# 2025 Translational Bioinformatics Year-in-Review

@proftatonetti @tatonetti.bsky.social

**Nicholas Tatonetti, PhD, FACMI**
Professor of Computational Biomedicine
Cedars-Sinai Medical Center

**2025 AMIA Summits - Pittsburgh**

# Disclosures

- NIH, FDA, DoD, Pfizer, AstraZeneca, Janssen, Amgen, PhARMA Foundation, CARI Health

- Co-Editor-in-Chief of BioData Mining

- I am influenced by my professional and personal network and experiences

- Biggest conflict: I am a geek for translational bioinformatics, methods that solve problems, computational medicine 😍

# Goals

- Review trends in the translational bioinformatics literature

- Create a "snapshot" of what the field is doing now (Spring 2025)

- Recognize innovative work and identify opportunities for the future

# Process

- Follow the literature throughout the year (i.e. my lab's #papers channel)

- Work with the talented and generous **AMIA Year-in-Review Committee**

- Triage all papers from a set of relevant journals since Jan 2024

  - Evaluate papers on a set of TBI criteria, score on:

    - Informatics Novelty, Application Importance, **Wow** Factor (total 0-9)

- I then take these scores and select papers to highlight in 1-5 slides

# Caveats

- Translational bioinformatics =

  **Informatics** methods that link
  **biological entities** (genes, proteins, cells, small molecules)
  to **clinical entities** (drugs, diseases, symptoms, etc.)
  — or vice versa.

- Covers the last 14 months (Jan 2024 - Mar 2025)

- Focused on human biology

- **What's NOT included:**

  - Amazing biology with straightforward informatics (PRS, looking at you 👀)

  - Amazing informatics but no link between the clinical and the molecular

  - Perspectives, reviews (for the most part)

*This is all thanks to...*

# The 2025 AMIA Year-in-Review Team!

2025 Translational Bioinformatics Year-in-Review Committee

# Final List

- 822 papers triaged; 221 reviewed and scored by the committee

- 25 presented here + 9 shout outs + 3 pieces of brain candy

  - Apologies for those I missed, misunderstood, or misjudged, biases/mistakes are all mine

- 6 TBI topics:

  - **Lose Control** - *Drug Discovery & Repurposing*

  - **Cruel Summer** - *Taylor-ed for you - Precision Medicine in Action*

  - **Good Luck, Babe** - *Bio Euphoria - Integrating Clinical & Molecular Data*

  - **What Was I Made For?** - *Biomarker Discovery & Validation*

  - **I Remember Everything** - *EHR, Real-world Evidence, & Epidemiology*

  - **Houdini** - *Emerging Therapeutics & Technologies*

- All authors are mentioned if ≤3, all first authors otherwise

- Slides will be posted to www.tatonettlab.org and linked to my Bluesky and other social media accounts
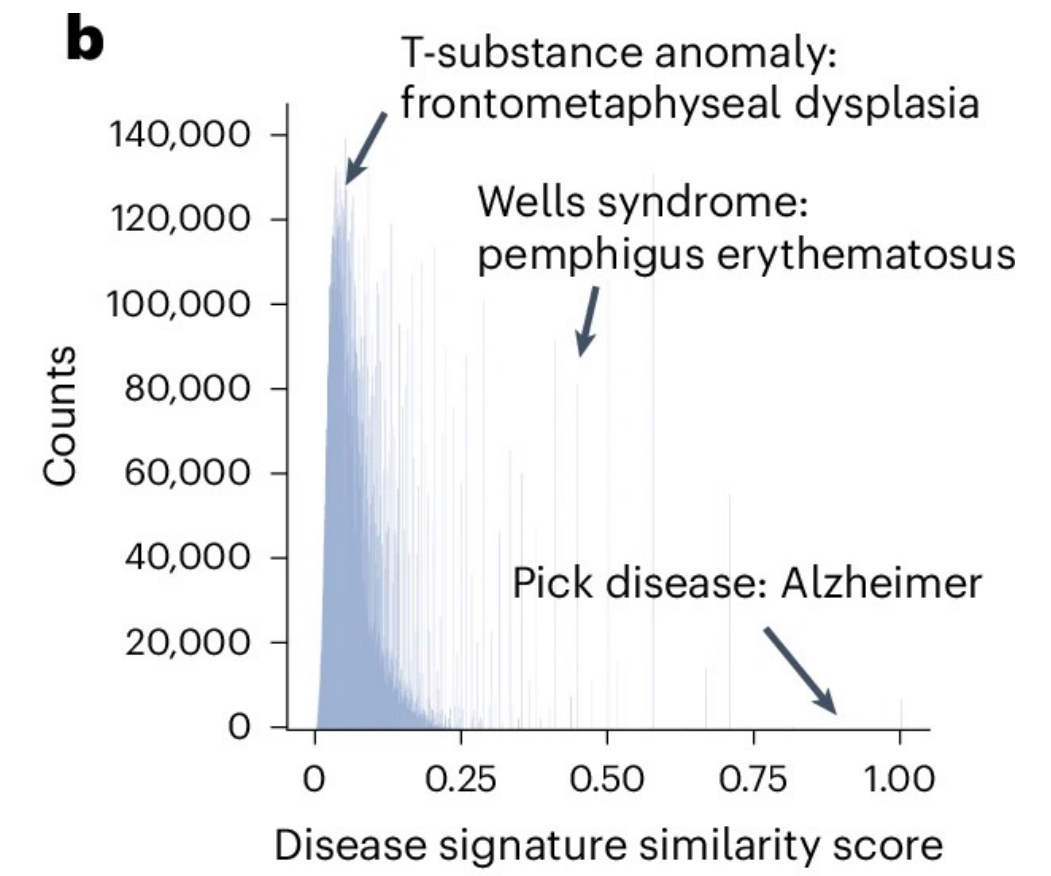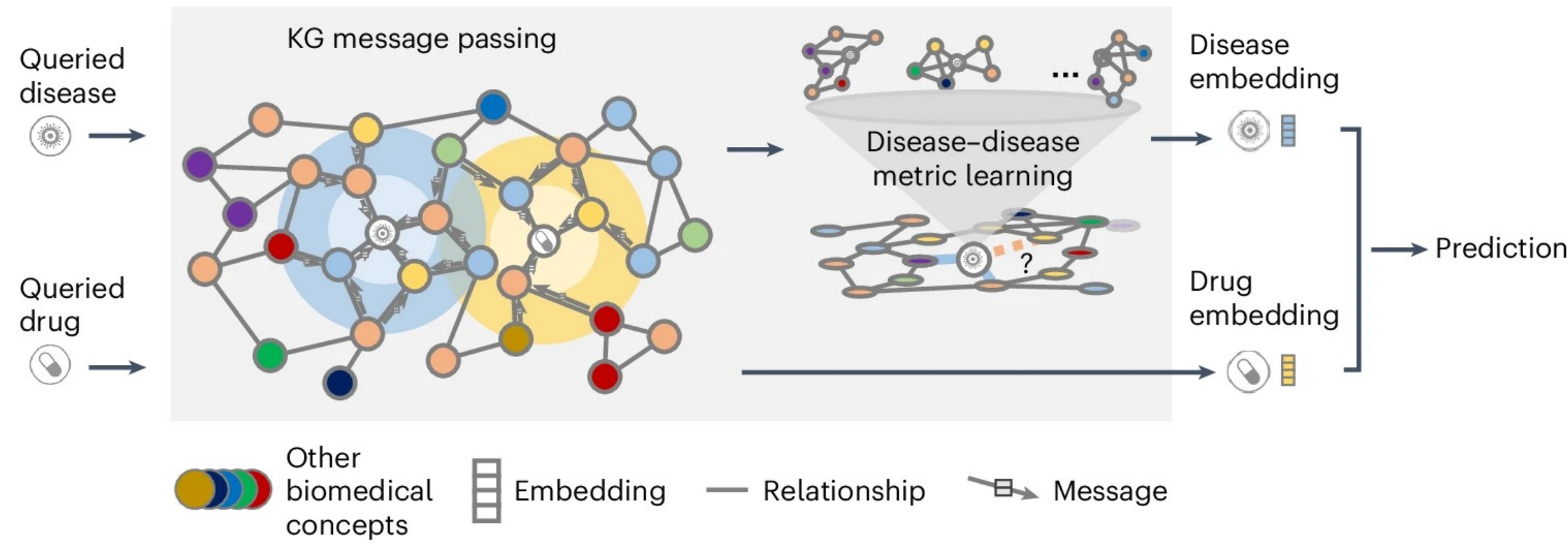
# Here we go…

# "Lose Control"

## Drug Discovery & Repurposing

# A foundation model for clinician-centered drug repurposing (Huang et al, *Nature Medicine*)
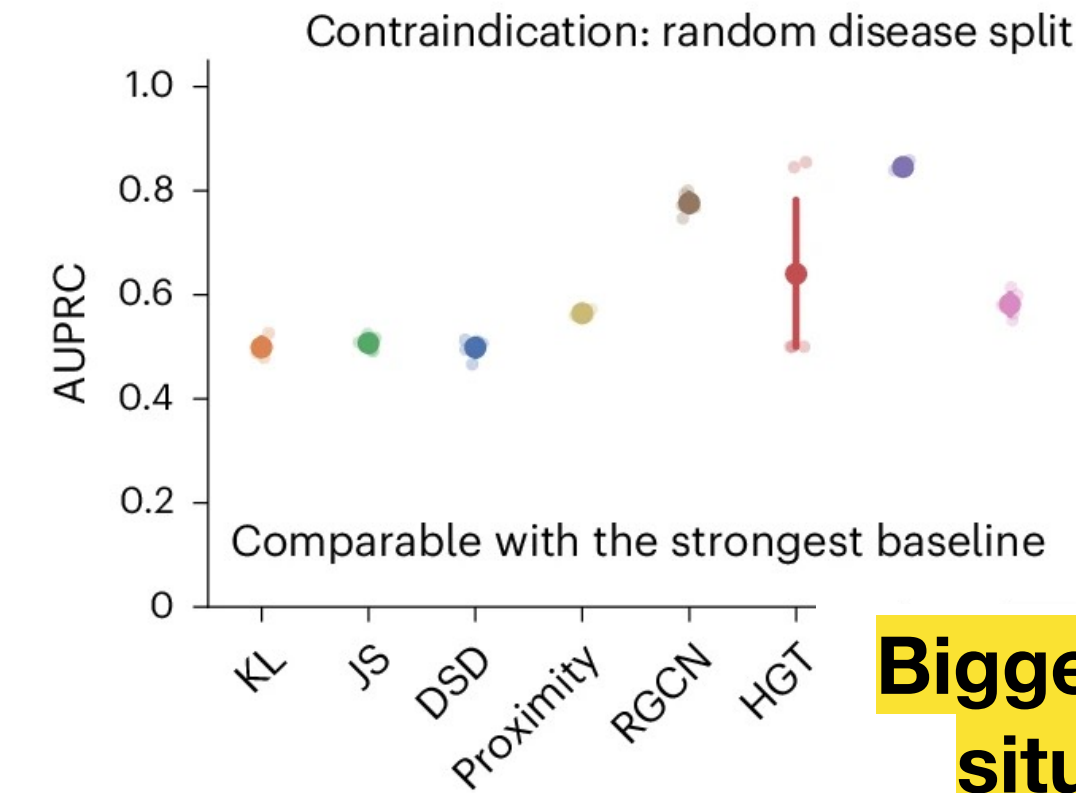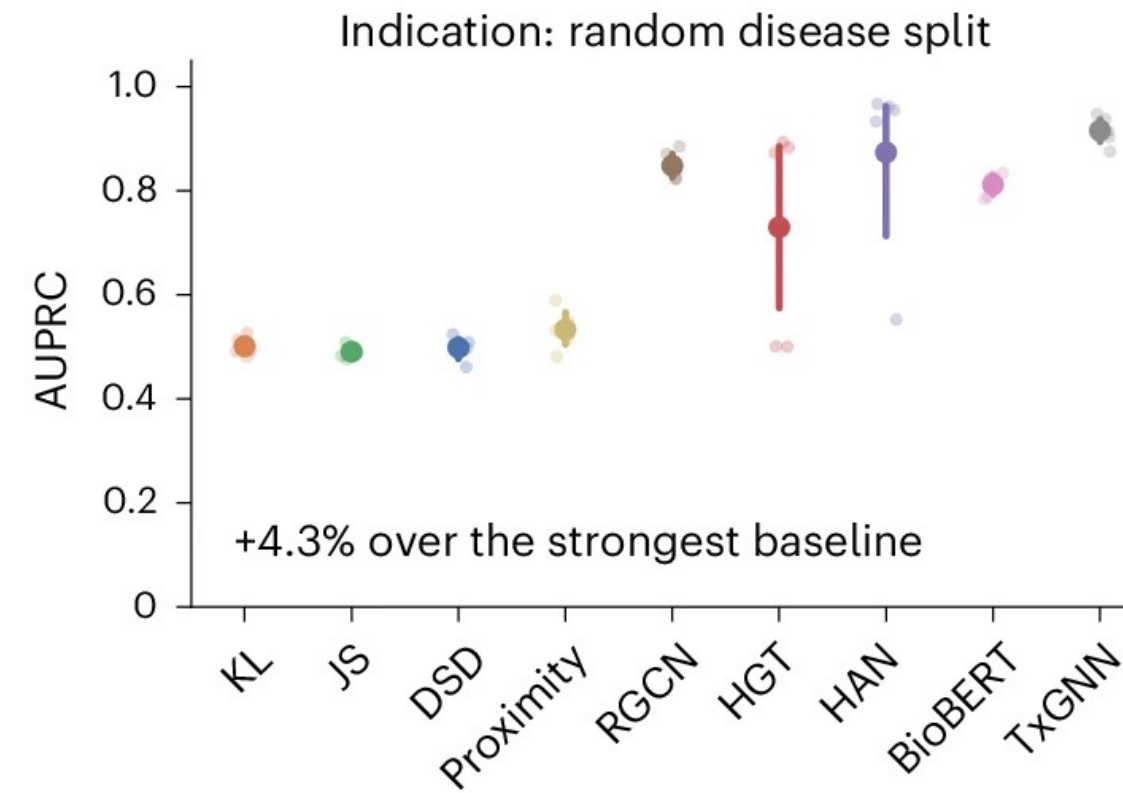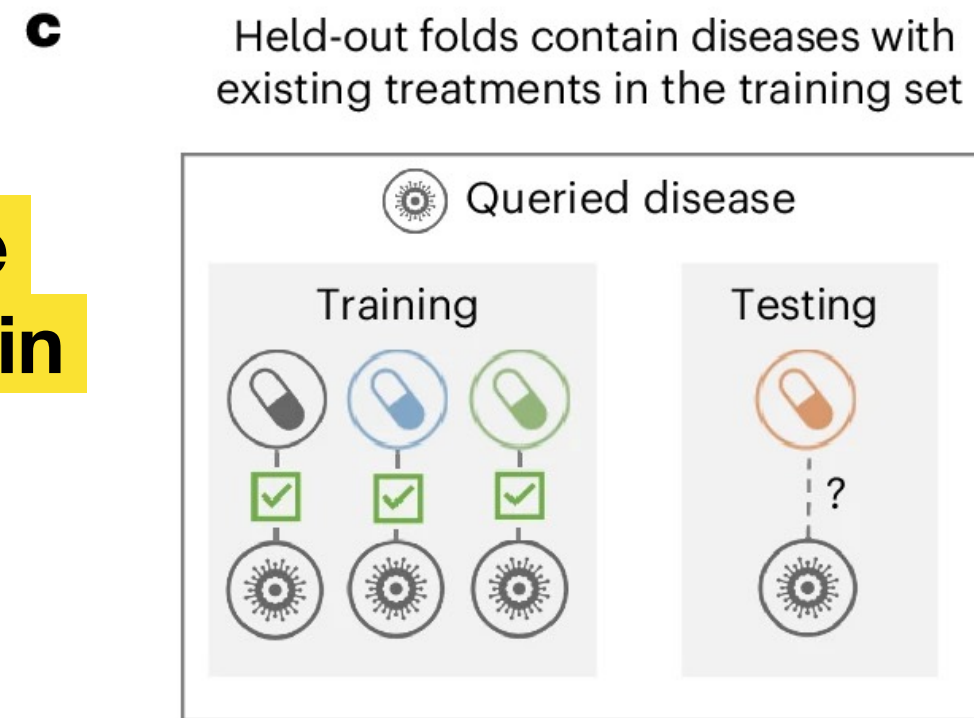
- Goal:

  - Previous drug repurposing strategies require some drugs for a disease to copy/learning from

  - Build a zero-shot algorithm that can make predictions even when no previous examples of drugs exist

- Method:

  - Embed diseases based on a knowledge graph of relationships between diseases, drugs, proteins, pathways, and clinical phenotypes — <u>allows diseases to share a embedding space and thus share information</u>

  - Use Graphical Neural Network (TxGNN) with message passing to produce embeddings for drugs or diseases (and for any of the other concepts for that matter)

  - Drugs and diseases are embedded in same space allowing them to be directly compared

- Result:

  - 49.2% improvement in drug indication prediction accuracy and 35.1% improvement in contraindication predictions over competing methods

  - Predictions aligned with real-world off-label use drug prescriptions

- Conclusion: Beautiful example of the power of building foundation models to improve performance for situations with little data available.

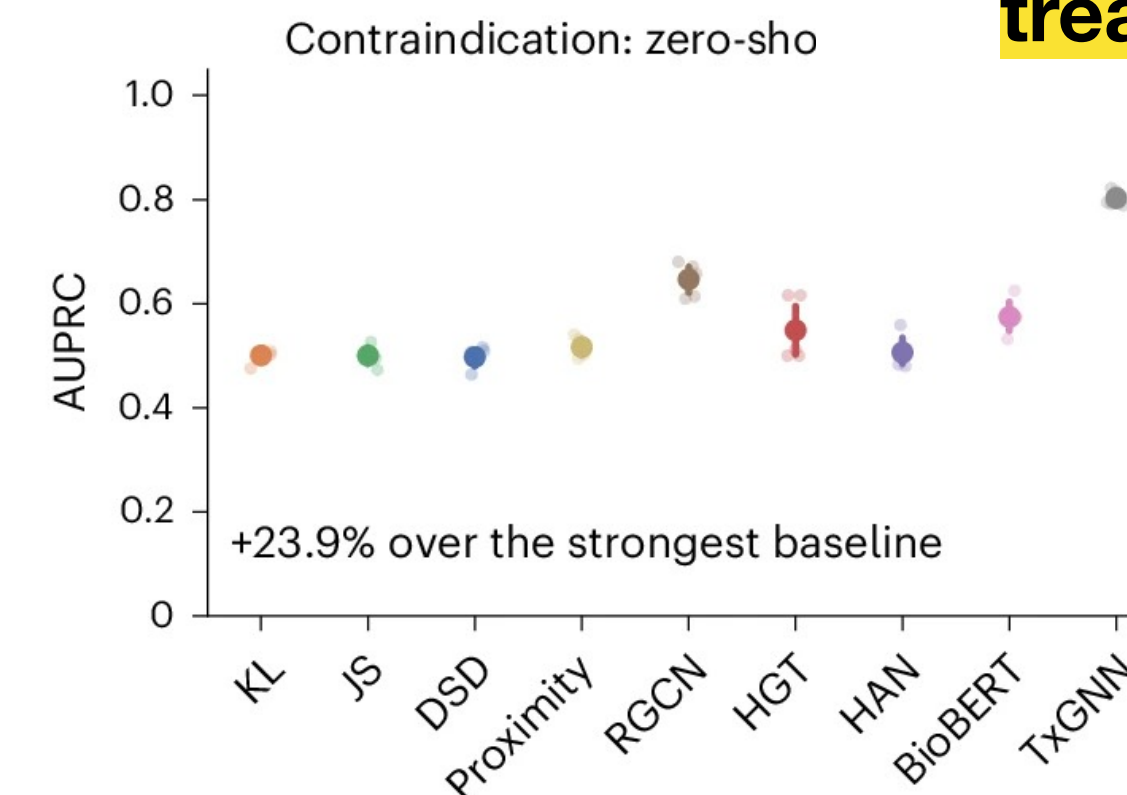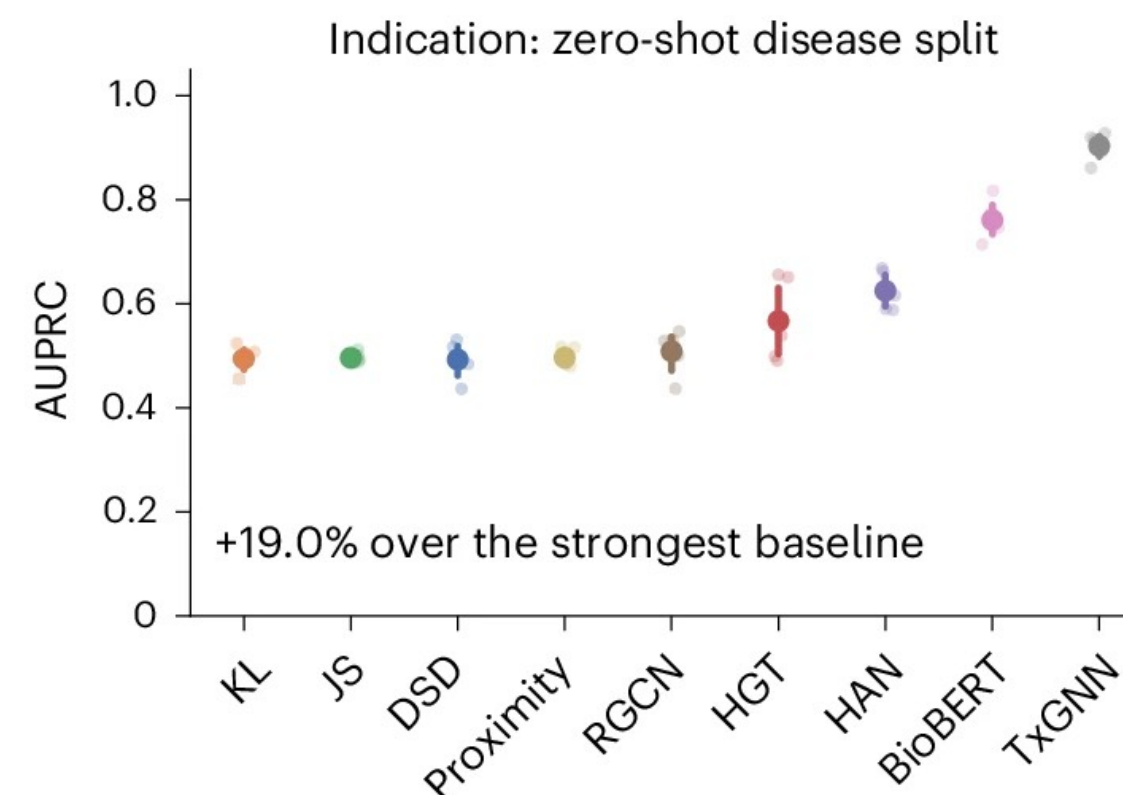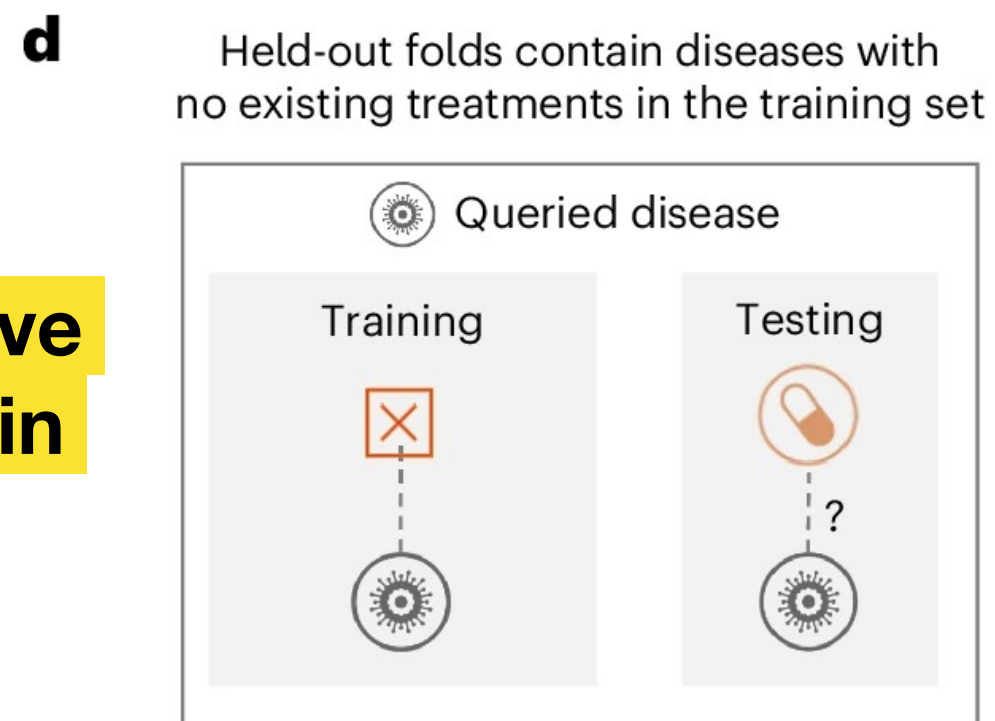GNN built on KG of drugs, disease, pathways, and more oh my!

