R. Kawaguchi, J. Liu, S. Ma, M.
Margolis, S. Mazariegos, J. R. Moran, D.
Quintero, M. Shi, M. Spector, R. Terwilliger,
K. J. Travaglini, B. Wamsley, S. Xiao, M. J.
Gandal, E. S. Lein, P. Roussos, N. Sestan, K.
P. White

M. J. Girgenti,
J. Zhang, D. Wang, D.
Geschwind, M. Gerstein

& PsychENCODE Consortium



212 members
Labs from 37 universities/institutions



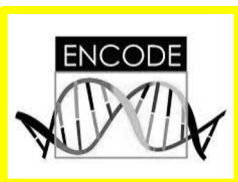
Team Science with >100 Authors, 20 Labs, 3 Consortia, 1 Funder

Joel Rozowsky[§], Jiahao Gao[§], Beatrice Borsari[§], Yucheng T Yang[§], Timur Galeev[§], Gamze Gursoy[§], Charles B Epstein[§], Kun Xiong[§], Jinrui Xu[§], Tianxiao Li[§], Jason Liu[§],

Keyang Yu^{7†}, Ana Berthel^{2,3†}, Zhanlin Chen^{8†}, Fabio Navarro^{2,3†}, Maxwell S Sun^{2,3†}, James Wright^{9†}, Justin Chang^{2,3†}, Christopher JF Cameron^{2,3†}, Noam Shores^{6†}, Elizabeth Gaskell^{6†}, Jorg Drenkow^{10†},

Jessika Adrian¹¹, Sergey Aganezov¹², François Aguet⁶, Gabriela Balderrama-Gutierrez¹³, Samridhi Banskota⁶, Guillermo Barreto Corona⁶, Sora Chee¹⁴, Surya B Chhetri¹⁵, Gabriel Conte Cortez Martins^{2,3}, Cassidy Danyko¹⁰, Carrie A Davis¹⁰, Daniel Farid^{2,3}, Nina P Farrell⁶, Idan Gabdank¹¹, Yoel Gofin⁷, David U Gorkin¹⁴, Mengting Gu^{2,3}, Vivian Hecht⁶, Benjamin C Hitz¹¹, Robbyn Issner⁶, Yunzhe Jiang^{2,3}, Melanie Kirsche¹², Xiangmeng Kong^{2,3}, Bonita R Lam¹¹, Shantao Li^{2,3}, Bian Li^{2,3}, Xiqi Li⁷, Khine Zin Lin¹¹, Ruibang Luo¹⁶, Mark Mackiewicz¹⁵, Ran Meng^{2,3}, Jill E Moore¹⁷, Jonathan Mudge¹⁸, Nicholas Nelson⁶, Chad Nusbaum⁶, Ioann Popov^{2,3}, Henry E Pratt¹⁷, Yunjiang Qiu¹⁴, Srividya Ramakrishnan¹², Joe Raymond⁶, Leonidas Salichos^{2,3,19}, Alexandra Scavelli¹⁰, Jacob M Schreiber²⁰, Fritz J Sedlazeck^{12,21,22}, Lei Hoon See¹⁰, Rachel M Sherman¹², Xu Shi^{2,3}, Minyi Shi¹¹, Cricket Alicia Sloan¹¹, J Seth Strattan¹¹, Zhen Tan^{2,3}, Forrest Y Tanaka¹¹, Anna Vlasova^{4,23,24}, Jun Wang^{2,3}, Jonathan Werner¹⁰, Brian Williams²⁵, Min Xu^{1,2}, Chengfei Yan^{2,3}, Lu Yu⁹, Christopher Zaleski¹⁰, Jing Zhang²⁶,

Kristin Ardlie, Michael Cherry, Eric Mendenhall, William Noble, Zhiping Weng, Morgan Levine, Alexander Dobin, Barbara Wold, Ali Mortazavi, Bing Ren, Jesse Gillis, Richard Myers, Michael Snyder, Jyoti Choudhary, Aleksandar Milosavljevic, Michael Schatz[#], Bradley Bernstein[#], Roderic Guigó[#], Thomas Gingeras[#], Mark Gerstein[#]



Labs,



Labs,



Labs



Image Credit:
Susanna Liu (Yale)

entex.encodeproject.org

Extra



Info about content in this slide pack

- **General PERMISSIONS**

- This Presentation is copyright Mark Gerstein, Yale University, 2019.
 - Please read permissions statement at www.gersteinlab.org/misc/permissions.html .
 - Feel free to use slides & images in the talk with PROPER acknowledgement (via citation to relevant papers or link to gersteinlab.org).
 - Paper references in the talk were mostly from Papers.GersteinLab.org.
-
- **PHOTOS & IMAGES.** For thoughts on the source and permissions of many of the photos and clipped images in this presentation see <http://streams.gerstein.info> .
 - In particular, many of the images have particular EXIF tags, such as `kwpotppt` , that can be easily queried from flickr, viz: <http://www.flickr.com/photos/mbgmbg/tags/kwpotppt>