Two held-out situations:

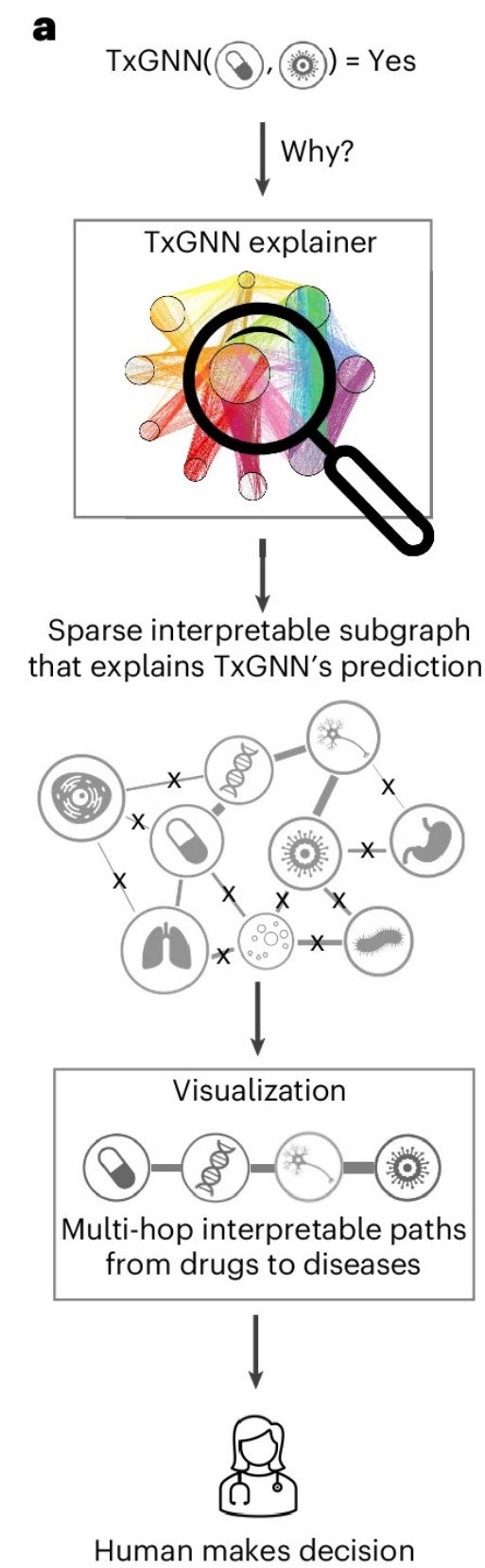Disease *does* have existing treatments in training

Disease *does not* have existing treatments in training

Bigger improvement in situation where no treatment exists

**Built an "explainer" to investigate the evidence for a drug repurposing candidate**

**Built an "explainer" to investigate the evidence for a drug repurposing candidate**

**Example of the path between a disease and drug that can be explored**

**Used EHR data to corroborate some drug repurposing hypotheses**

**a** Medication information in EMRs → Inclusion criteria → Calculation of log(OR) for all drug–disease pairs → Evaluation against FDA-approved indications → Evaluation against log(OR)

TxGNN / Medical records / Patient ID: / Drug ID: / Disease ID:

478 diseases with ≥ 1 patients
1,290 drugs with ≥ 10 patients
1,272,085 patients with at least 1 drug and at least 1 disease

FDA-approved indications / All drug–disease pairs

Density / log(OR)

0.73 ↔ 3.51
0.45 ↔ 2.14
0.21 ↔ 1.21
...
0.01 ↔ 0.15

Prediction log(OR)

**b** Pie chart — race:
Asian 4.7%, Unknown 19.0%, Black 12.7%, Other 21.5%, White 42.0%

**c** Pie chart — sex:
Male 40.1%, Female 59.9%

**d** Number of diseases EHR Data:
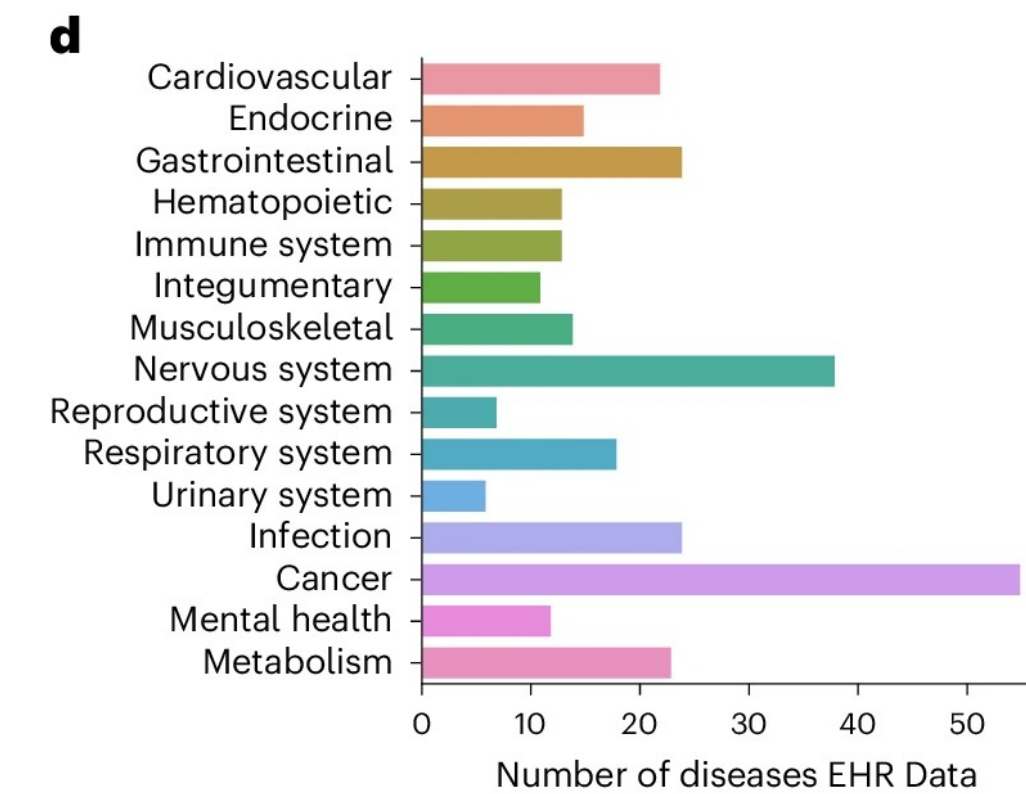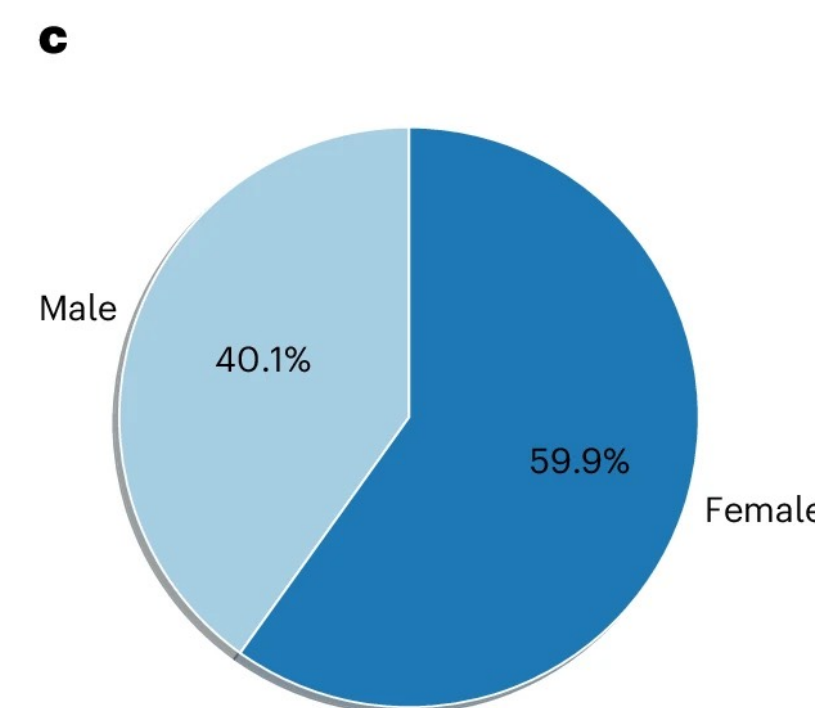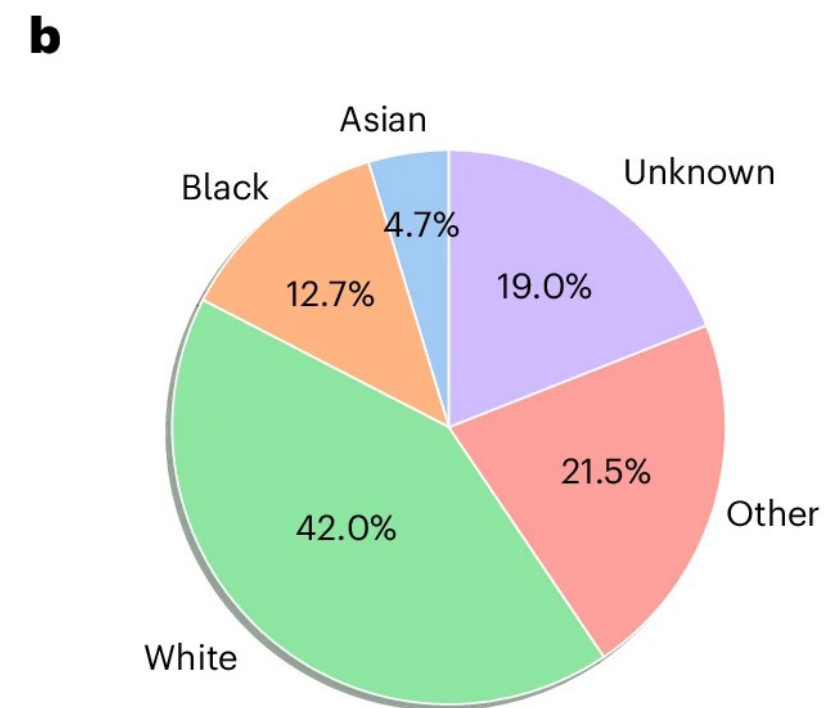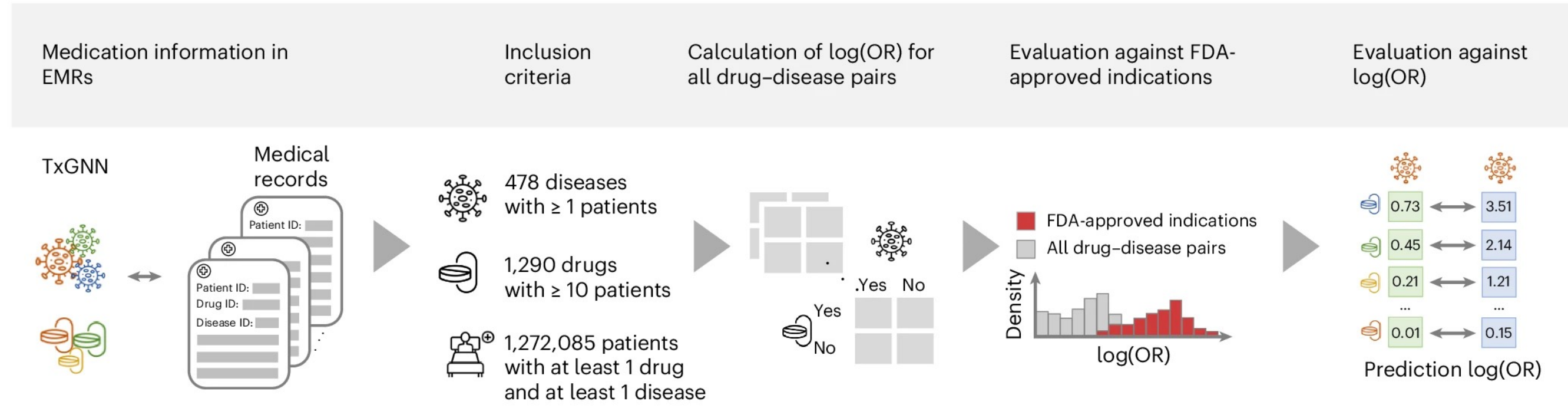Cardiovascular, Endocrine, Gastrointestinal, Hematopoietic, Immune system, Integumentary, Musculoskeletal, Nervous system, Reproductive system, Respiratory system, Urinary system, Infection, Cancer, Mental health, Metabolism

**See generally a shift between the OR for approved drugs and contraindications**

**e** Density vs log(OR):
FDA-approved indications (red), Contraindications (green), All drug–disease pairs (grey)

**f** log(OR) — Predictions for 478 diseases:
FDA-approved indications, Contraindications
Top 1 drug (~2.25), Top 5 drugs (~1.9), Top 5% drugs (~1.3), Bottom 50% drugs (~1.1), +107%

**g** Wilson's disease:
Predicted indication likelihood vs log(OR)
Deferasirox TxGNN: 0.79 log(OR): 5.26

# ADMET-AI Enables Interpretable Predictions of Drug-Induced Cardiotoxicity (Mukherjee, Swanson et al, *Circulation*)

- Goal: Predict drug-induced cardiotoxcity pre-clinically and identify causal factors

- Method:

  - Use ADME-AI to generate 41 ADMET properties

  - Feed those into Extreme Gradient Boosting to predict cardiotoxicity

- Result:

  - AUROC = 0.72 w/ top predicted features: CYP2D6 metabolism, Nrf2-antioxidant response, aromatase inhibition

  - Classify drugs into three categories: safe, high risk, withdrawn

- Conclusion: Predicting drug effects continues to be very hard.

**Drugs ranked and labels by concern in DCITRank**

**Slightly better performance**

**Most important features make sense**

CYP2D6 metabolizes:
β-Blockers → **Metoprolol, Carvedilol, Propranolol** (heart failure, hypertension)
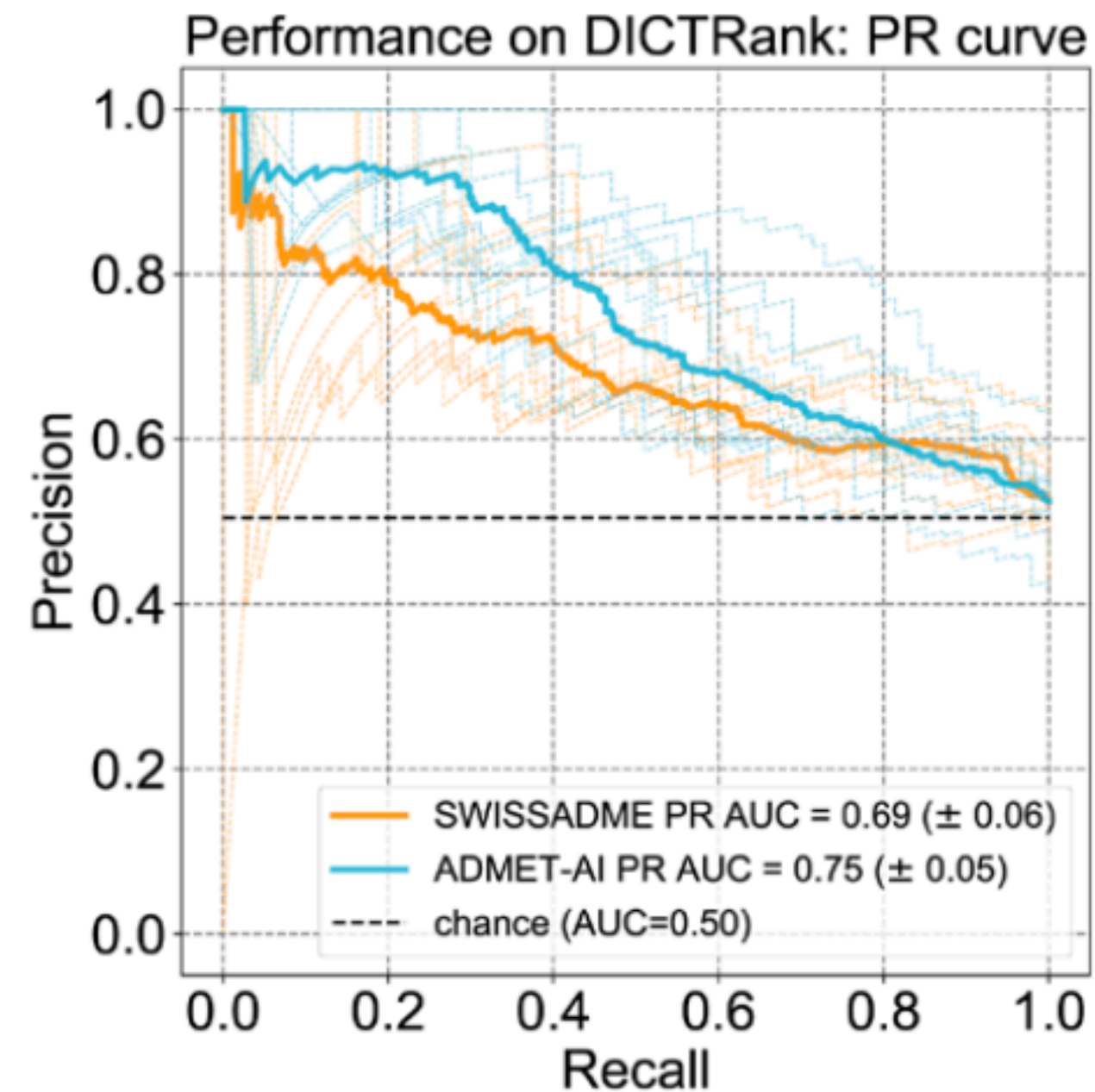Antiarrhythmics → **Flecainide, Encainide** (arrhythmias)
Calcium channel blockers → Some **Dihydropyridines** (e.g., Nifedipine)
ACE inhibitors & ARBs → Some undergo minor CYP2D6 metabolism

**Cute way to show that these are important features**



E — Shapley values to explain model co-efficients

G — DICT concern: none

H — DICT concern: most

I — Withdrawn drugs

# Pan-cancer proteogenomics expands the landscape of therapeutic targets (Savage et al, *Cell*)

- Goal: To identify new therapeutic targets by integrating pan-cancer proteogenomic data

- Method:

  - Multifaceted and multimodal computational analysis strategy

    - Link to drug target databases; use synthetic lethal data to find pairs of cancer drivers that could be targets; identify hyper expressed or overactive proteins; prioritized neoantigens
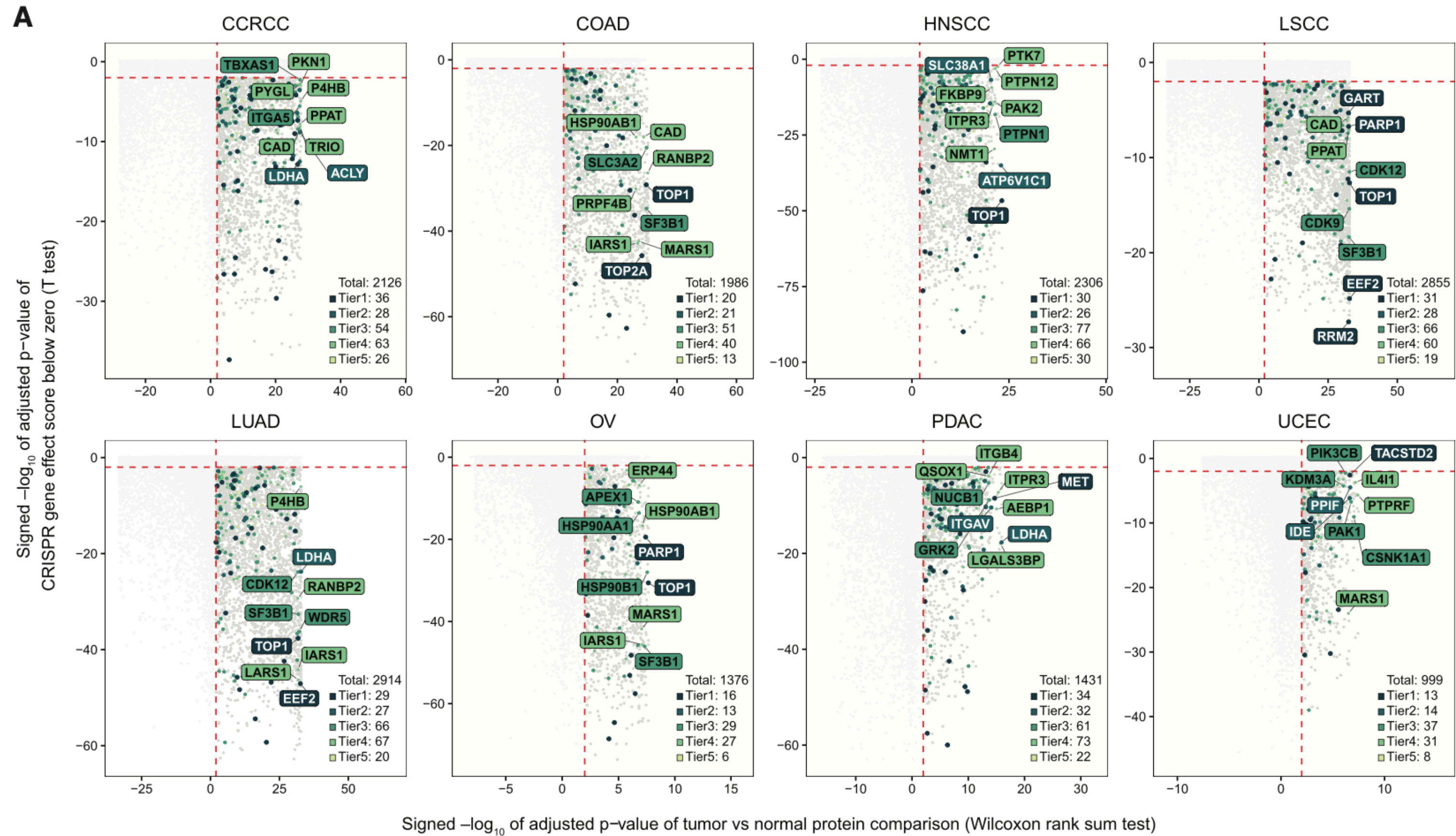
  - Couple with experimental validation (binding affinity, SL screens, etc)

- Result: Identified and characterized 2,863 druggable proteins across five target tiers

- Conclusion: Multimodal data integration makes a stronger case for these putative targets.

#IS25    #YIR25    𝕏 @proftatonetti    🦋 @tatonetti.bsky.social    https://doi.org/10.1016/j.cell.2024.05.039

Identify upregulated genes as potential targets by cancer type

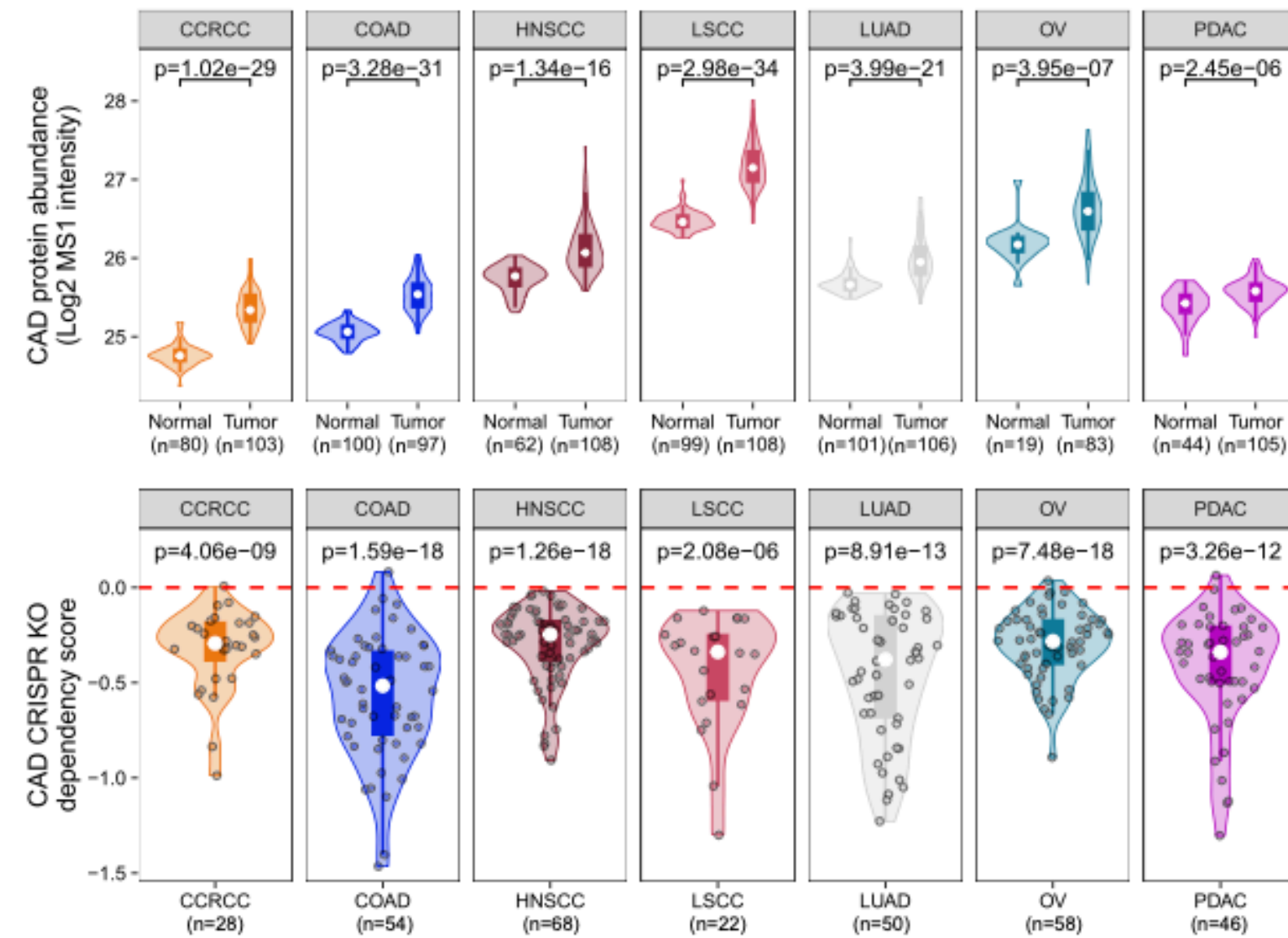Non-pan-essential targets shared by 5+ cancer types

**Confirmed putative targets are upregulated in tumor (vs. normal) as predicted**
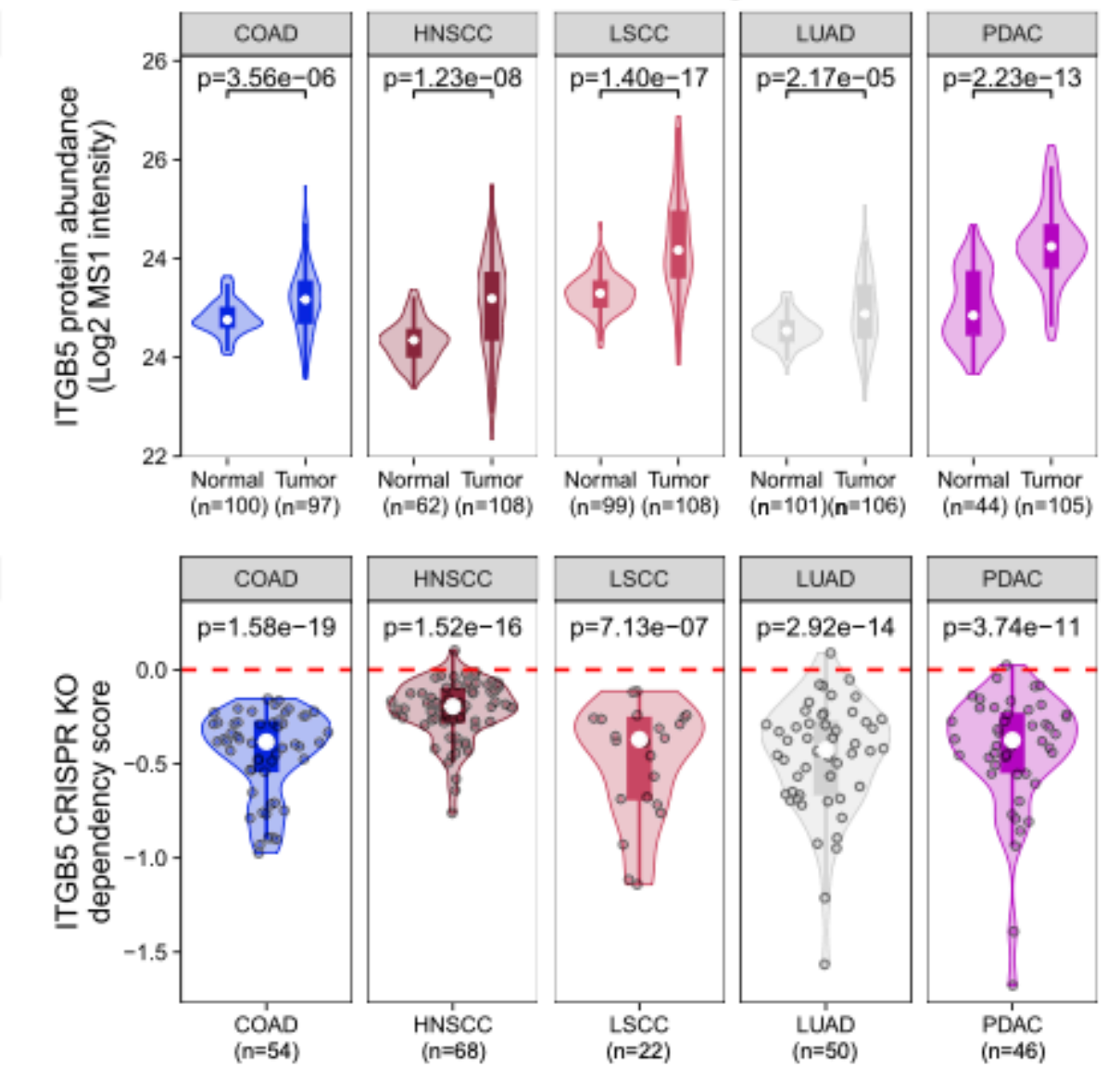

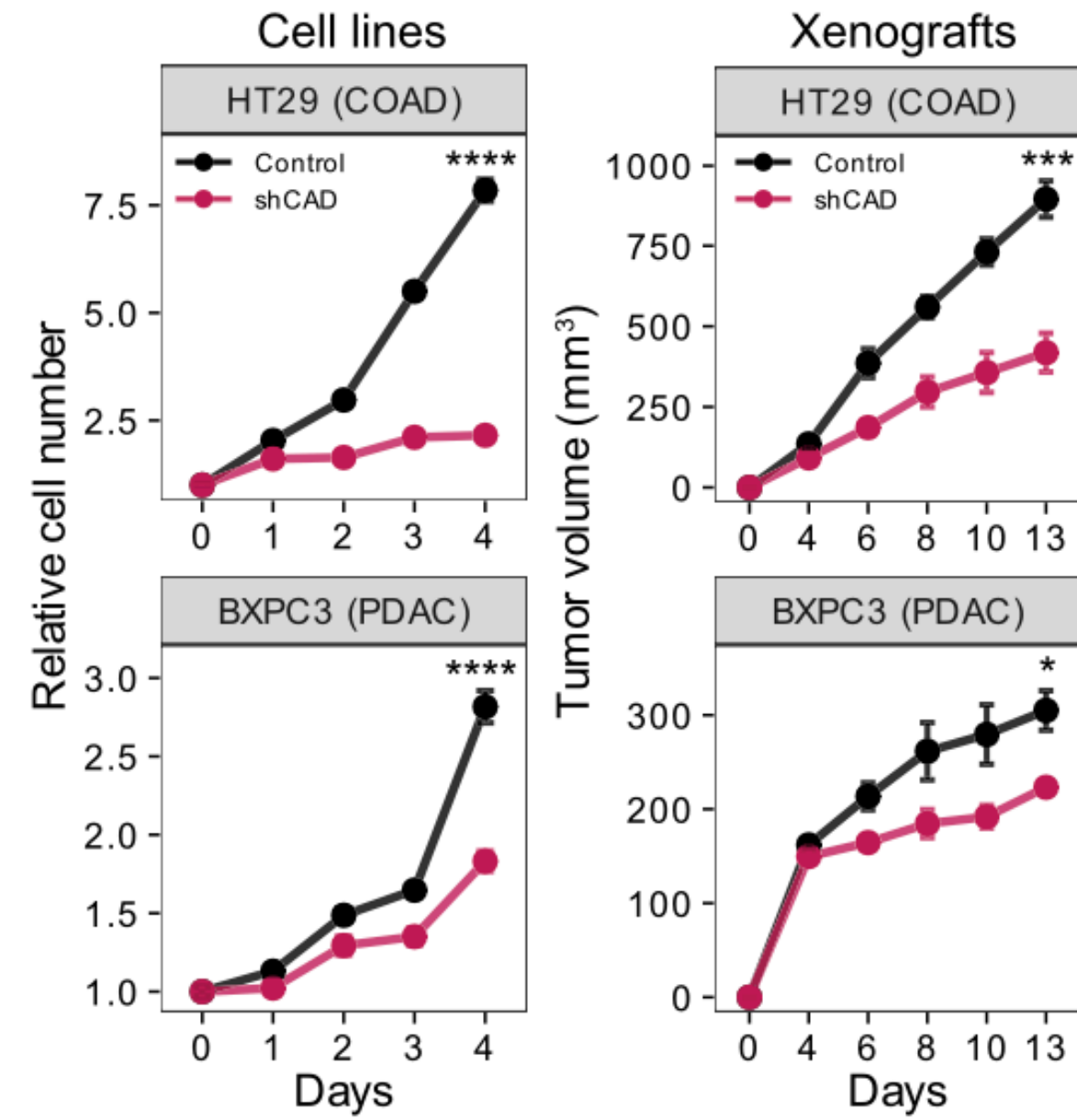
**E**  Tier 4 target: CAD

**F**  Tier 4 target: PAK2

**G**  Tier 4 membrane target: ITGB5

**And so with CRISPR that these cells are dependent on these genes**
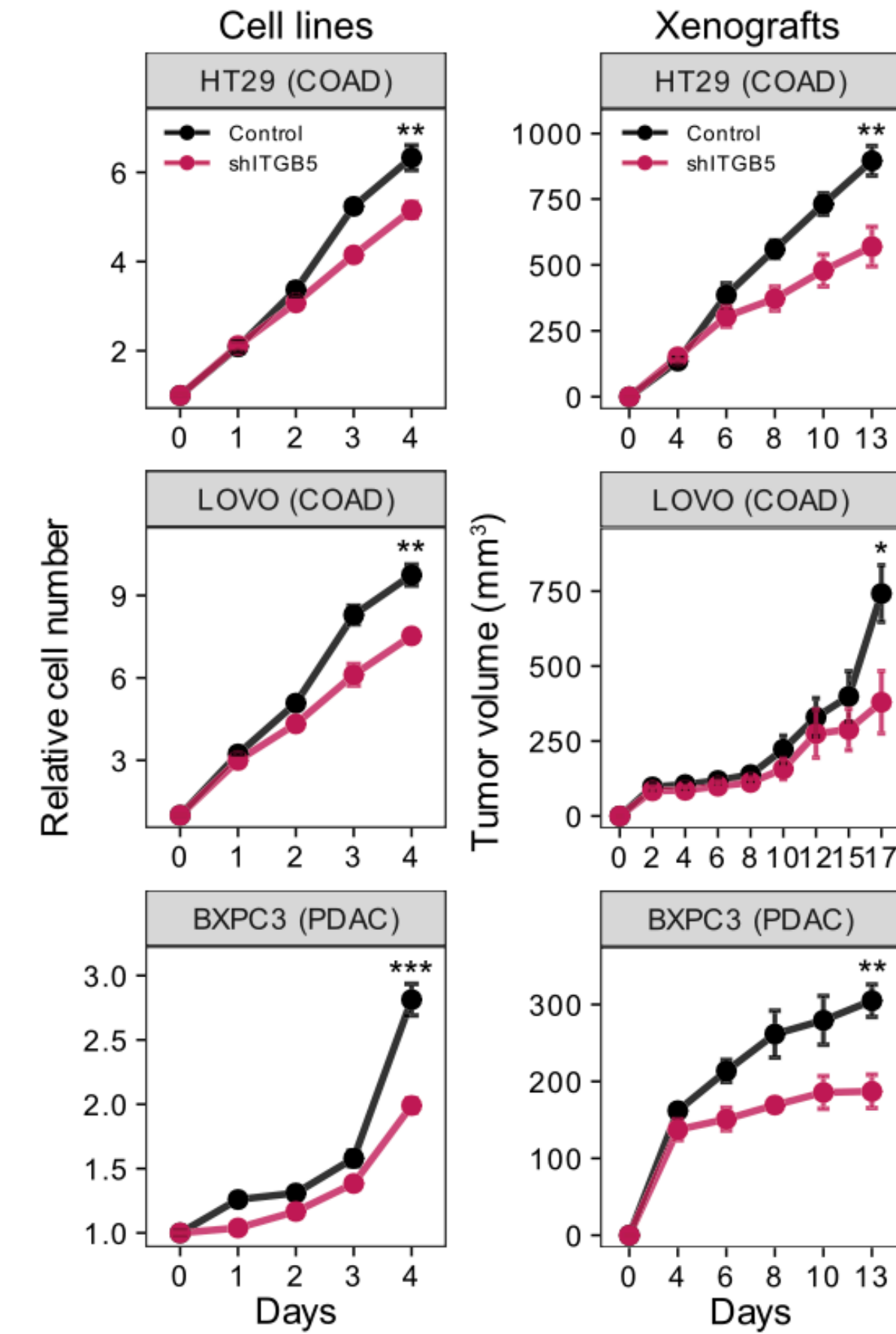
Knock down experiments in mouse xenografts

Knocking down these targets inhibits growth in pretty much every test

# Trainee Spotlight



**PUSHKALA JAYARAMAN**
ICAHN SCHOOL OF MEDICINE MOUNT SINAI, NY

# AutoEdge-CCP: A novel approach for predicting cancer-associated circRNAs and drugs based on automated edge embedding (Chen et al, *PLoS Comp Bio*)

- Goal: Develop a computational framework to predict associations between circular RNAs (circRNAs), drugs, and cancer

- Method:

  - Multi-source heterogenous network (circRNA, drugs, and cancer)

  - Use GNNs to build embeddings

  - Uses a Learn-to-Rank (LTR) framework to learn relationship between circRNA and cancer (and drug-cancer)

- Result:

  - Big improvements over baselines (AUROC of 0.989 vs 0.700)

  - Edge embeddings offer mechanistic interpretations

- Conclusion: The future is joint…embeddings

**AutoEdge-CCP performs better overall than other methods**

**Classifying cancer type**

**Predicting missing (held out) edges**

**Both node and edge embeddings are needed**

**Novel drug-cancer predictions held up in computational docking experiments**



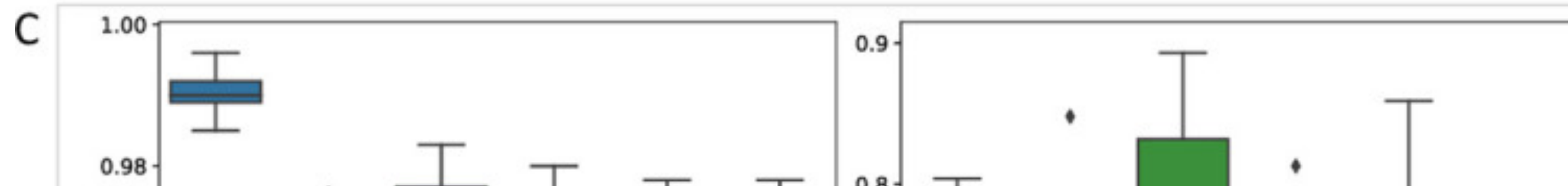| Cancer | Target | | | Binding |
| | Protein | PBD ID | Reference | energy (Kcal/mol) |
|---|---|---|---|---|
| Esophageal Squamous Cell Carcinoma | TGF-beta receptor type-2 (TGFBR2) | 5E8Y | [36] | -7.06 |
| | Cellular tumor antigen p53 (TP53) | 4ZZJ | [37] | -5.85 |
| | Polyunsaturated fatty acid lipoxygenase (ALOX12) | 3D3L | [38] | -4.87 |
| Colorectal cancer | Mothers against decapentaplegic homolog 4 (SMAD4) | 1G88 | [40] | -5.99 |
| | Catenin beta-1 (CTNNB1) | 1P22 | [41] | -4.59 |
| | DNA mismatch repair protein Mlh1 (MLH1) | 6WBB | [42] | -4.47 |

# scRank infers drug-responsive cell types from untreated scRNA-seq data using a target-perturbed gene regulatory network (Li et al, *Cell Reports Medicine*)

- Goal: Identify drug-responsive cell populations without requiring post-treatment transcriptomic data

- Method:

  - Use scRNA-seq data to generate a treatment naive regulatory network

  - Simulate drug perturbation by removing the target from the network (<u>assumes inhibition and only one target</u>)

  - Use manifold alignment and network diffusion to measure the global and local effects on the network

- Result:

  - Outperforms existing methods (e.g. Augur and DEG) in simulations and real data

  - Identified well known cell type-drug associations from the literature

  - Experimental validation three different disease contexts

- Conclusion: I didn't realize we could still do single cell analysis without deep learning

**Input is a bunch of single cell RNAseq data and knowledge on drug targets**

**The manifold alignment maps to a lower dimension and measures distance — gives you a global effect**

**Output is a ranking of cell types**

**The diffusion metric maps out to the drugs neighbors, giving you a local effect**

**A** Input
▶ Untreated scRNA-seq data
▶ Drug direct target

Drug
Inhibit
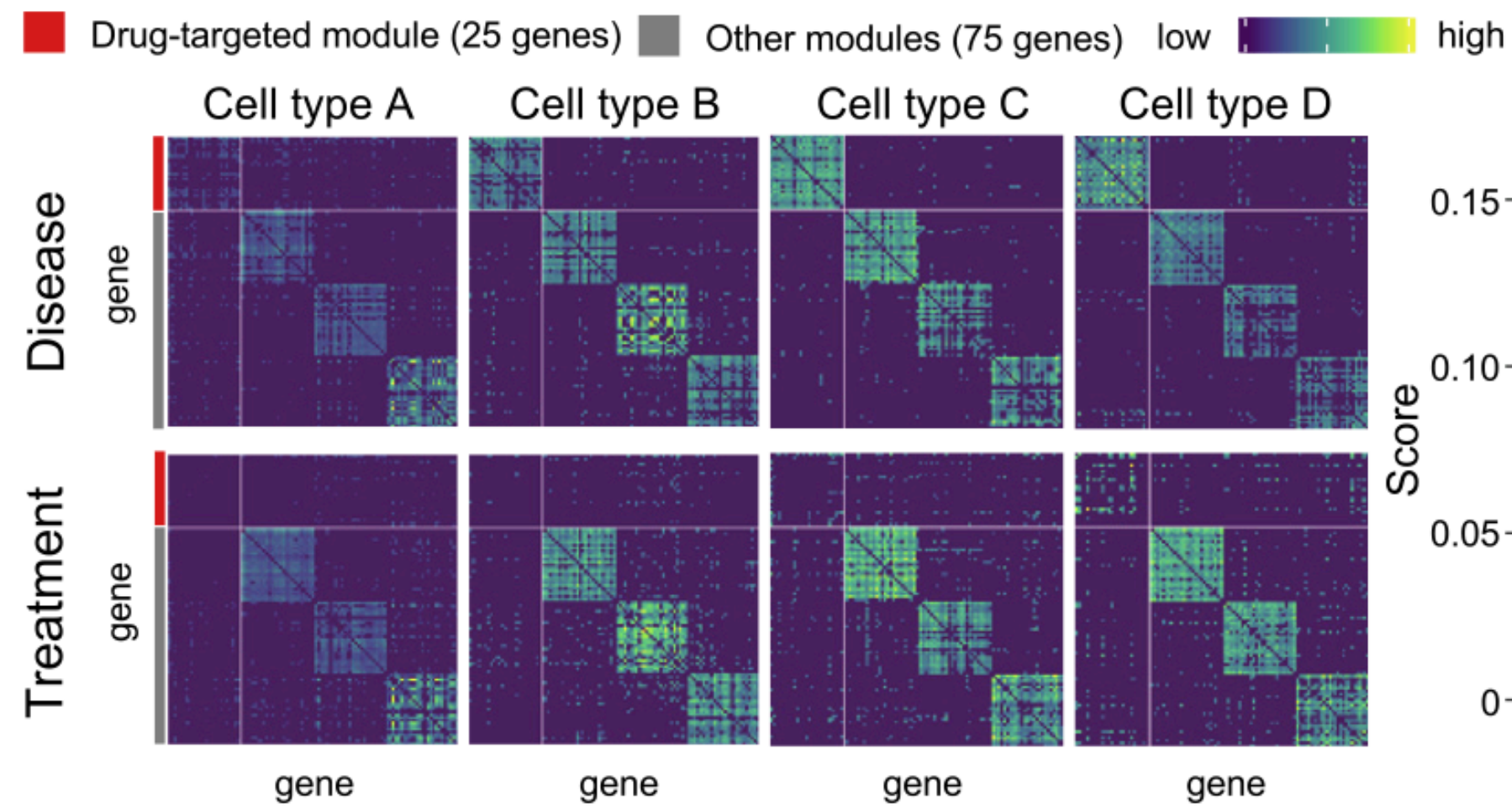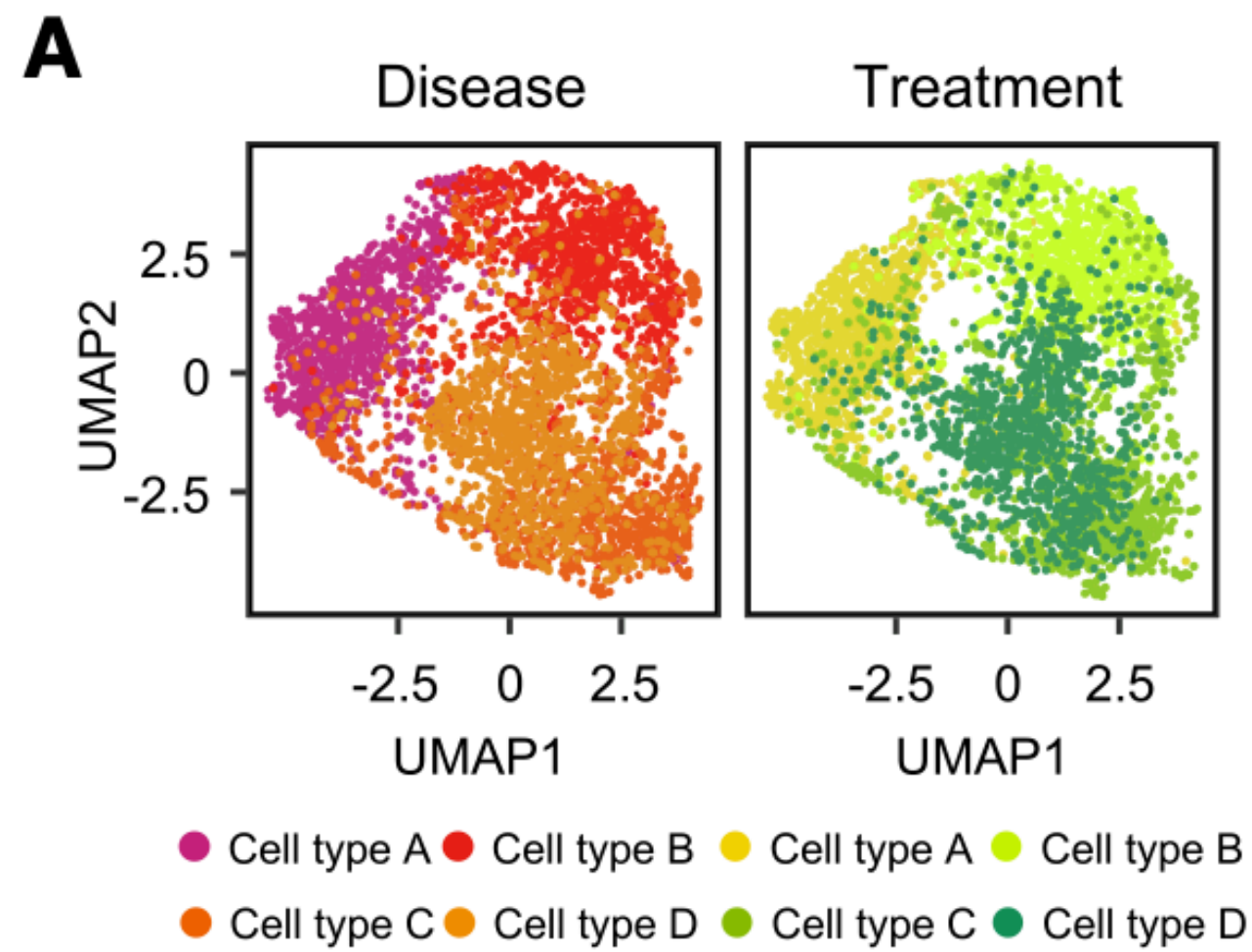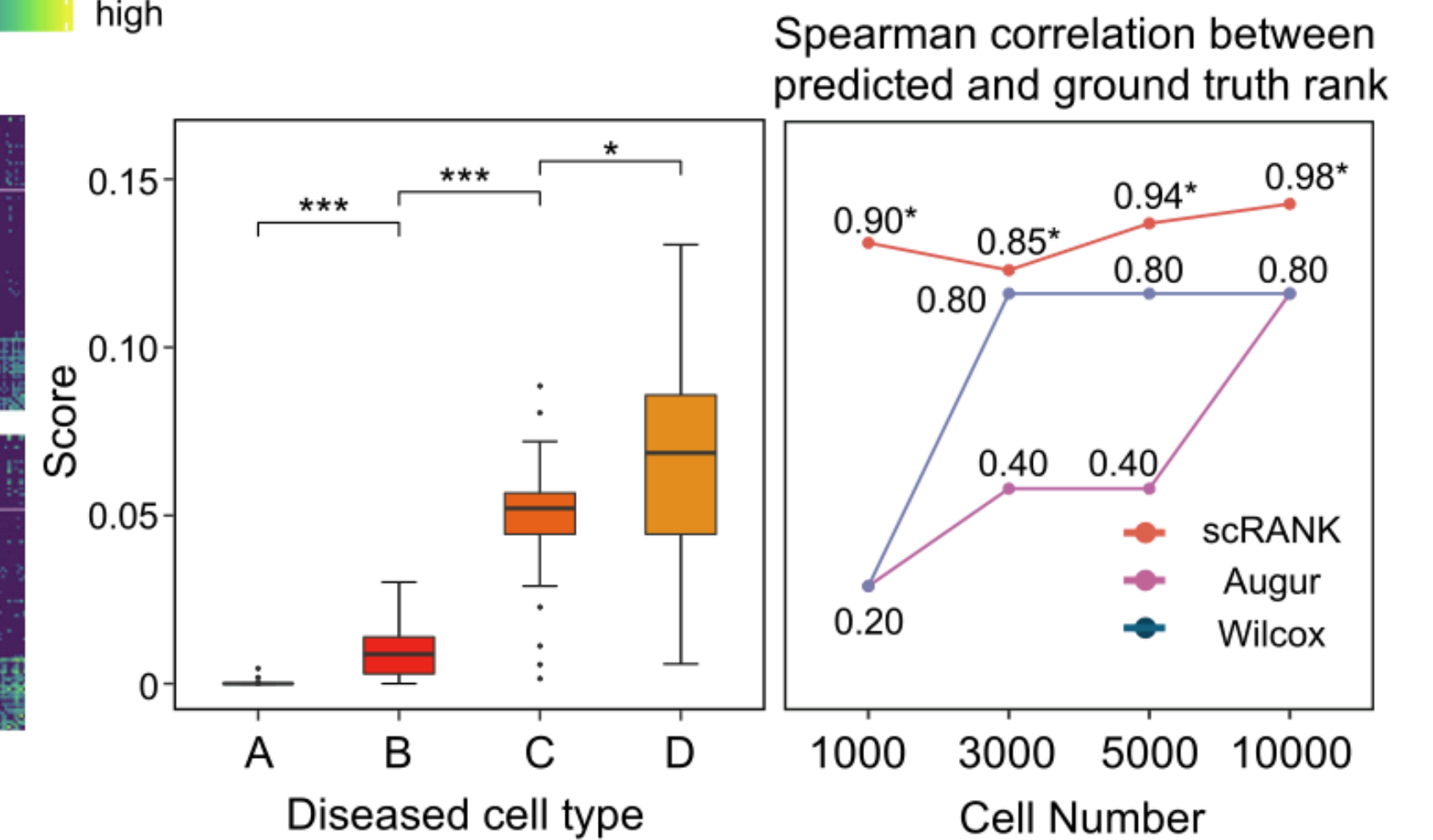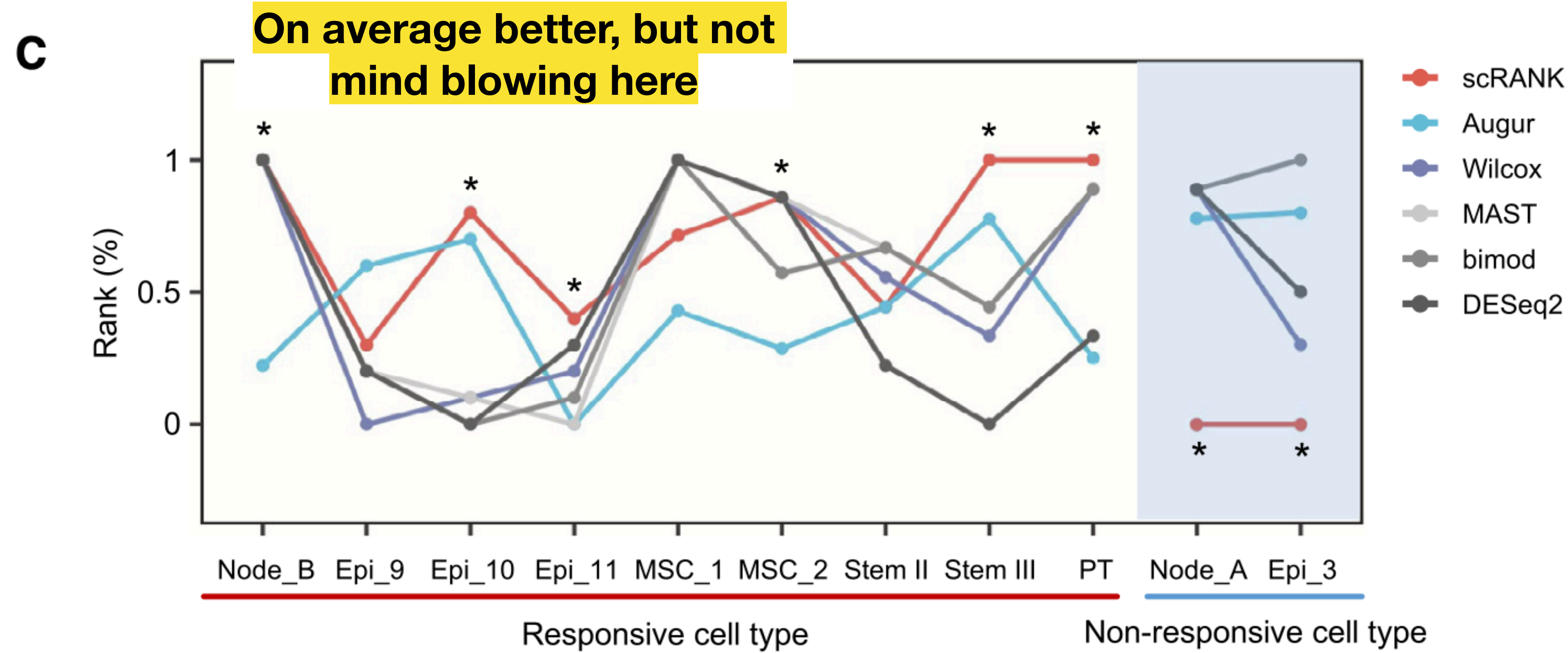Cell type 1   Cell type 2
Cell type 3   Cell type 4
Target

**B** Cell type ranking based on tpGRN
▶ Untreated cell type-specific GRN
▶ Cell type prioritization
Gene

Target
Untreated GRN
Gene in untreated GRN
Gene in tpGRN
Manifold Alignment
G1 distance  G1
G2 distance  G2
Dim2
Dim1

**C** Output

| Cell type | Rank |
|---|---|

**D** (i) Construct GRN
[Drug target genes] + [HVGs] + [TF genes]
Genes

| | Cells | 1 | 3 | · |
| | | 2 | 0 | · |
| | | · | · | · |

PCR
Genes
Genes
Cell type-specific GRN

**E** (ii) Calculate distance of gene nodes between GRNs
Untreated GRN
Manifold space shared by two GRNs
Manifold Alignment
tpGRN
Dim 2
Dim 1
Drug target

Accumulated distance of target-centered neighbor nodes
Disease   Target perturbation
Gene 1
Gene 2
Gene 3
Euclidean distance
Gene node embedding

**F** (iii) Quantify the target-centric diffusion effect
2-hop propagation

Accumulated distance of target-centered neighbor nodes

$$D_d W_{out} / Deg_d + \sum_n^N D_n W_{in} + \sum_n^N D_n W_{out} / Deg_n$$

Target   1-hop nodes   2-hop nodes

Target node
Co-expressed downstream gene node

Simulations

scRANK better captures which cell types are affected

Notice modules are "knocked out"

**Validation in real data**



**On average better, but not mind blowing here**

**Better on average, too**

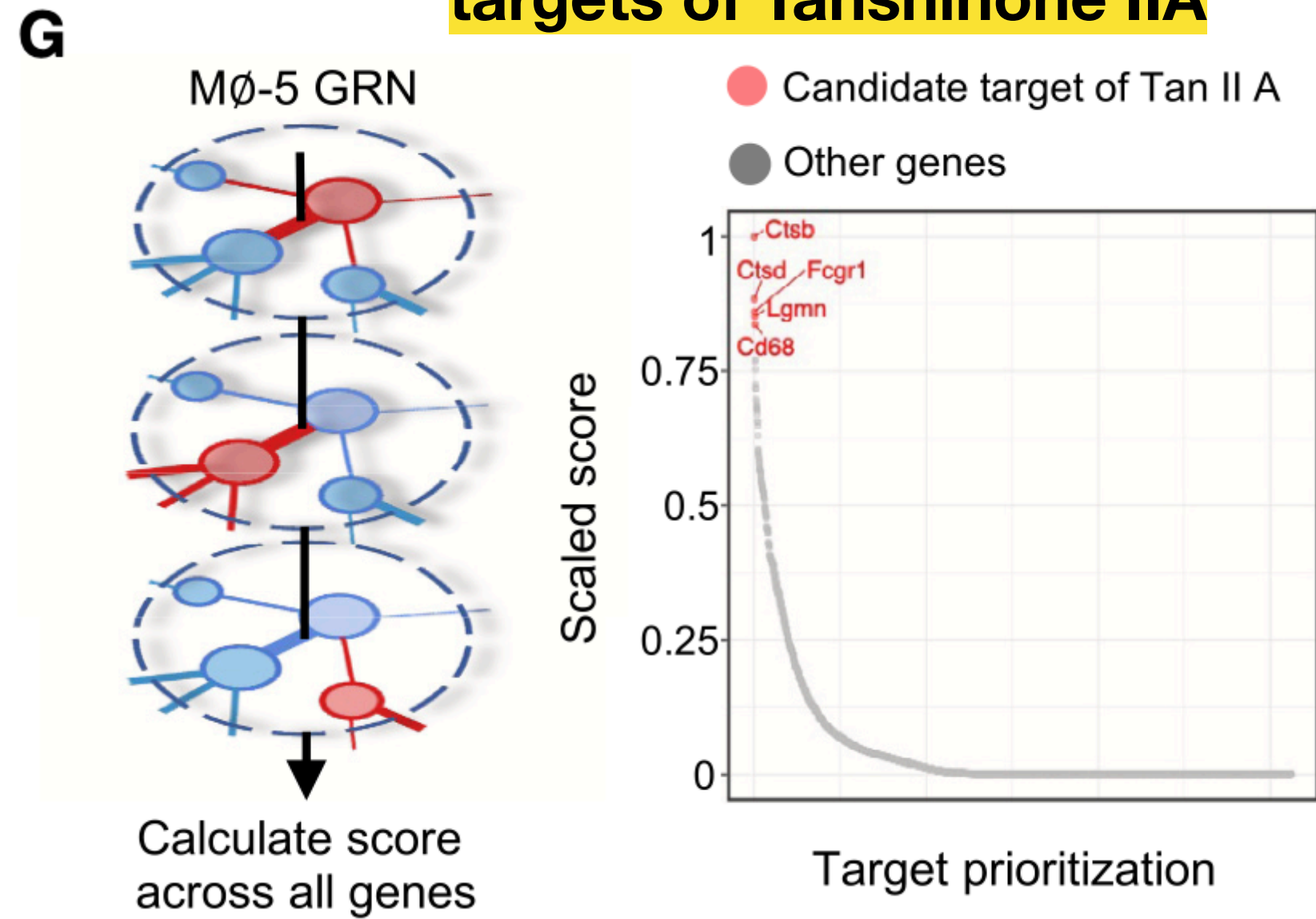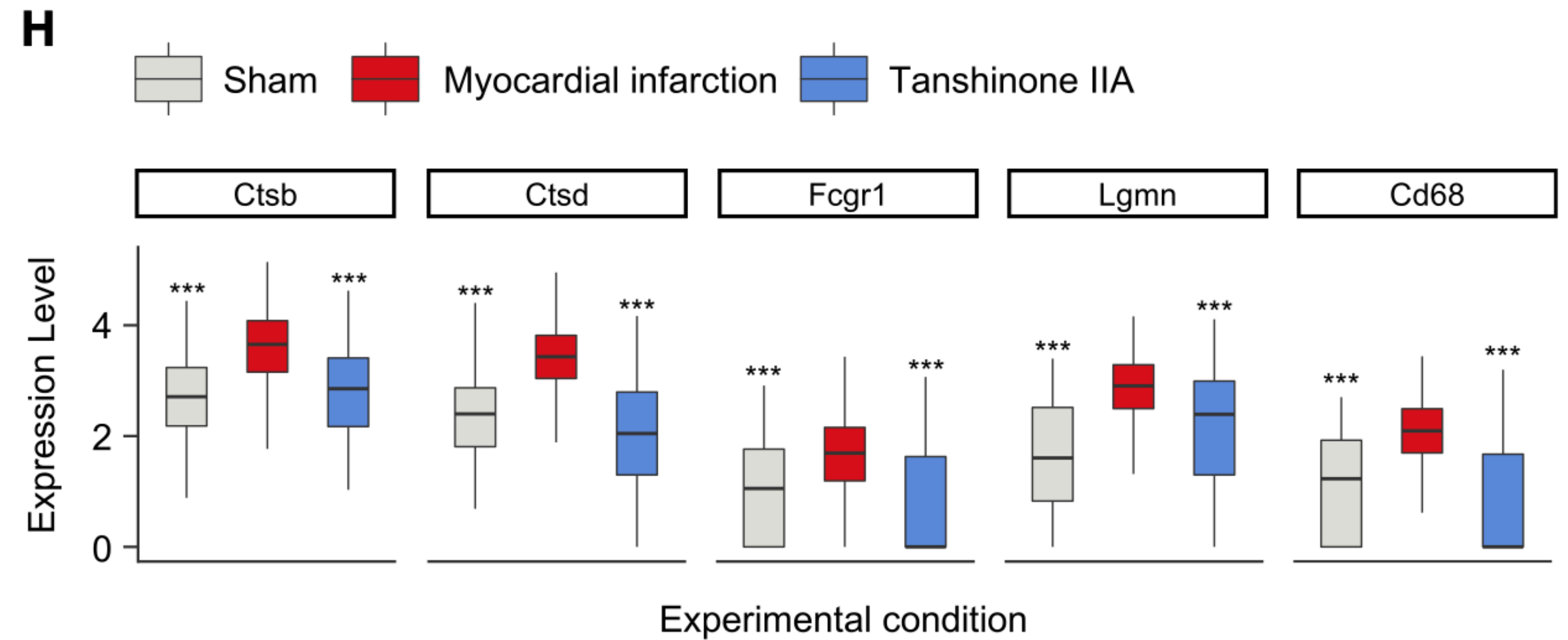*Much* better and ranking non-responsive cell types lower

Diving in on myocardial infarction and tanshinone IIA

Can use scRANK to predict off-targets of Tanshinone IIA

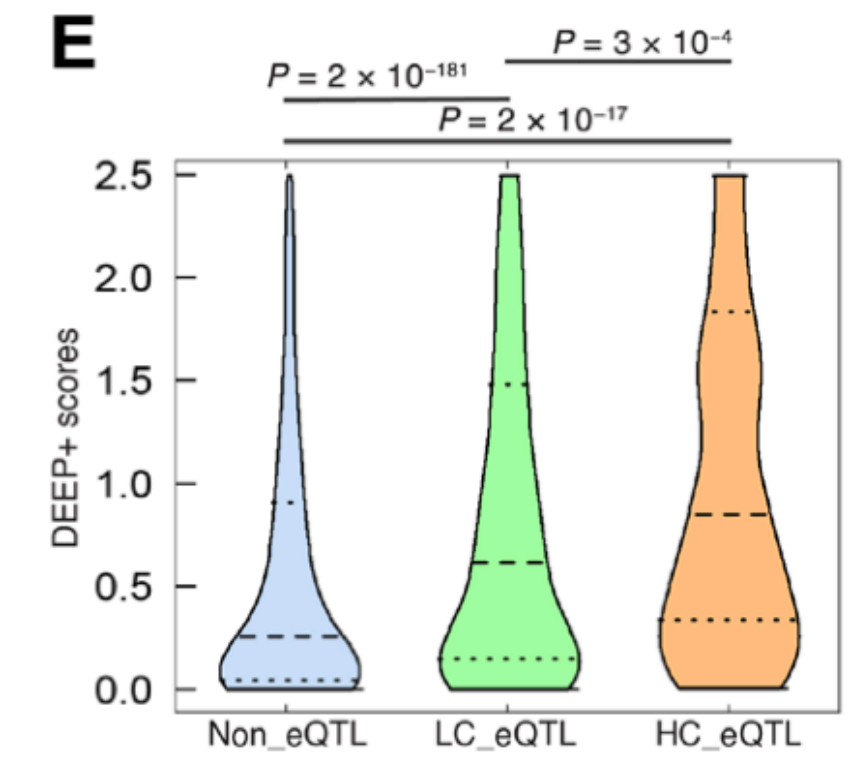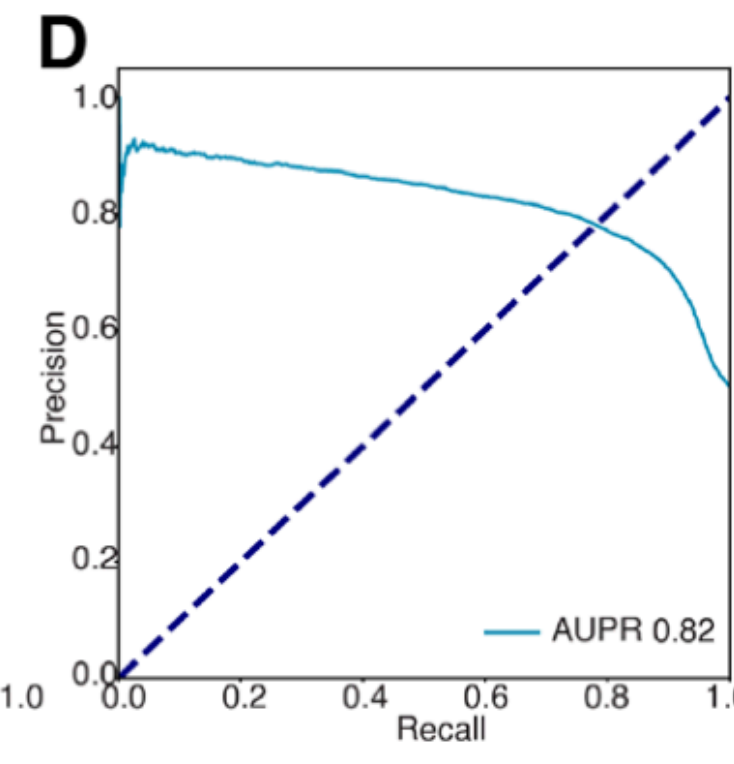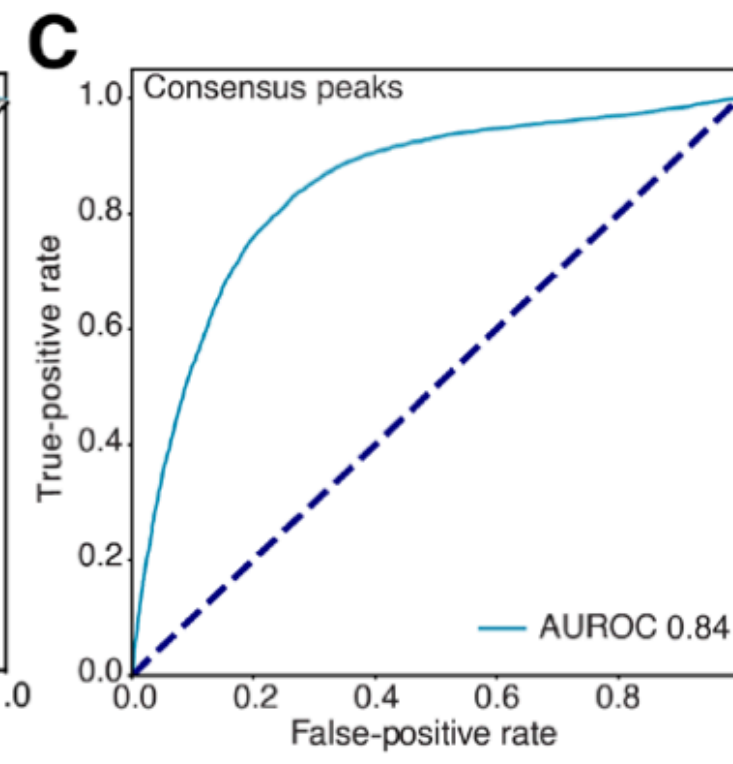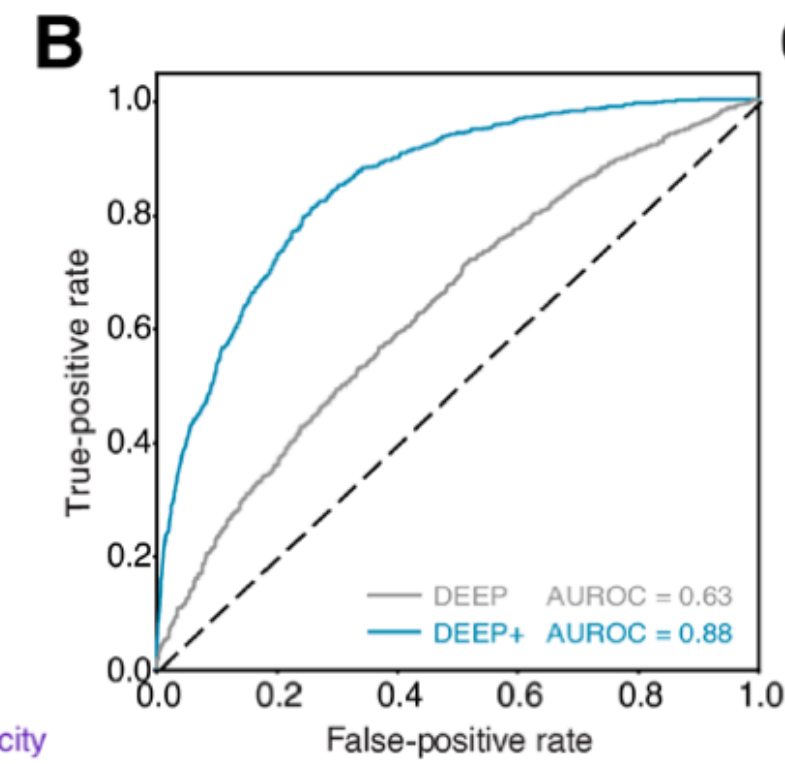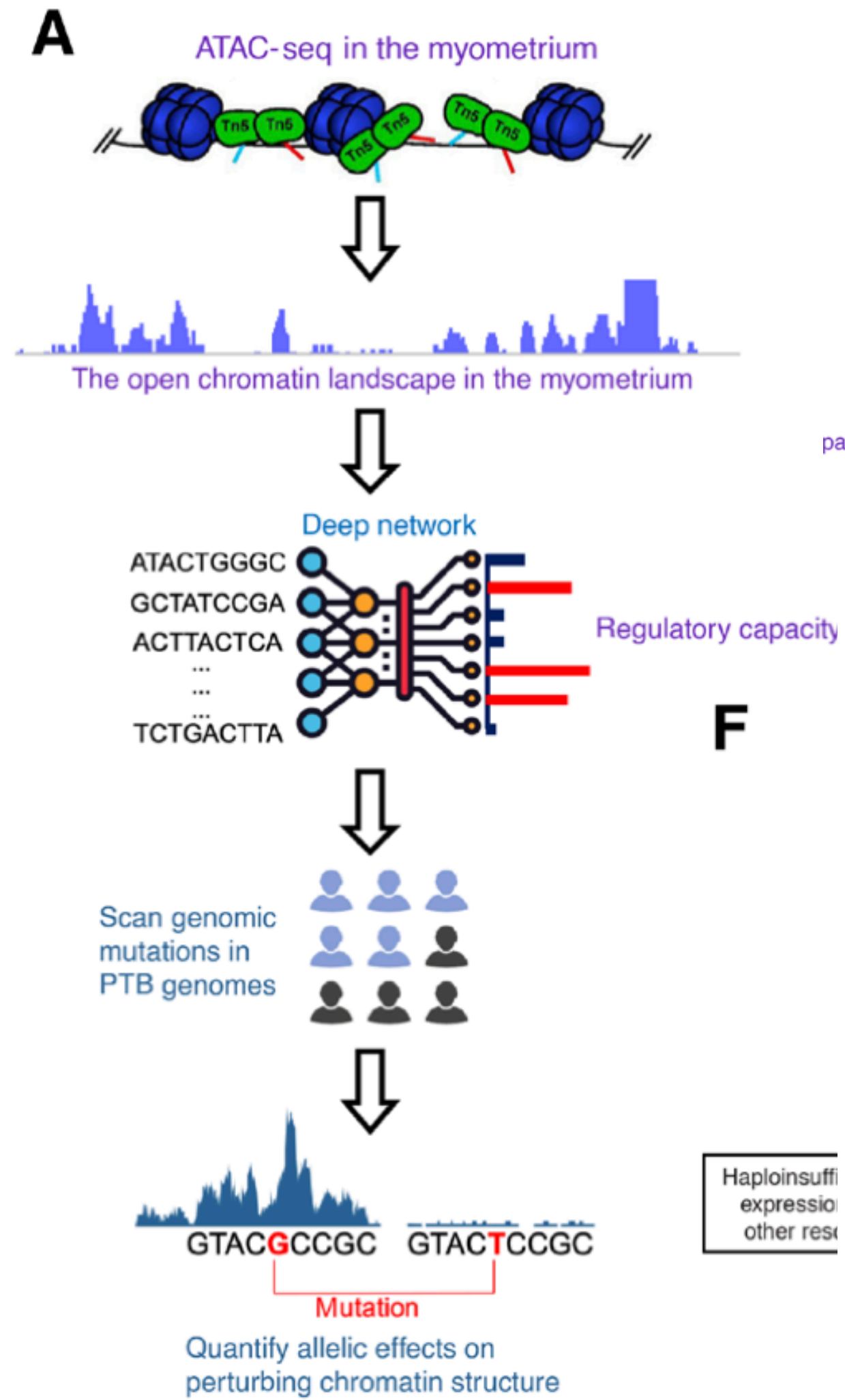And the expression of these are inhibited by Tanshinone IIA in experiments

# Integrative analysis of noncoding mutations identifies the druggable genome in preterm birth (Wang et al, *Science Advances*)

- Goal: Identify genetics of preterm birth; use it explain variation of response to progestin therapy; use it to identify new drug candidates

- Method:

  - Introduce DEEP+ to analyze genomic variants for effect on chromatin accessibility (genetic associations have been in noncoding regions)

  - Develop a bayesian method (BEAR) to integrate output of DEEP+ with epigenetic and transciptomic data to compute a posterior for each genomic locus

    - >High BEAR score means: strong GWAS, disrupts chromatin accessibility, affects dosage-sensitive genes, alters gene expression in the uterus

- Result:

  - Found ~1k genomic loci with high BEAR scores, including those linked to previously unidentified genes affecting muscle relaxation and inflammatory pathways

  - Mutation burden predicted response to progestin therapy

  - Discovered and validated new small molecule to treat spontaneous preterm birth
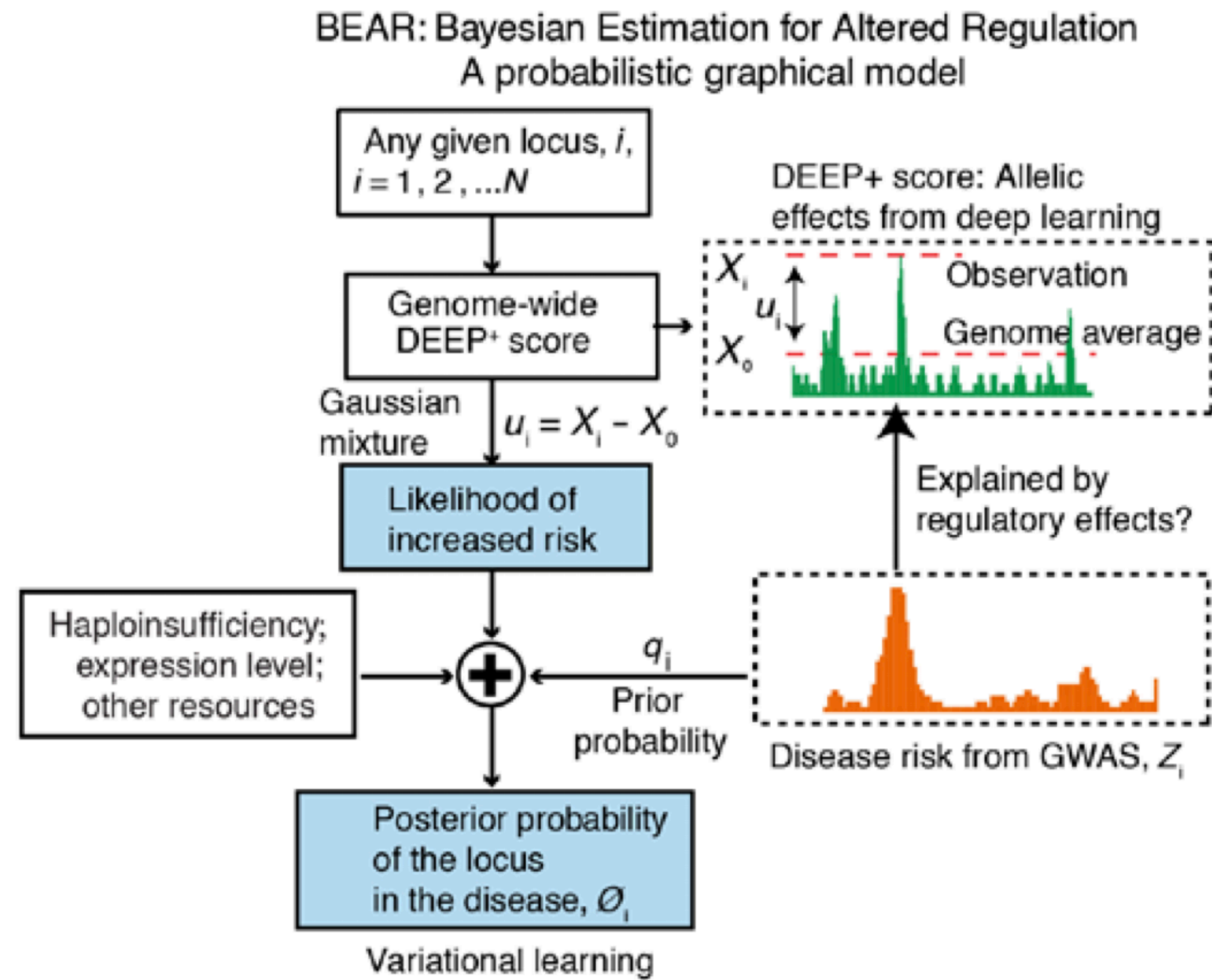
- Conclusion: 👏

**DEEP+ setup**

**DEEP+ improves on DEEP
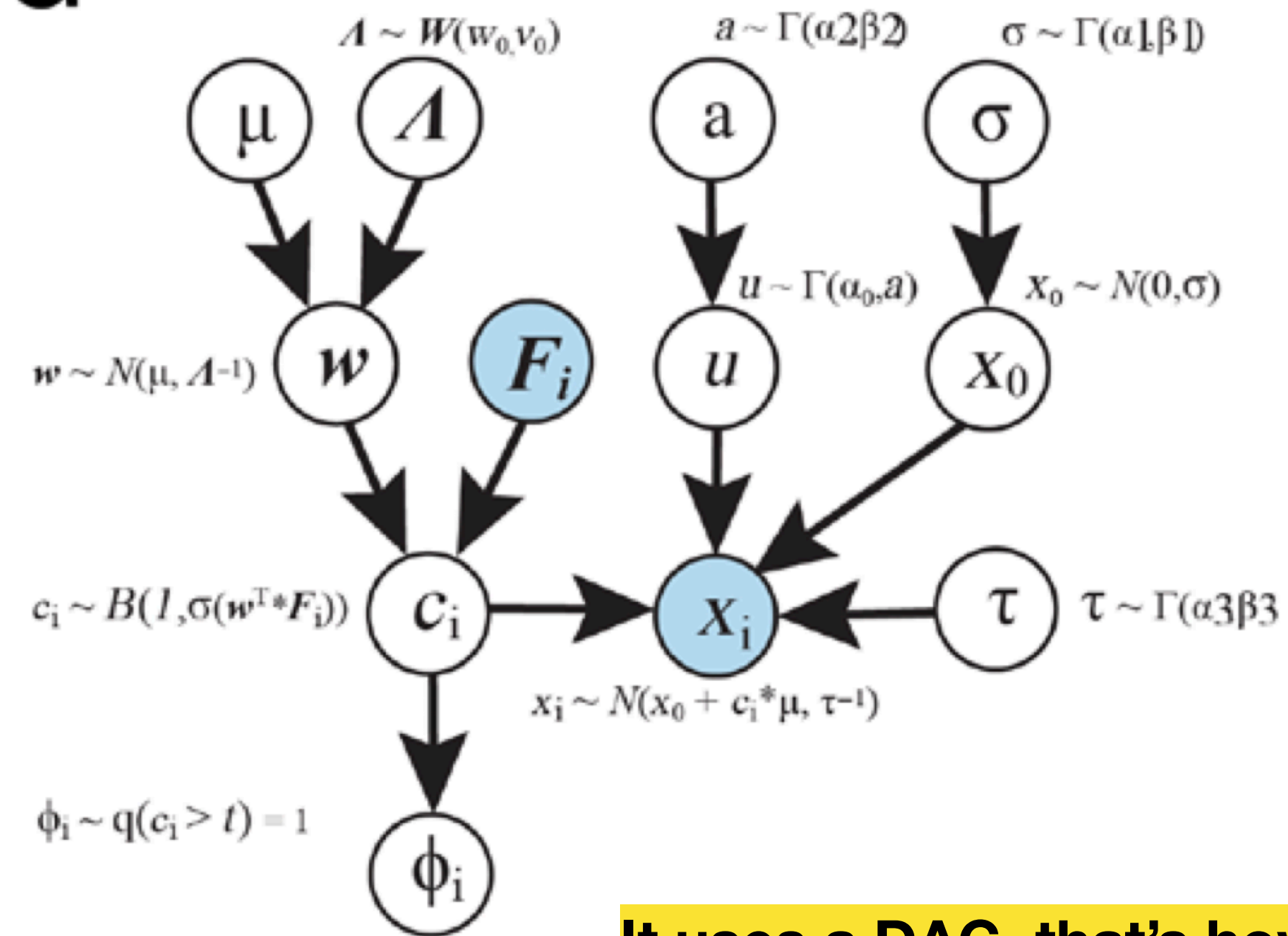(s/CNN/ResNet/g and some other stuff)**



**High eQTLs have higher DEEP+ scores**
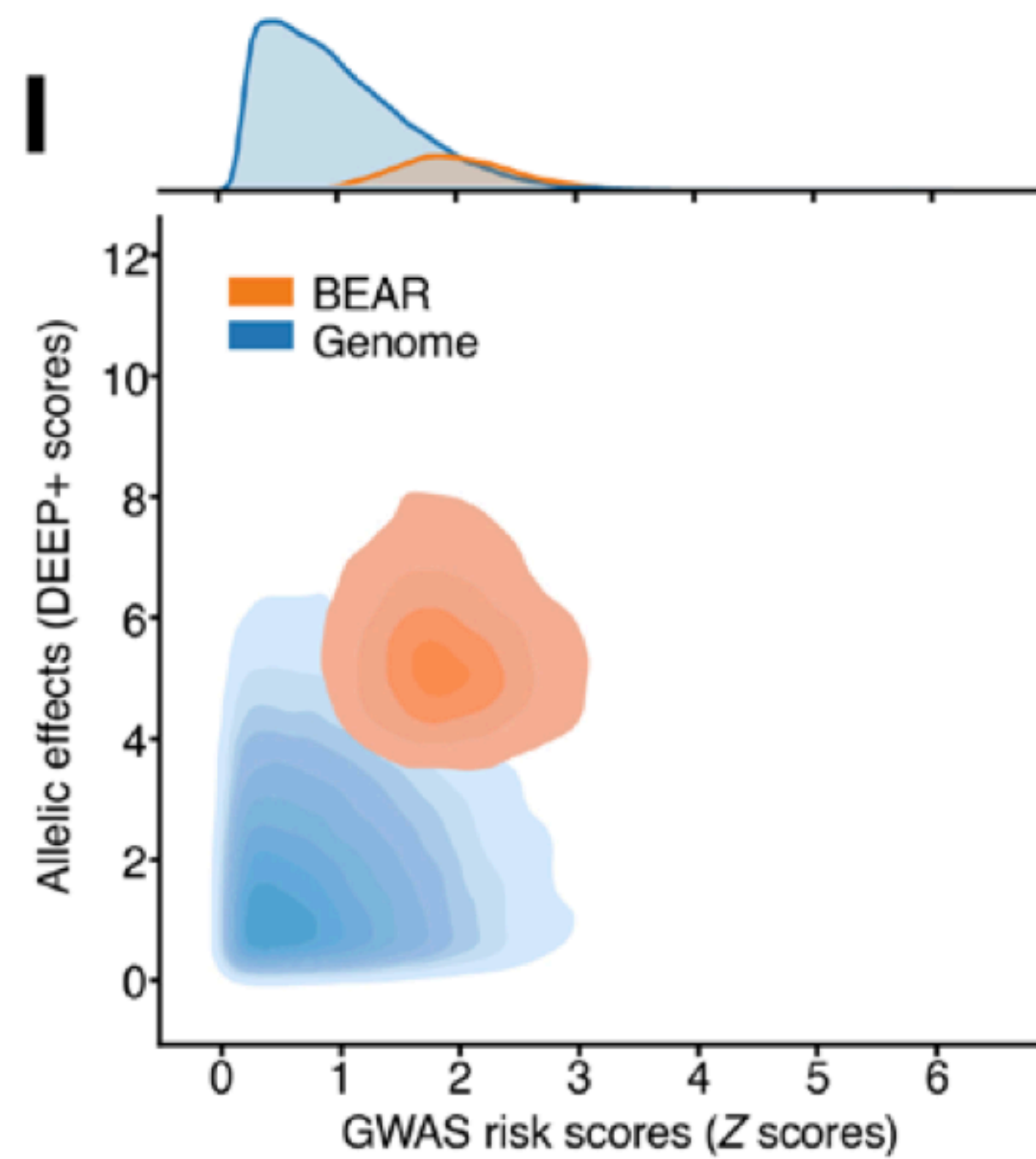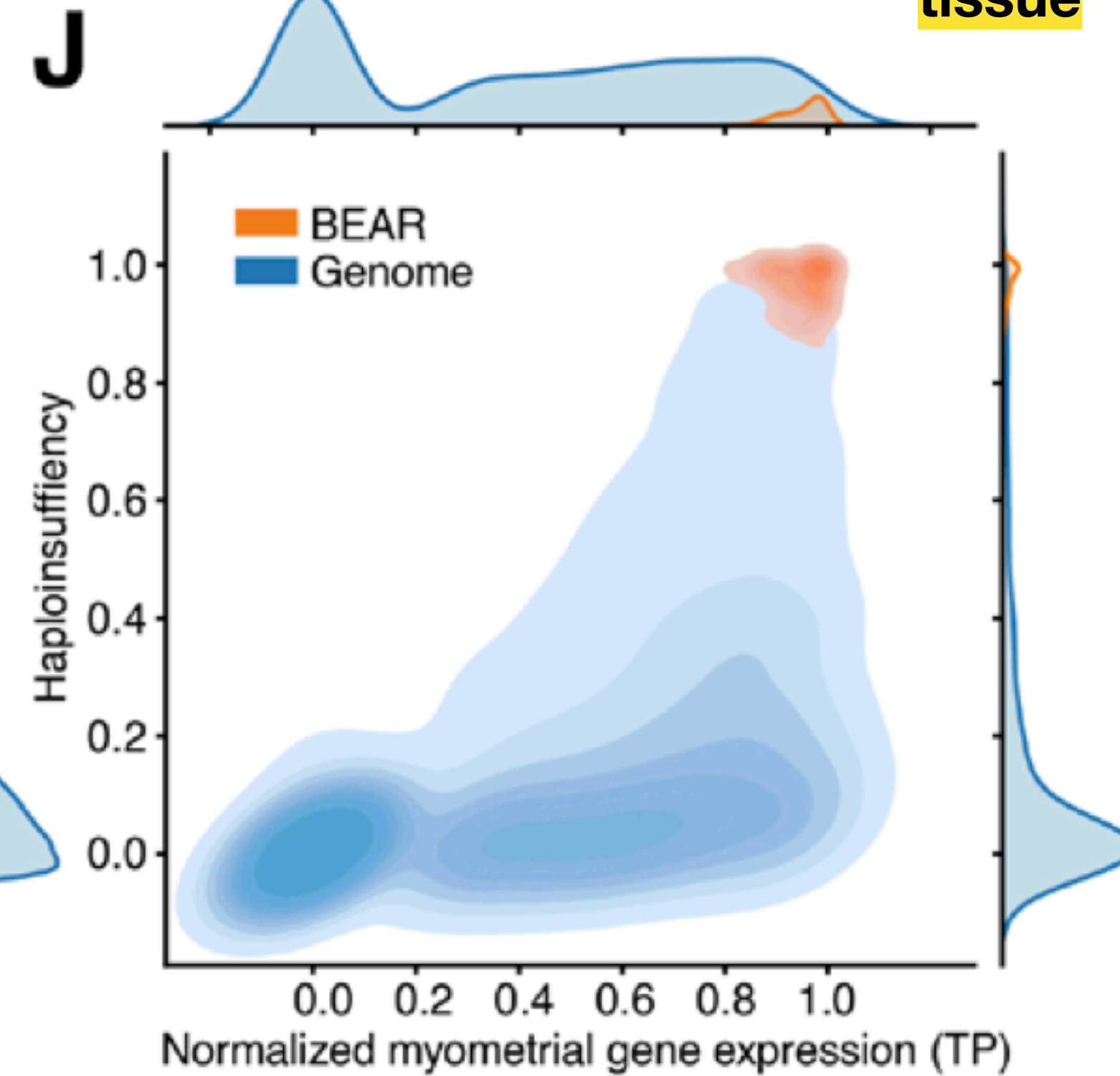
BEAR is a bayesian multimodal integration method



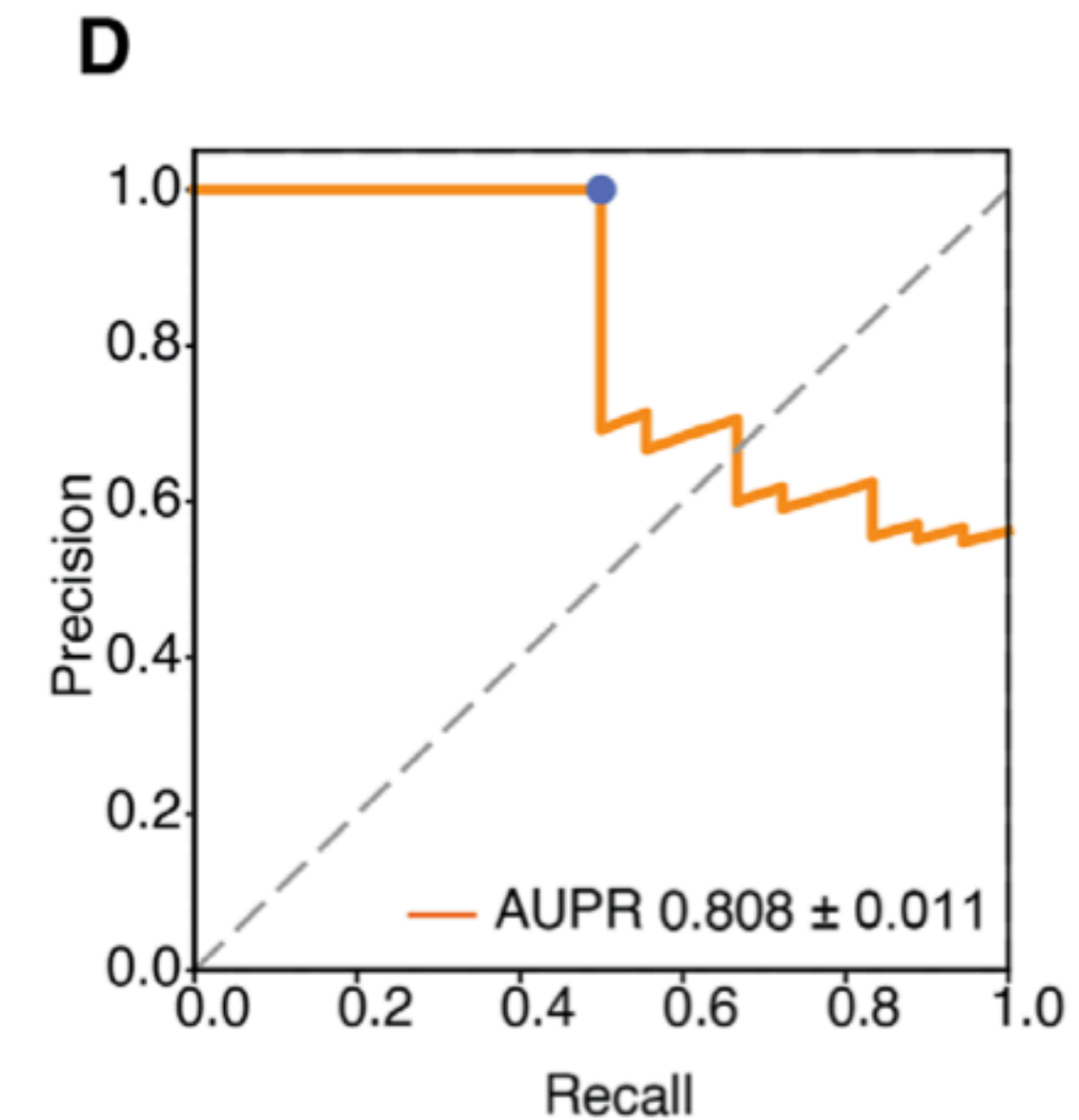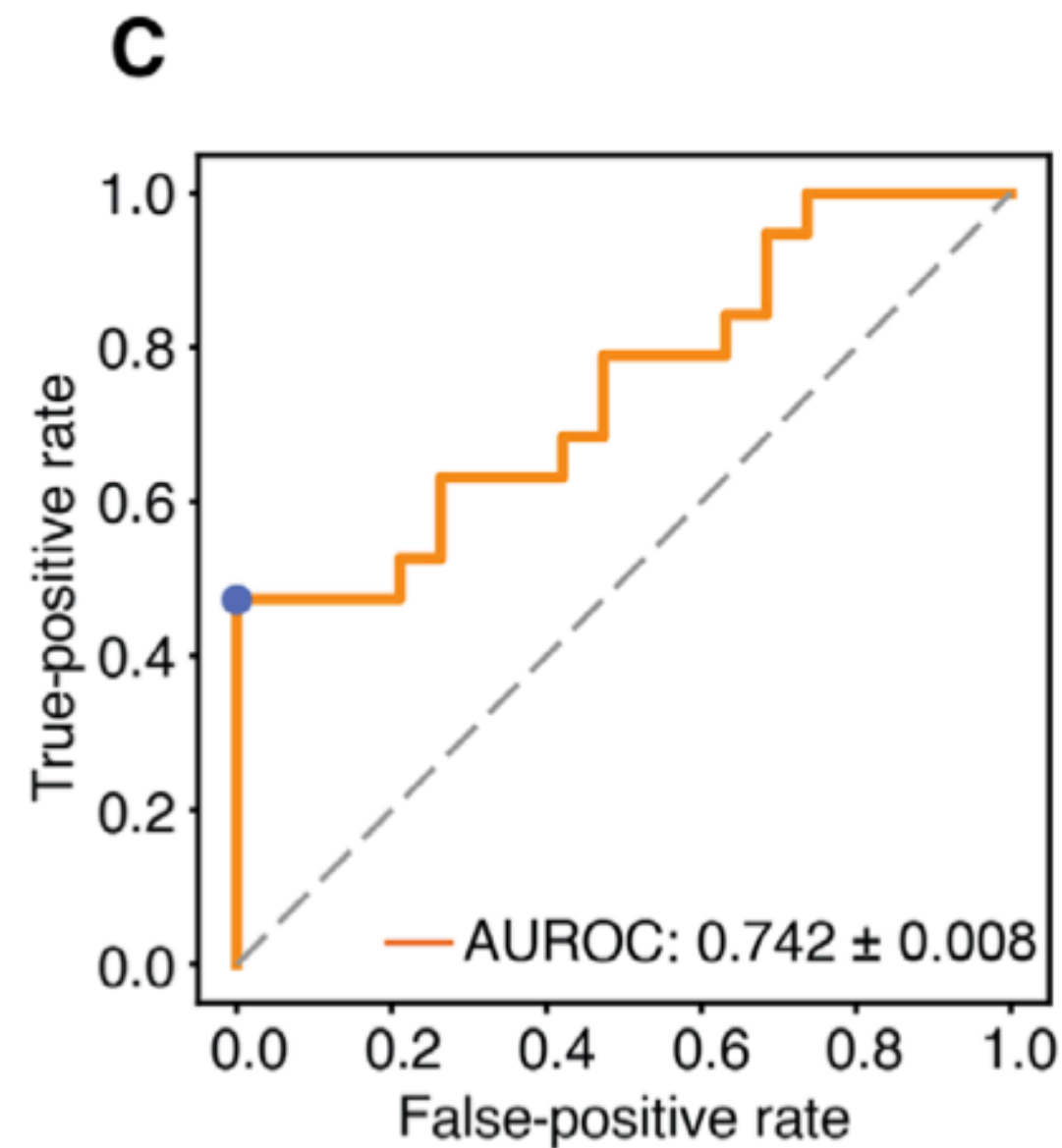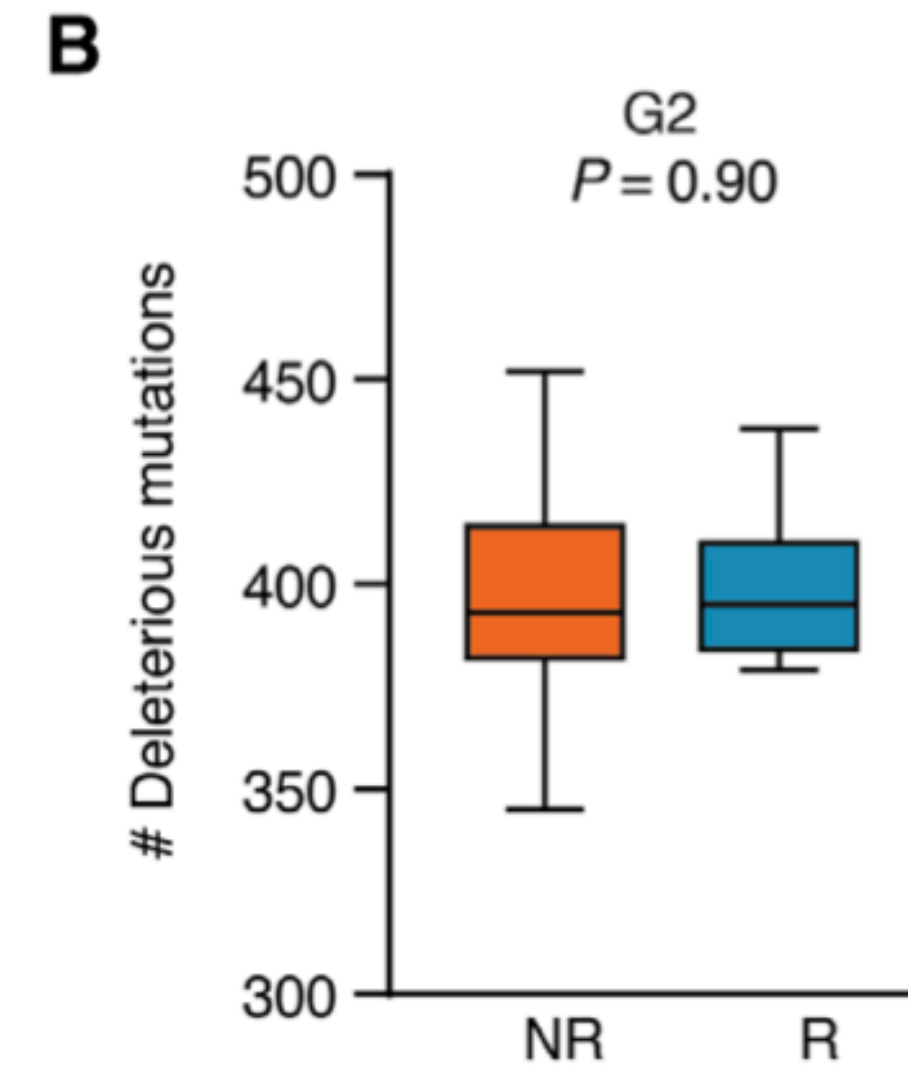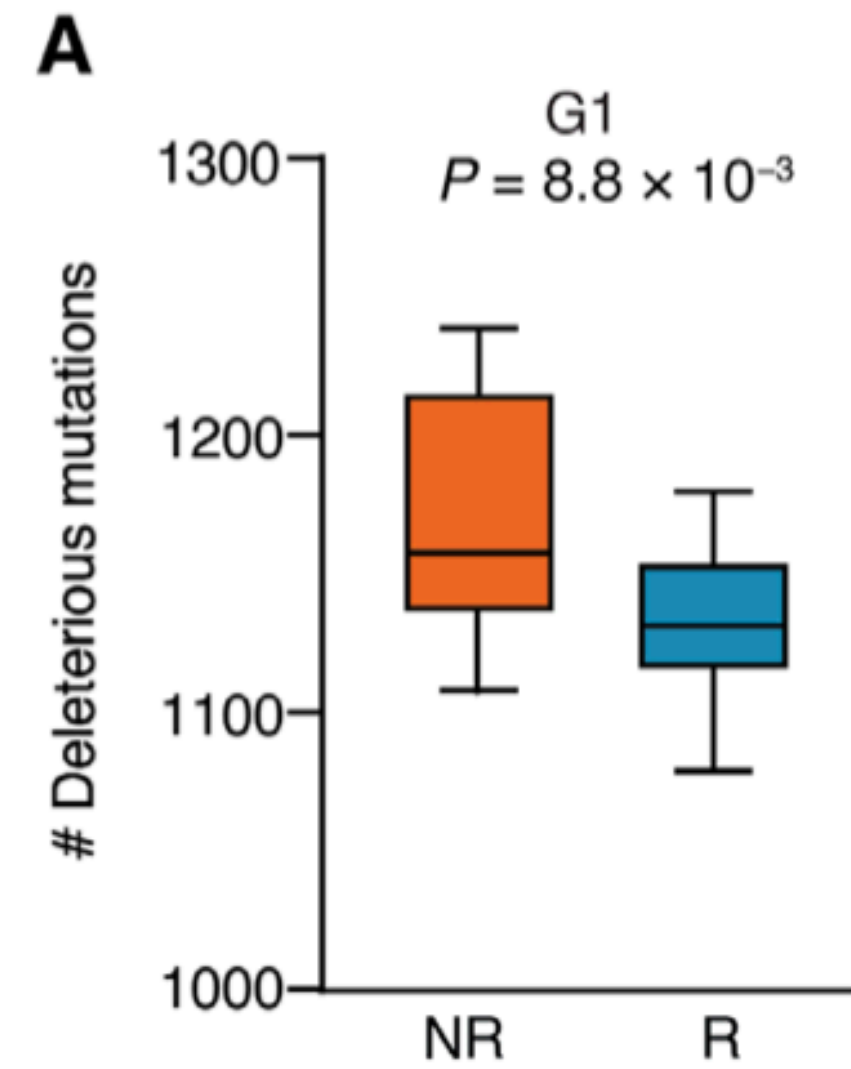It uses a DAG, that's how you known it's Bayesian

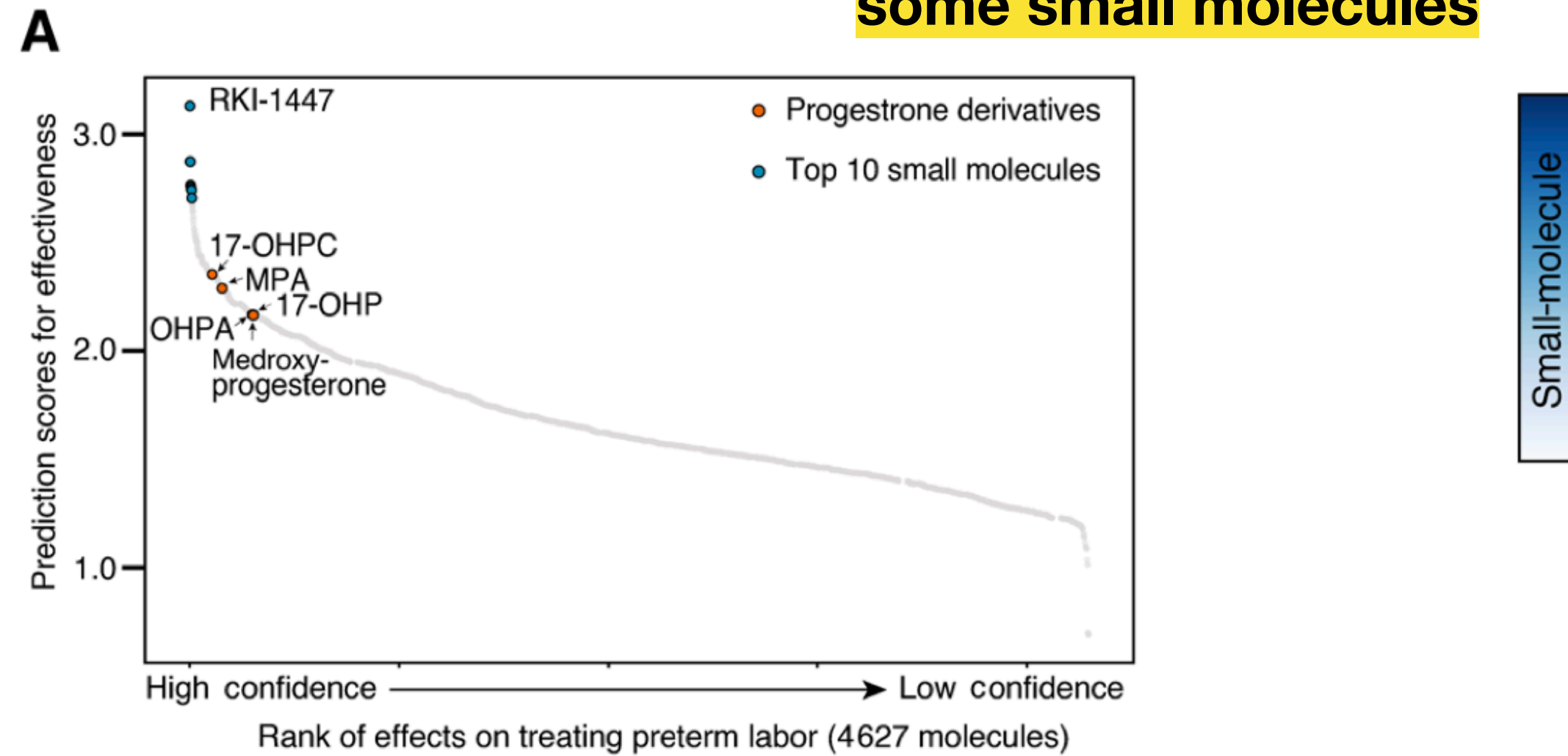**BEAR loci have higher DEEP+ scores and strong GWAS risk associations**

**… and high haploinsufficiency scores and higher expression in myometrial tissue**

Chromatin accessibility mutation burden is predictive of progestin response

**Drug discovery rediscovered progesterone derivatives and also some small molecules**
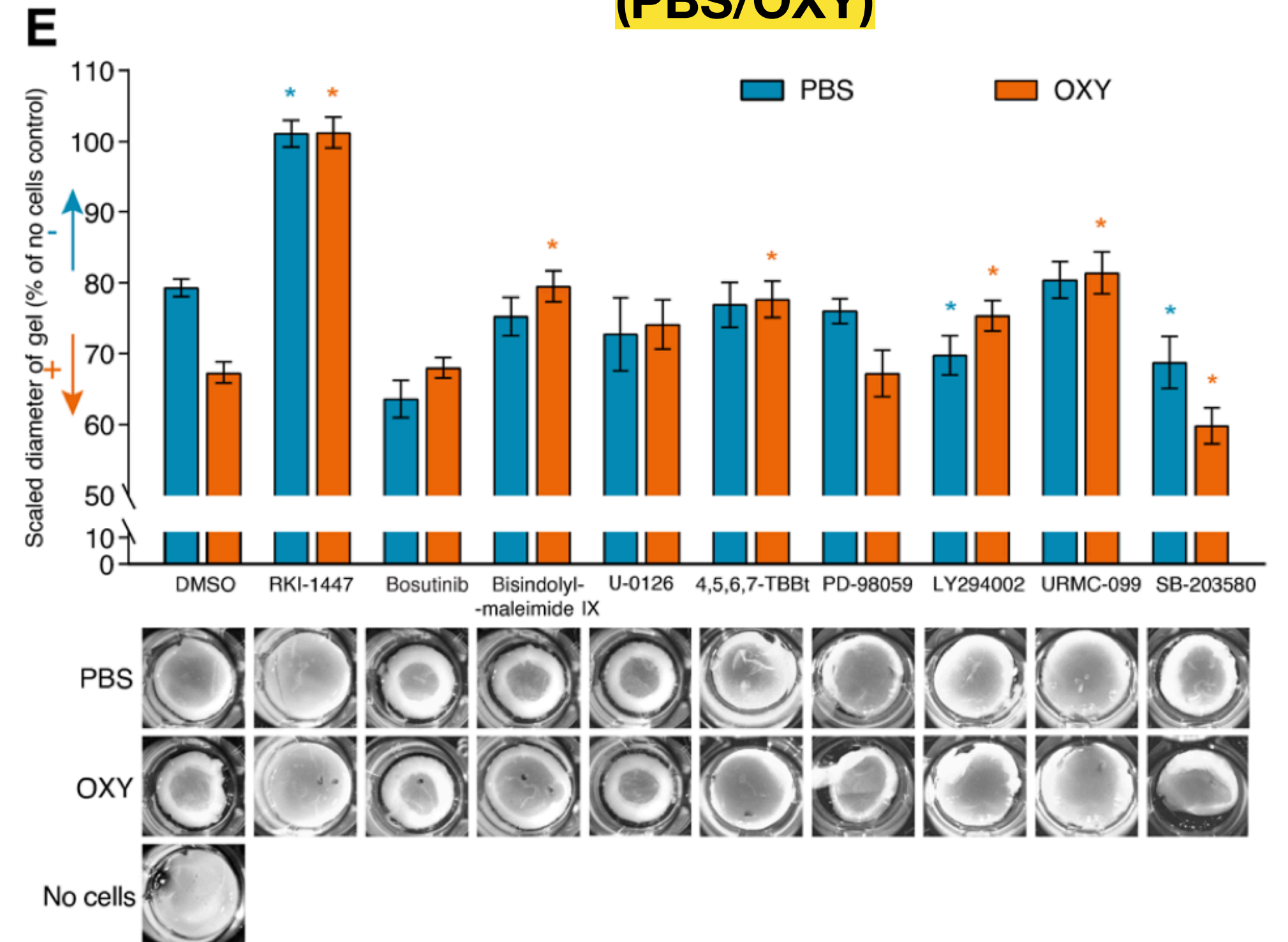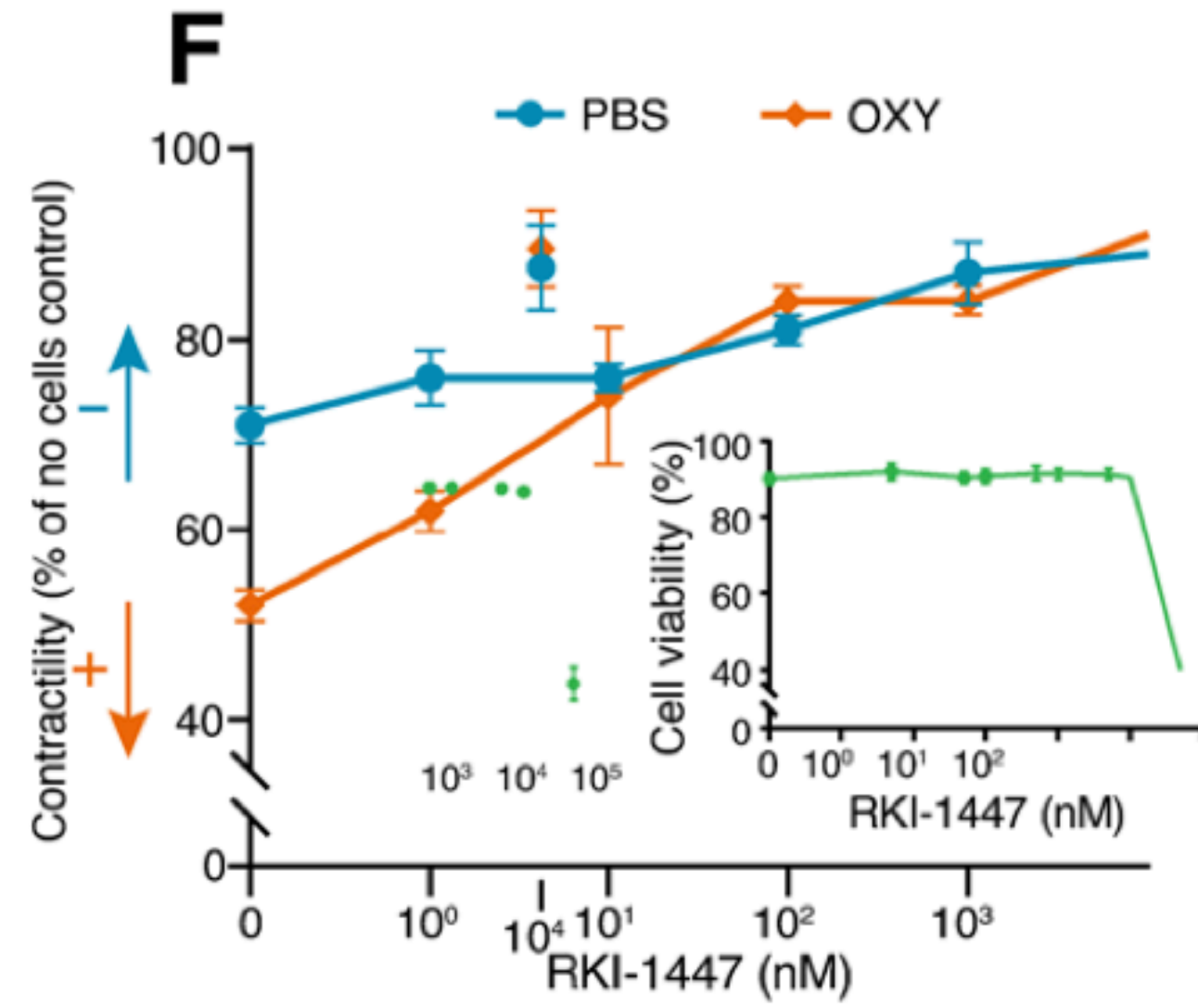
**Small molecule candidates reduced contracts in both physiological states (PBS/OXY)**



**A**

Prediction scores for effectiveness

- RKI-1447
- Progestrone derivatives
- Top 10 small molecules

17-OHPC
MPA
17-OHP
OHPA
Medroxy-progesterone

High confidence → Low confidence

Rank of effects on treating preterm labor (4627 molecules)

Small-molecule

**B**

| Drug name | Score | Percentile |
|---|---|---|
| Aspirin | 2.529 | 0.61% |
| 17-OHPC | 2.332 | 3.00% |
| 17-OHP | 2.11 | 6.92% |
| Nifedipine | 1.814 | 27.93% |
| Terbutaline | 1.689 | 37.71% |
| Ibuprofen | 1.45 | 69.28% |
| Naproxen | 1.388 | 79.43% |
| Atosiban | 1.129 | 99.58% |

**D**
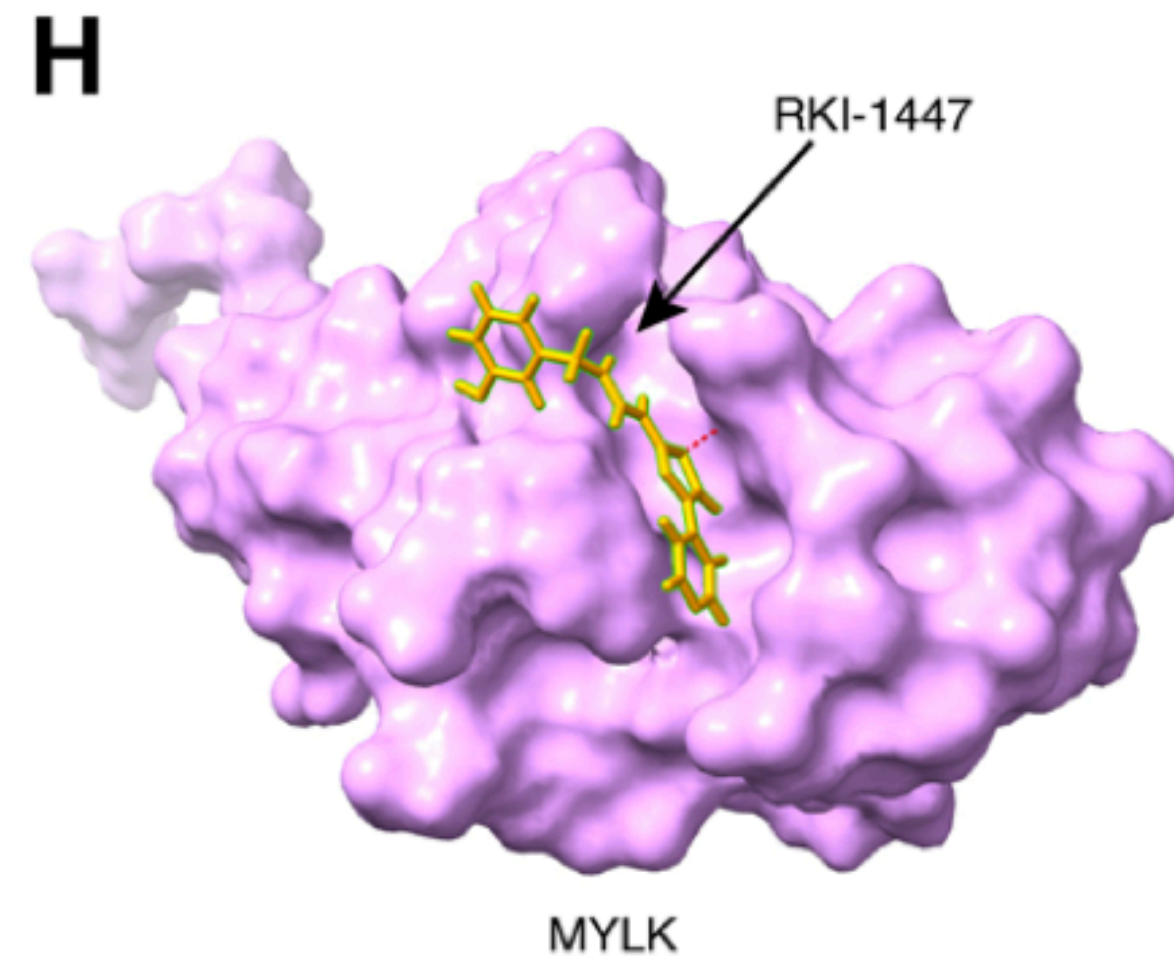
| Drug name | Function |
|---|---|
| RKI-1447 | ROCK kinase inhibitor |
| Ephedrine-HCL | Adrenoceptor agonist |
| Bosutinib | Bcr-Abl kinase inhibitor |
| Bisindolylmaleimide IX | Pan-PKC inhibitor |
| U-0126 | MEK inhibitor |
| 4,5,6,7-TBBt | CK2 inhibitor |
| PD-98059 | MEK inhibitor |
| LY294002 | PI3K inhibitor |
| URMC-099 | MLK inhibitor |
| SB-203580 | p38 MAPK inhibitor |

**E**

Scaled diameter of gel (% of no cells control)

PBS / OXY

DMSO, RKI-1447, Bosutinib, Bisindolyl-maleimide IX, U-0126, 4,5,6,7-TBBt, PD-98059, LY294002, URMC-099, SB-203580

PBS / OXY / No cells

**RKI-1447 dose response curve (reduction of contractions is increase Y)**

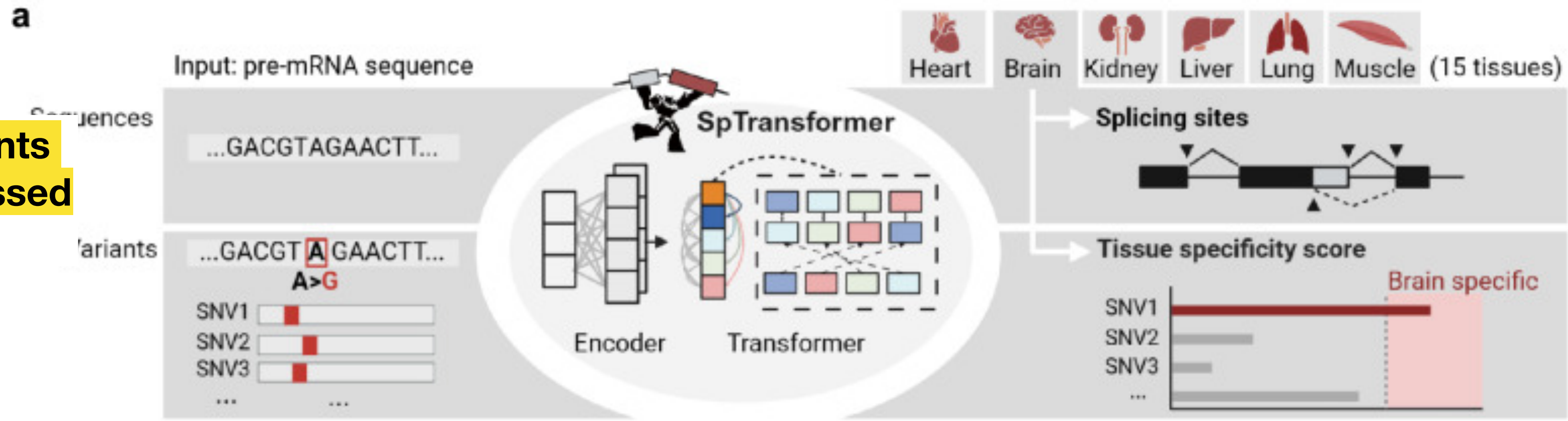**Fits in the pocket, so nice!**

# "Cruel Summer"

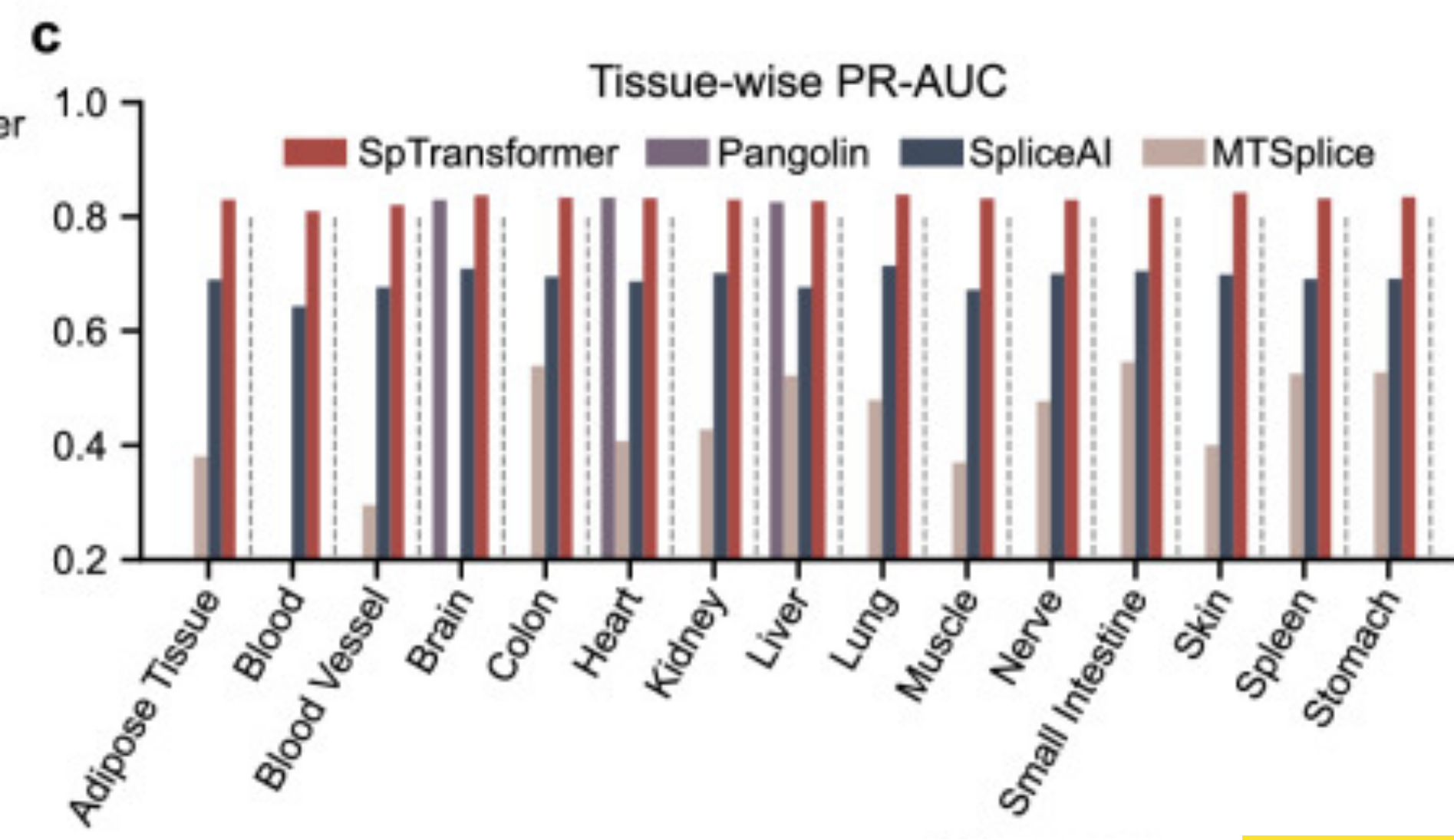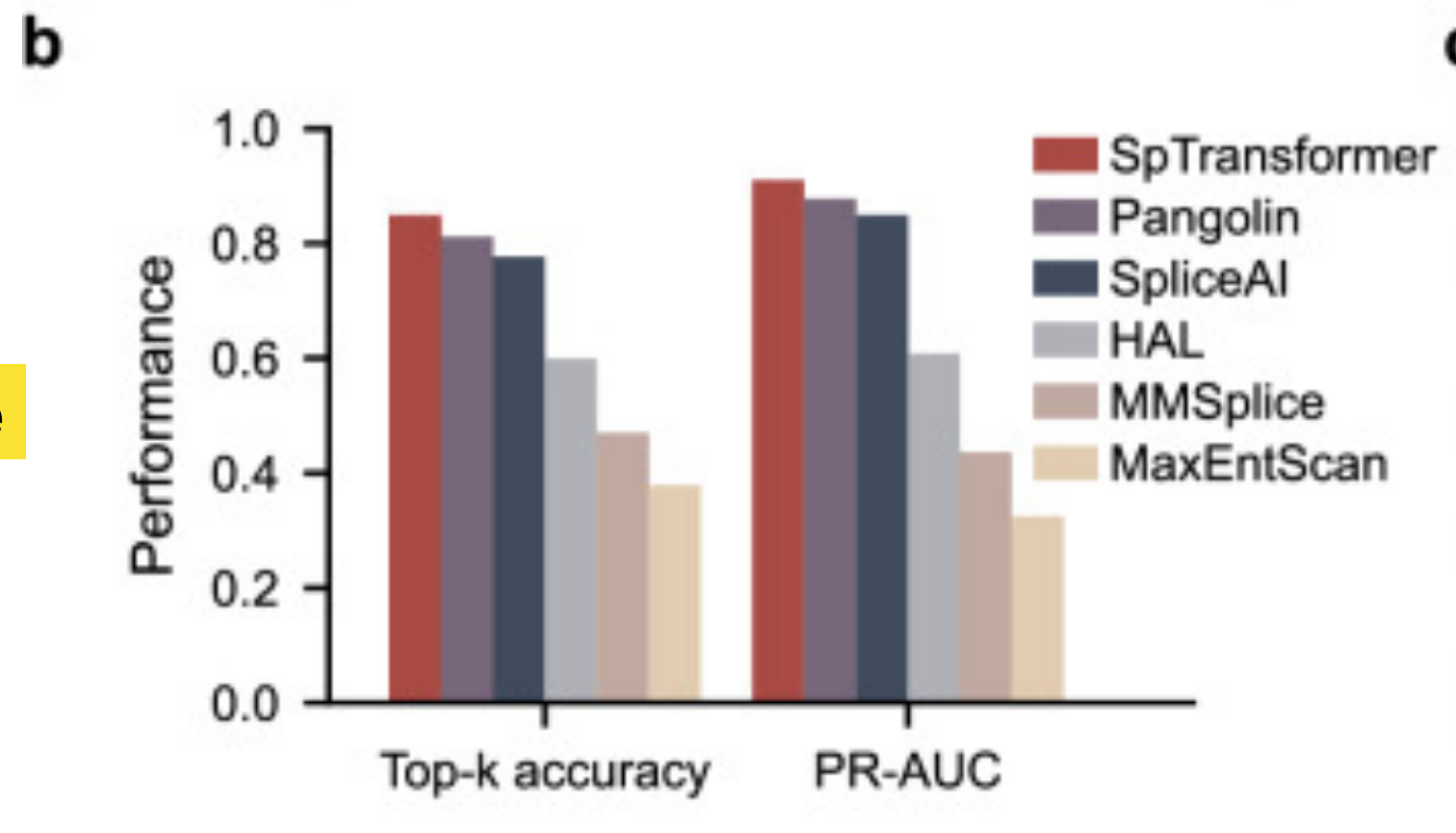## Taylor-ed for You—Precision Medicine in Action

# SpliceTransformer predicts tissue-specific splicing linked to human diseases (You et al, *Nature Communications*)

- Goal: To predict tissue-specific RNA splicing variation

- Method:

  - Use a Sinkhorn Transformer — allowing for increased sequence context by using a structured attention mechanism

  - Integrates both GTEx (for fine-grained tissue specificity) and cross-species RNASeq data (to improve generalization)

  - Can generate a ΔSplice score for each variant to quantify its effect on splicing

- Result:

  - Found that 60% of intronic and synonymous mutations have high ΔSplice scores

  - Outperforms existing methods on splice site prediction and even more so for tissue specificity

  - Found tissue specific splicing alterations are enriched for diseases of those tissues

- Conclusion: Adding sequence and tissue context greatly improves understanding of variation
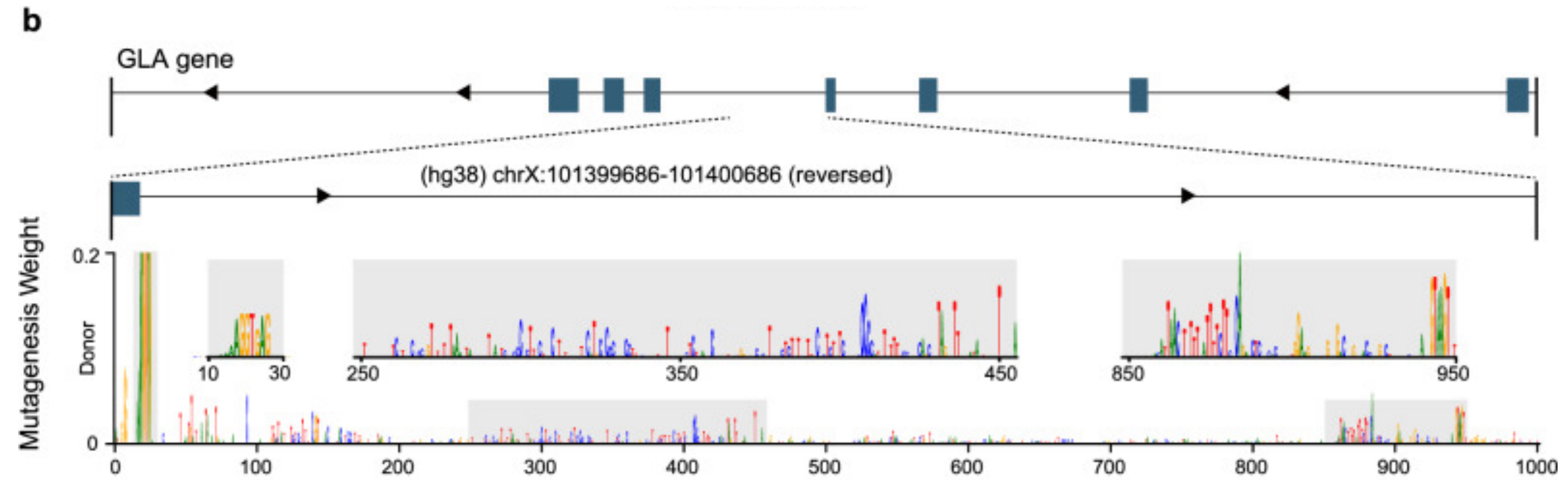
Sequence and variants are encoding and passed to transformer
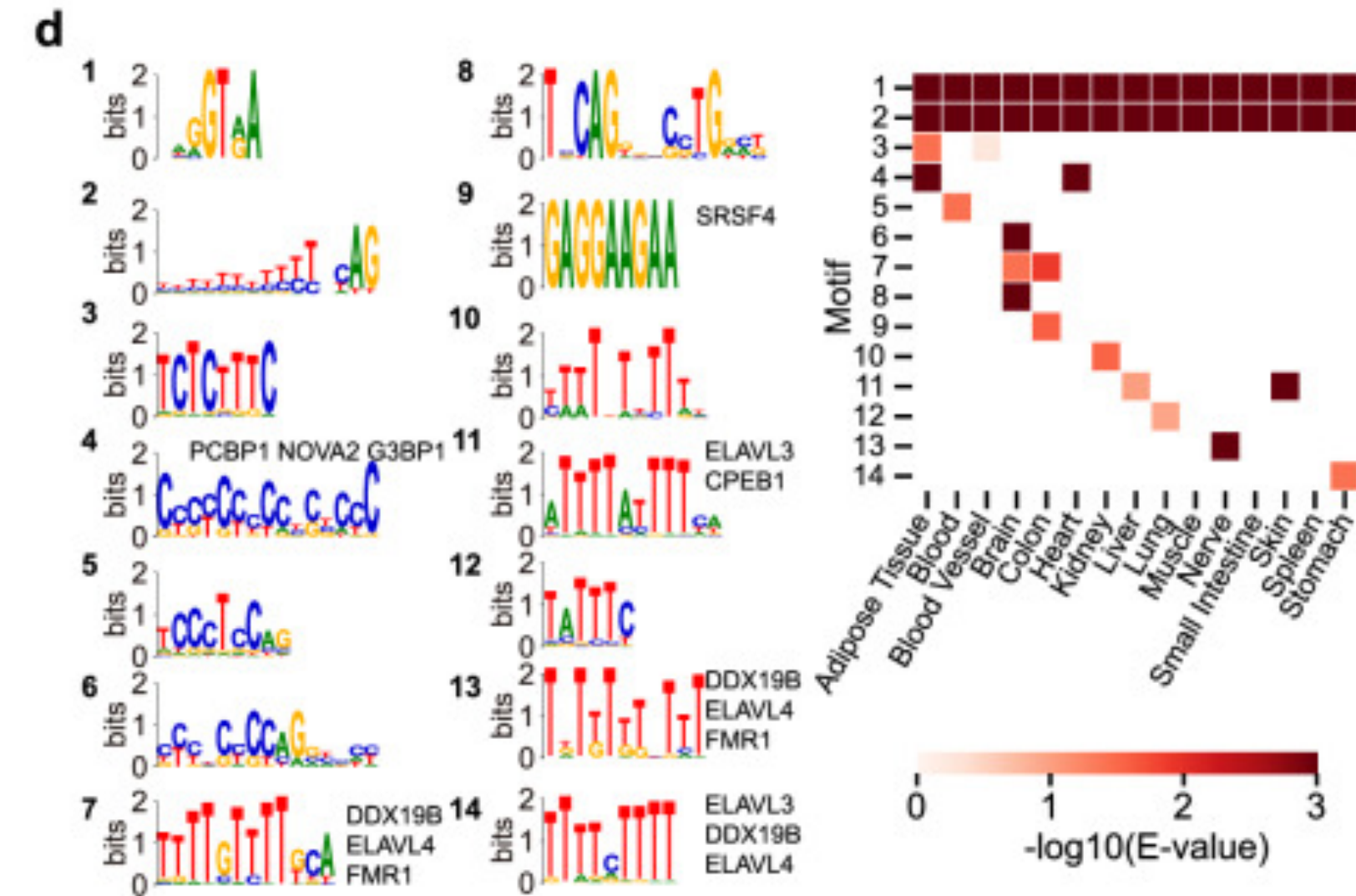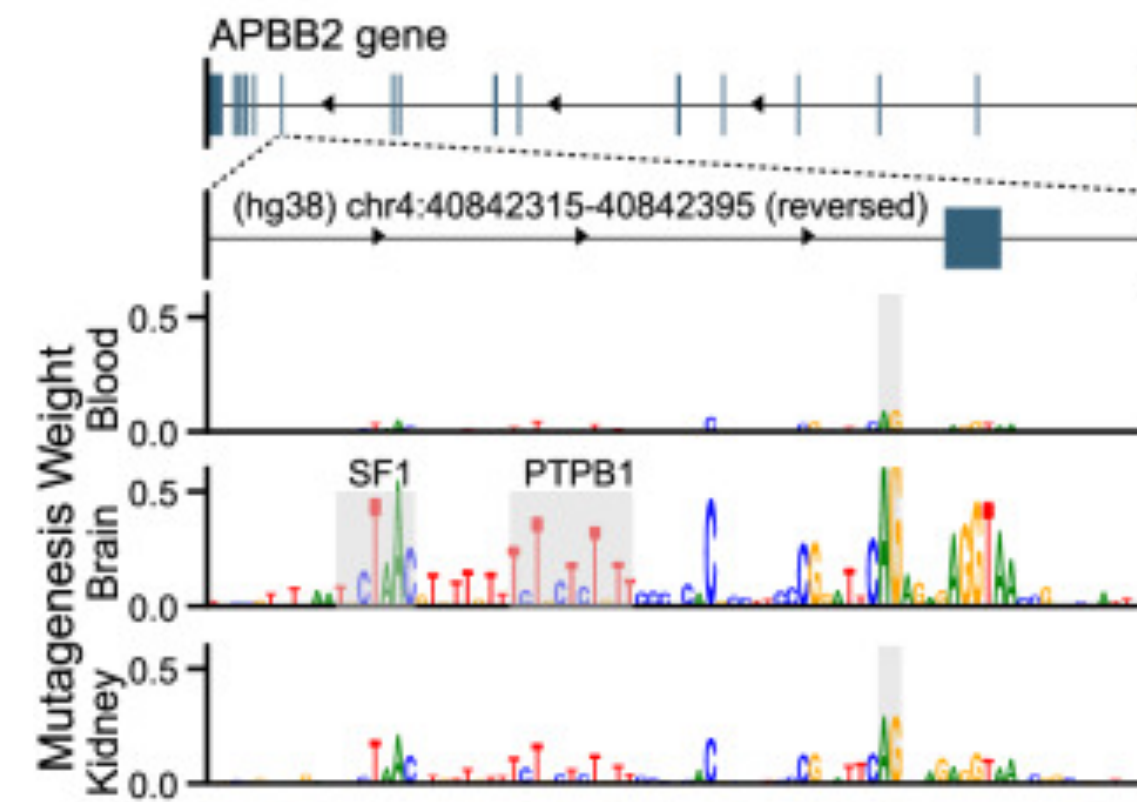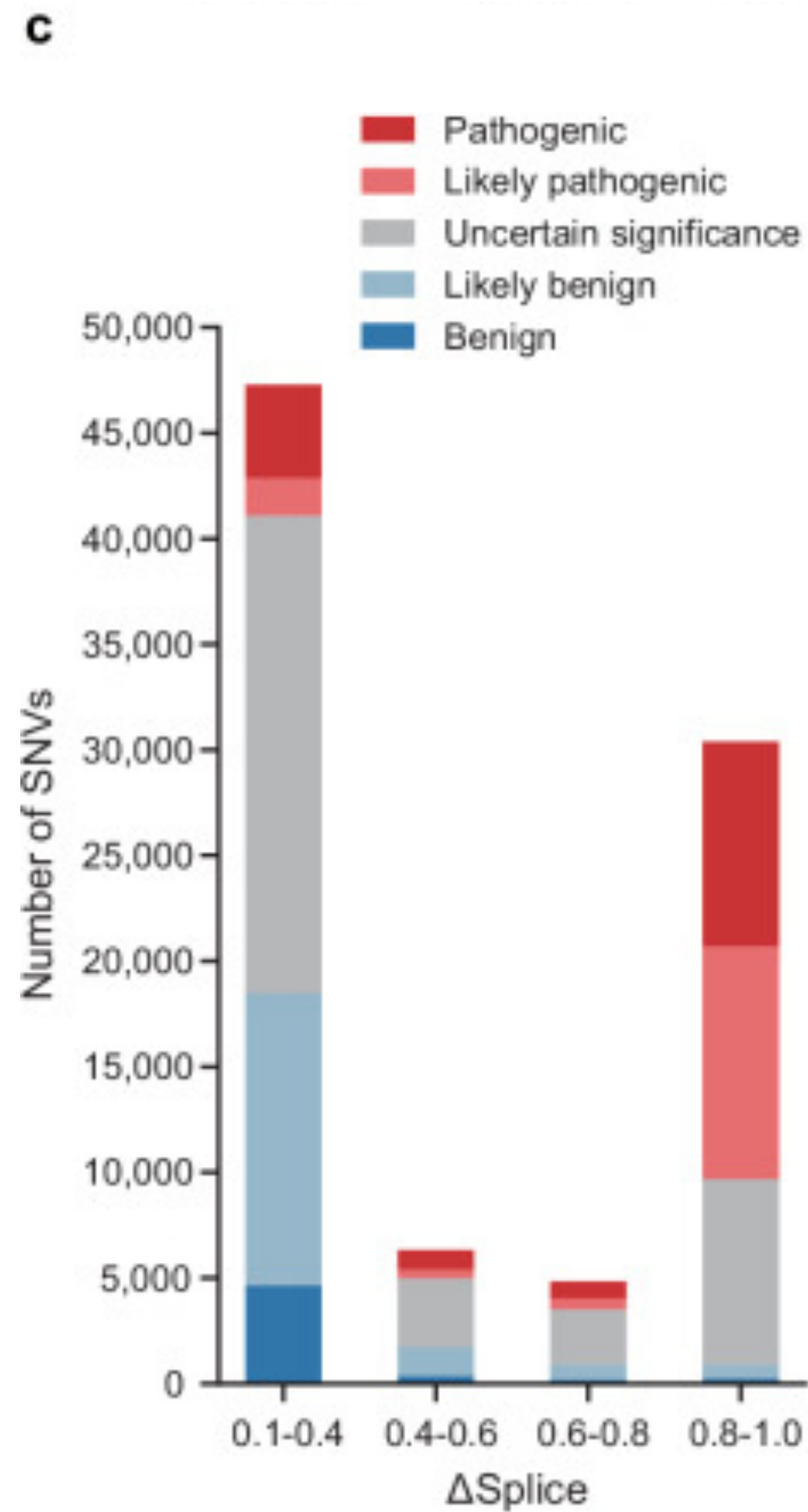
Better performance overall

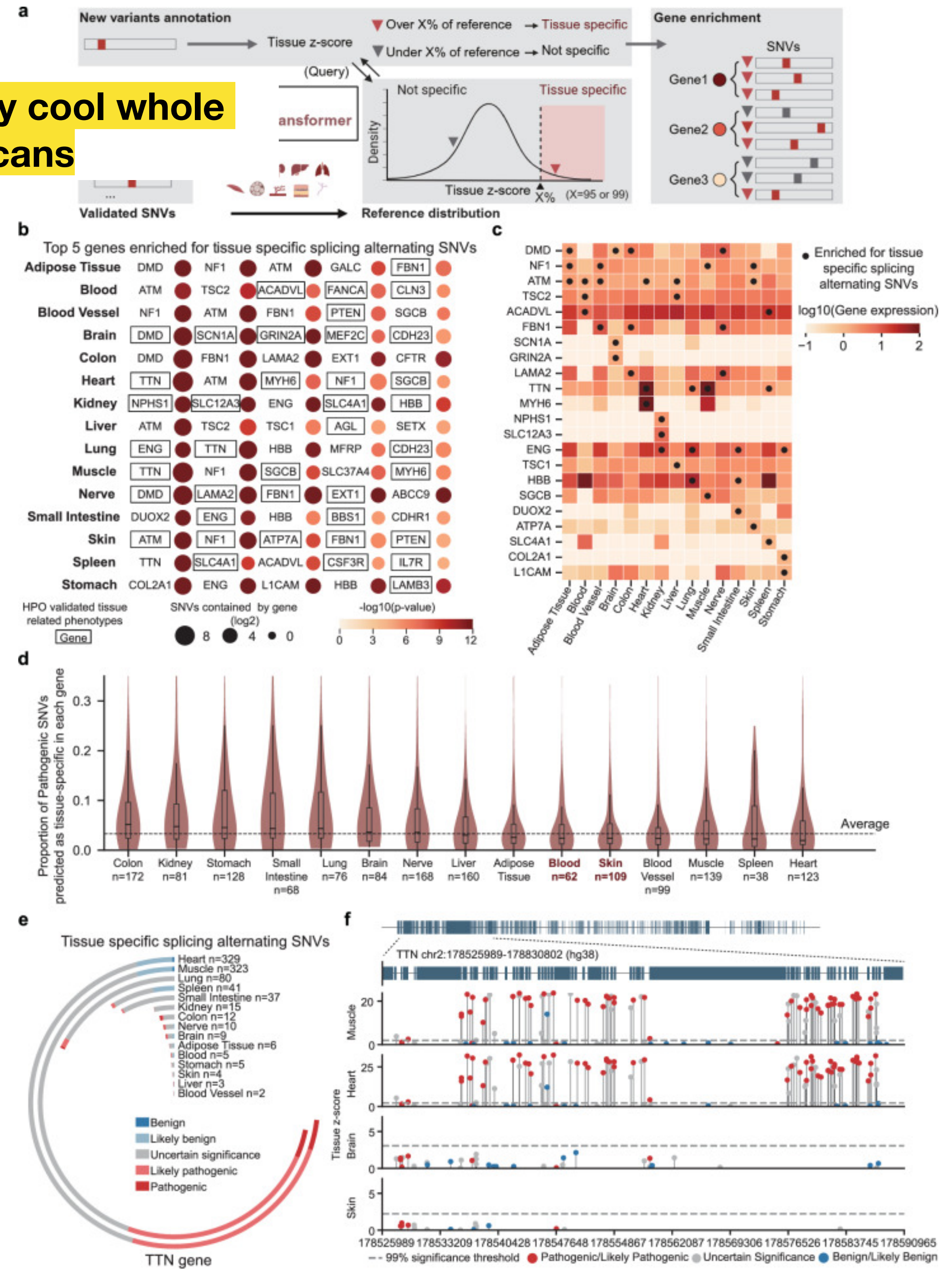... and especially for tissue specificity

SpTransformer can score the regulatory regions for their affect on splicing

Can do some pretty cool whole genome scans

Large ΔSplice scores are more likely to be pathogenic

# Integrating imaging and genomic data for the discovery of distinct glioblastoma subtypes: a joint learning approach (Guo et al, *Scientific Reports*)
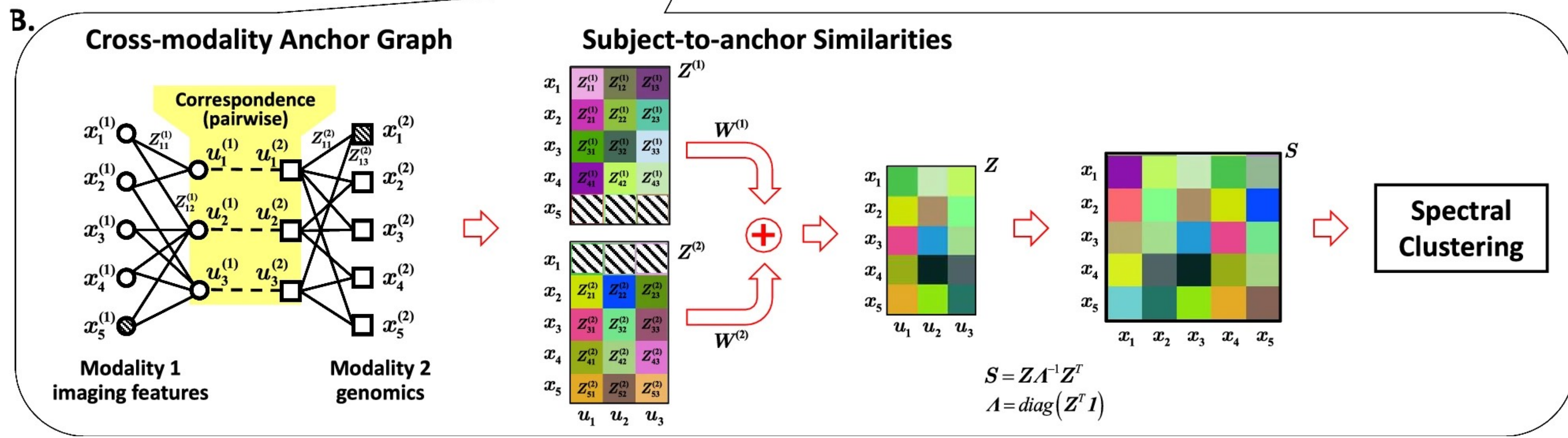
- Goal: Improve power of glioblastoma studies by integrating both imaging and genomic data

- Method:

  - Joint learning model that integrates MRIs and mutations in 27 key genes

  - 571 glioblastoma patients

- Result:

  - Found three subtypes of glioblastoma with significant survival differences

  - Found image-genomic correlation between RB1 and PTEN pathways linked to diffusion metrics

- Conclusion: Multimodal data continues to empower scientific investigation

**Found significant survival differences in the clusters**

**Found significant survival differences in the clusters**

**Clusters have distinct imaging morphologies**

A.

Survival Probability

Discovery Cohort
HR (Subtype 1 & Subtype 2): 1.31 (0.94 - 1.84)
HR (Subtype 1 & Suptype 3): 1.64 (1.17 - 2.31)
HR (Subtype 2 & Subtype 3): 1.34 (0.95 - 1.90)
95% confidence interval

Groups
Subtype 1
Subtype 2
Subtype 3

p = 0.00733

Subtype 1
Subtype 2
Subtype 3

B.

Survival Probability

Subtype 1
Subtype 2
Subtype 3

Can identify interaction between imaging and genetics



Imaging features are things like: Perfusion-related features, contrast enhancement patterns, diffusion measurements, etc.

# Discovering the gene-brain-behavior link in autism via generative machine learning (Kundu et al, *Science Advances*)

- Goal: Identify brain structural changes linked to 16p11.2 CNVs (deletions & duplications) — most important genetic risk factor for Autism

- Method:

  - 206 individuals with MRIs, genetics, and behavioral assessments

  - Generative model (called 3D Transport-Based Morphometry) that quantifies the distance between patient and the reference (average of all)

  - Subsequent supervised learning using penalized latent discriminant analysis to predict variant status

- Result:

  - High accuracy at CNV prediction

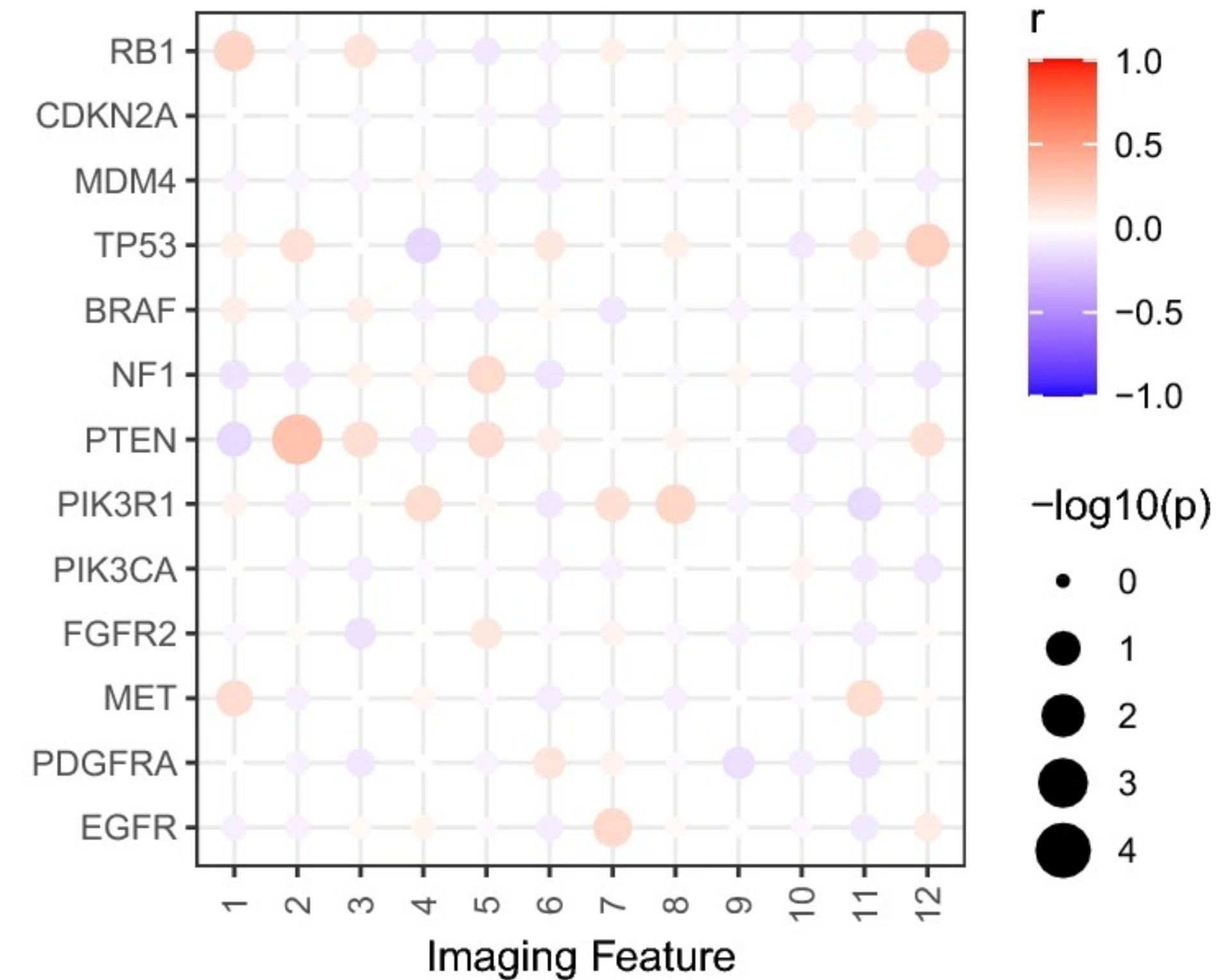  - Deletion carriers -> brain overgrowth

  - Duplication carriers -> brain undergrowth

  - Strong explanation of behavioral phenotypes

- Conclusion: Good demonstration that you can have too much information. By converting raw MRIs into these morphometry distances they were able to get better signal. There's a lesson here for all of us.

Converting to transport space highlights the differences

**TBM shows what happens to the brain under different variant conditions**

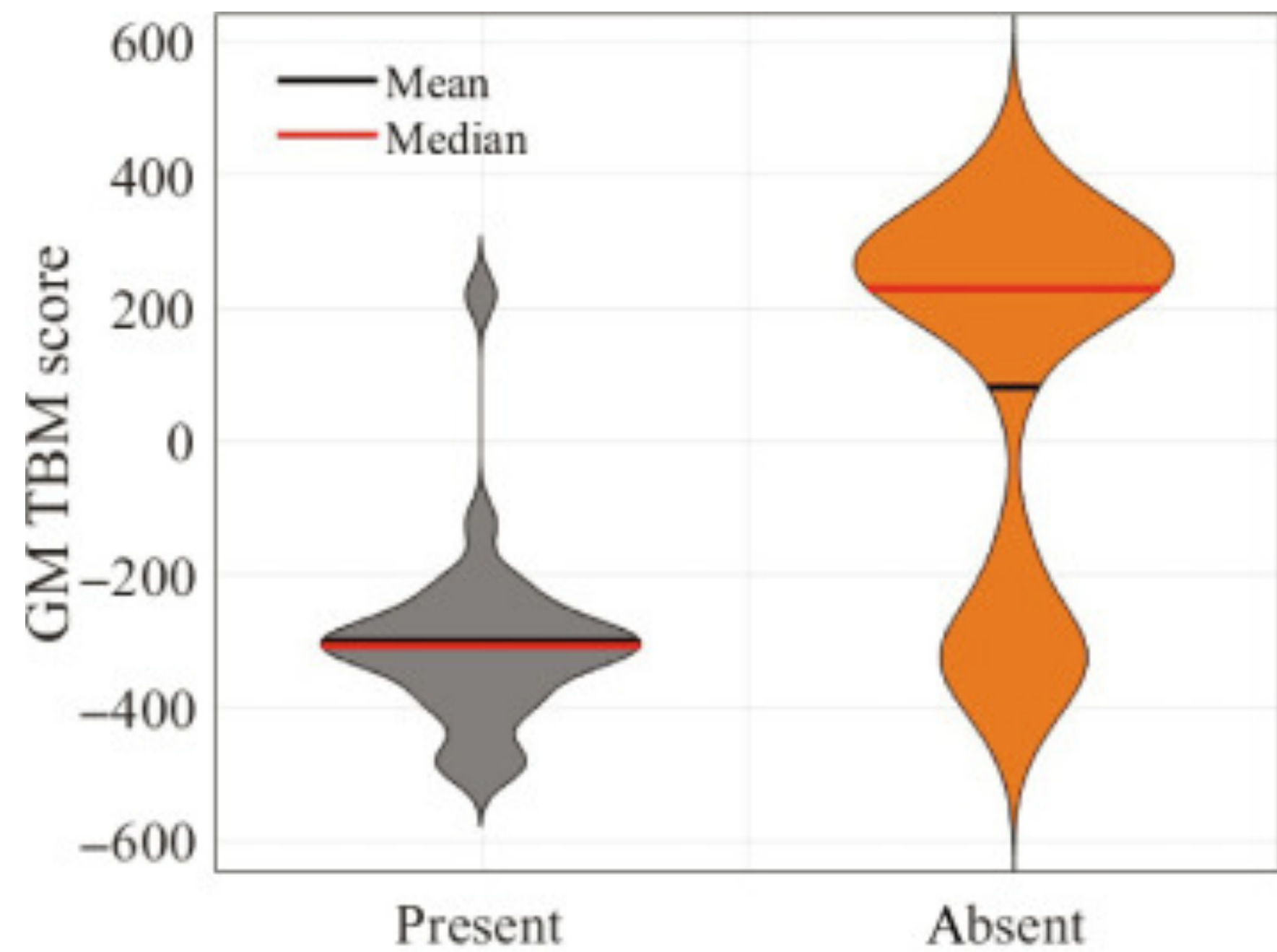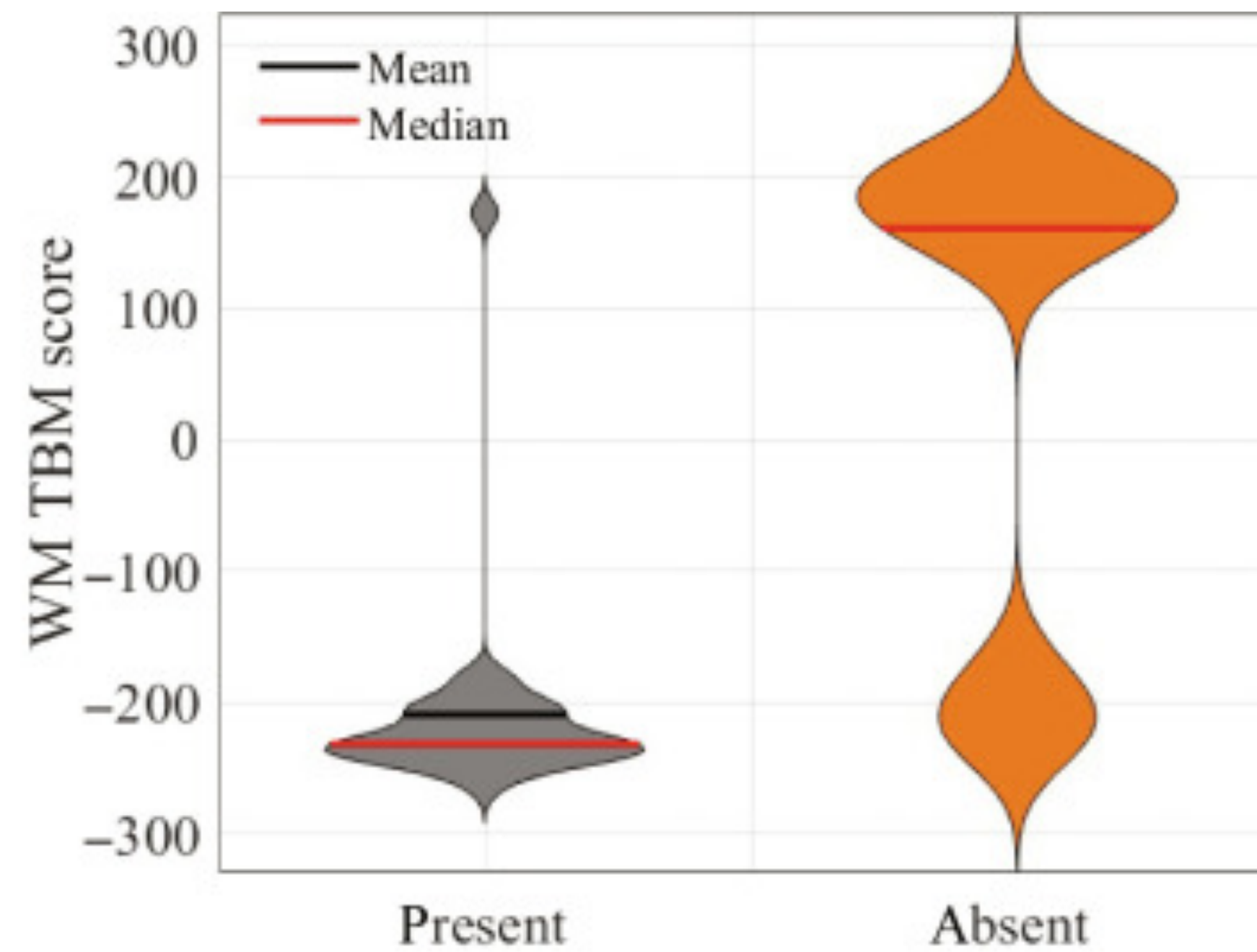**TBM scores are good discriminators of the presence of an articulation disorder**



**A** gray matter

**B** white matter

# Trainee Spotlight



**VEDA PRIYA
PULIGUNDLA**
UNIVERSITY OF
MINNESOTA

# Deep learning predicts DNA methylation regulatory variants in specific brain cell types and enhances fine mapping for brain disorders ( Zhou et al., 2025 )  https://www.science.org/doi/10.1126/sciadv.adn1870

**Goal:**

- Identify cell type-specific DNAm variants  & Improve fine mapping of brain disorder risk loci

**Method:**

- Developed **INTERACT model** (CNN + Transformer) : A Deep Learning Model
- Trained on **single-nucleus DNAm data**
- **In silico mutagenesis** for variant effects
- **LD-score regression** for heritability analysis

**Results:**

- **High accuracy** (AUC = 0.99) in DNAm prediction
- Identified **regulatory variants for schizophrenia, depression, and Alzheimer's Disease (AD)**
- **Fine mapping** reveals cell type-specific causal loci
- **rs74504435 linked to EGFR in astrocytes** (potential AD target)
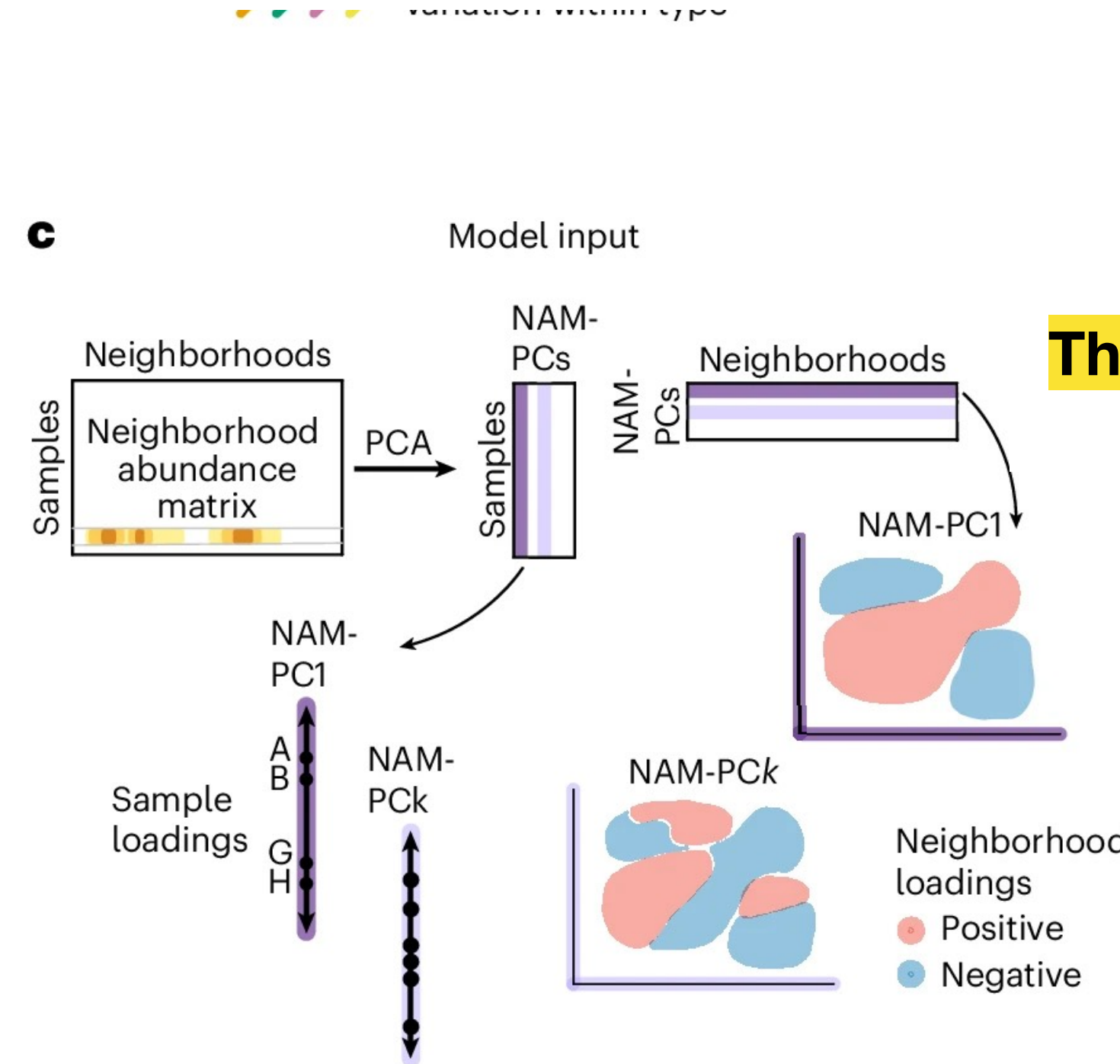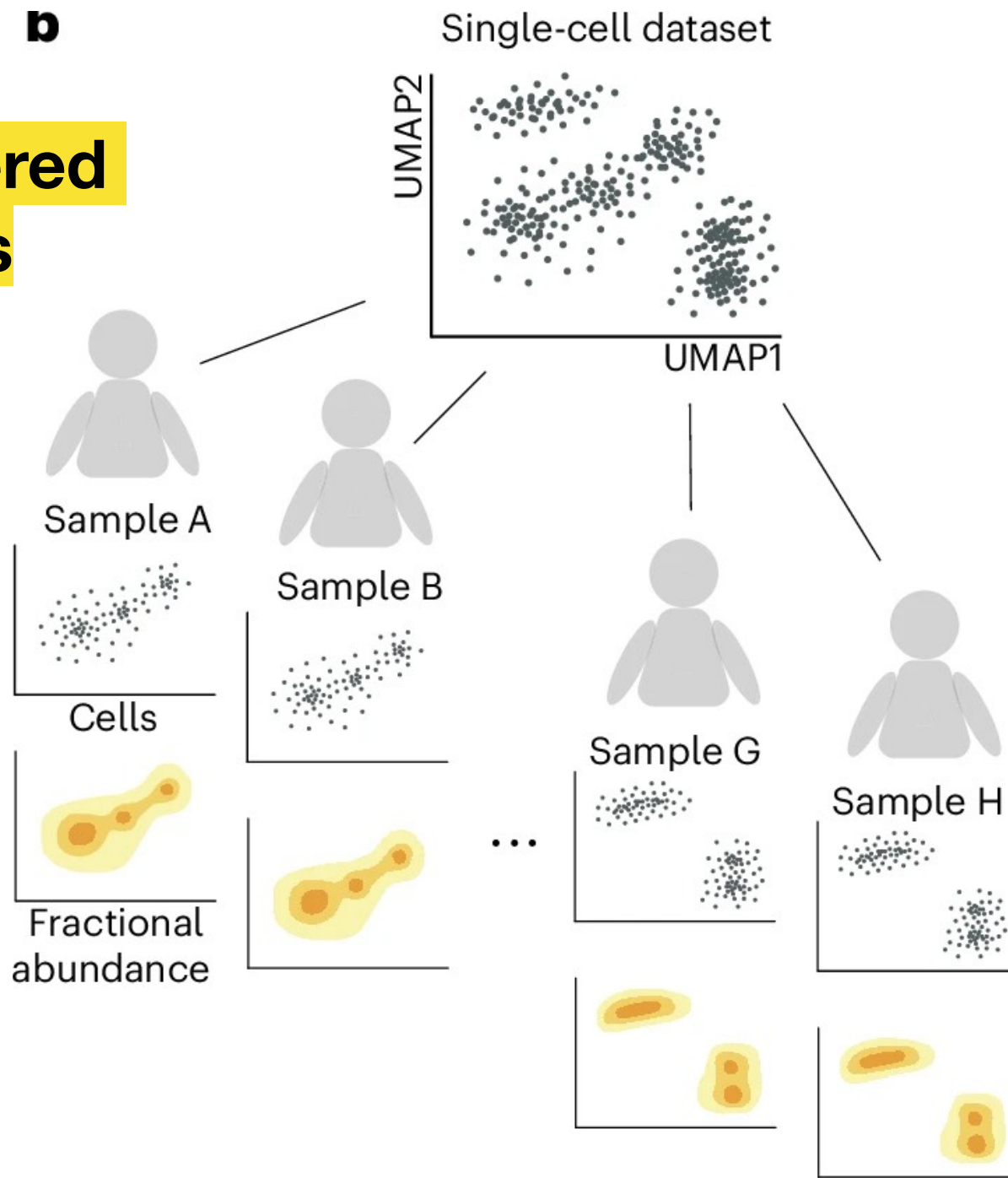
**Conclusion:**

- Deep learning improves genetic risk mapping

# Identifying genetic variants that influence the abundance of cell states in single-cell data (Rumker et al, *Nature Genetics*)

- Goal: Identify genetic variants that control <u>cell state</u> abundance ("csaQTLs")

- Method:

  - Introduce GeNA a statistical method to identify csaQTLs

    - Does not require states to be pre-defined

    - Uses Neighborhood Abundance Matrix (NAM) and PCA

  - Apply to dataset of ~1k individuals (800k peripheral blood mononuclear cells)

- Result:

  - Five significant csaQTLs identified, primarily affecting natural killer (NK) and myeloid cell states

  - rs3003-T was linked to increased NK cells expressing tumor necrosis factor

- Conclusion: I'll admit — not a "QTL" I had thought of before

**First cells are clustered in neighborhoods**

**Then do PCA to reduce dimension**

**Use a chi-squared test to look for any deviations that correlate with allelic dose**

b Single-cell dataset
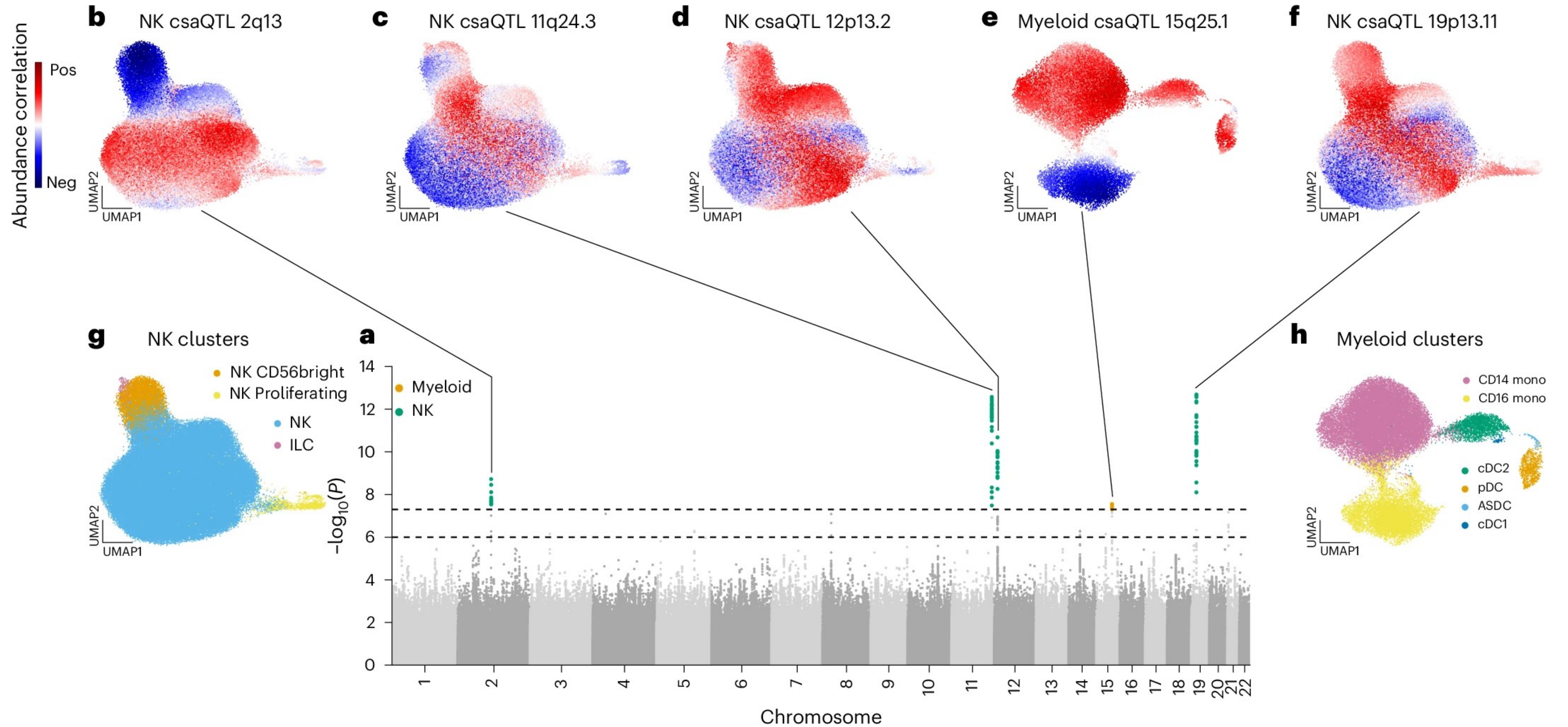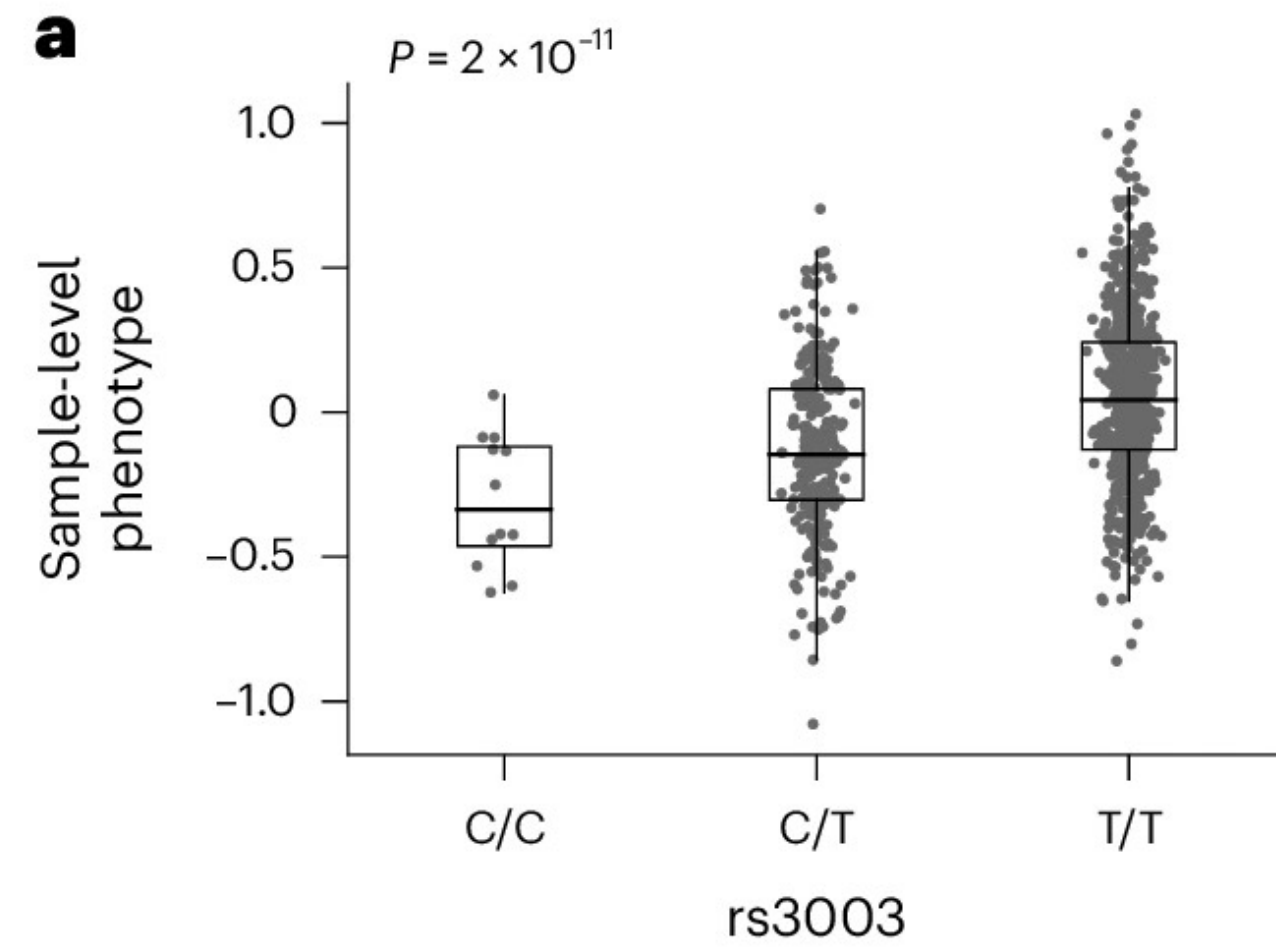
Sample A

Sample B

Sample G

Sample H

Cells

Fractional abundance

c Model input

Neighborhoods

Neighborhood abundance matrix

PCA

NAM-PCs

Neighborhoods

NAM-PC1

NAM-PCk

Neighborhood loadings
- Positive
- Negative

Sample loadings

NAM-PC1

NAM-PCk

d Association test per allele

$X_1 = \left(\dfrac{\widehat{\beta_1}}{\text{s.e.}}\right)^2 \sim \chi_1^2$

$Y = \sum_{i=1}^{k} X_i \sim \chi_k^2$

$X_k = \left(\dfrac{\widehat{\beta_k}}{\text{s.e.}}\right)^2 \sim \chi_1^2$

Test desired set of alleles

(for example, genome-wide survey)

NAM-PC1-k    Allele dose

e Output: csaQTL

Defines and identifies genotype-associated abundance shifts

Phenotype value per individual

Abundance shift
- Expanded
- Depleted

Cell states

**Discovered csaQTLs**

**Correlation between allelic dose and abundance of given cell state**

**Deeper dive on rs3003**

GeNA output



**a** $P = 2 \times 10^{-11}$

Sample-level phenotype (y-axis): 1.0, 0.5, 0, −0.5, −1.0

rs3003 (x-axis): C/C, C/T, T/T

Detected association between genotype and cell-state abundance shift phenotype

**b** Neighborhood-level phenotype

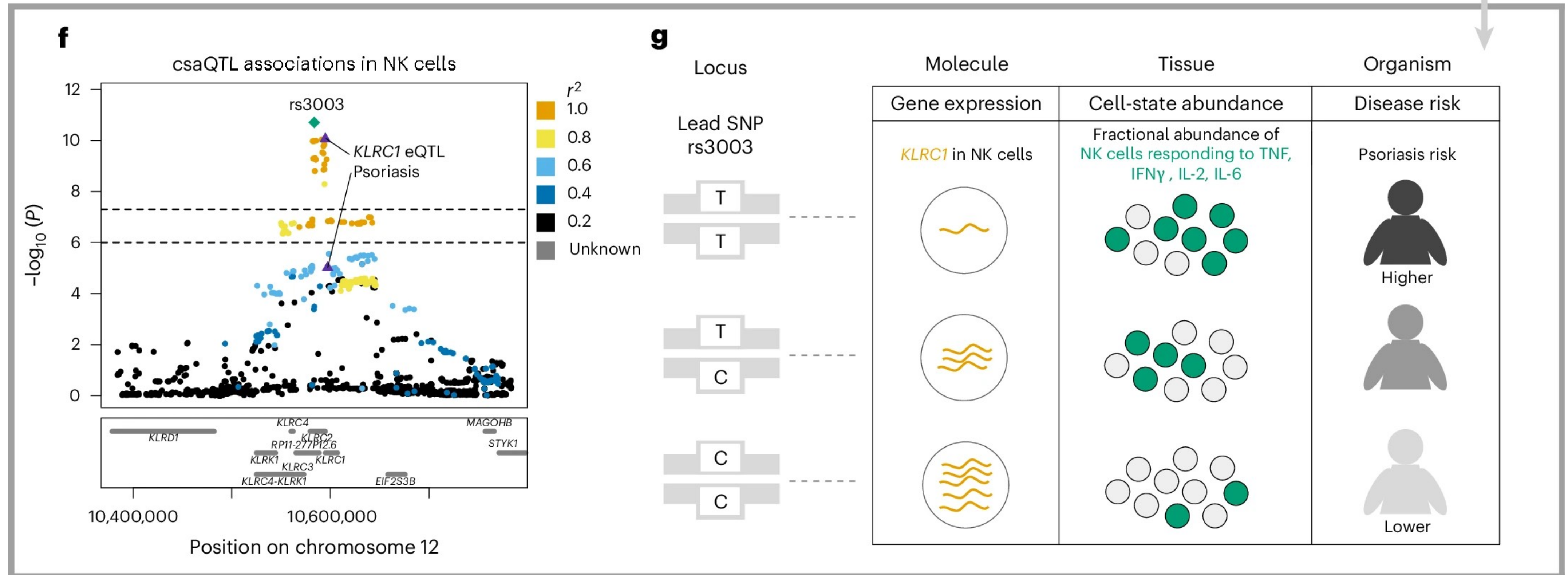Abundance correlation: Positive / Negative

UMAP2 / UMAP1

Definition of cell-state abundance shift associated with genotype, at neighborhood resolution

**Does/Abundance relationships**

**Deeper dive on rs3003**

**Potential psoriasis mechanism**



Colocalization of molecular, tissue and organism-level traits

**f** csaQTL associations in NK cells

**g** Locus / Molecule / Tissue / Organism

Gene expression — Cell-state abundance — Disease risk

Lead SNP rs3003

KLRC1 in NK cells — Fractional abundance of NK cells responding to TNF, IFNγ, IL-2, IL-6 — Psoriasis risk
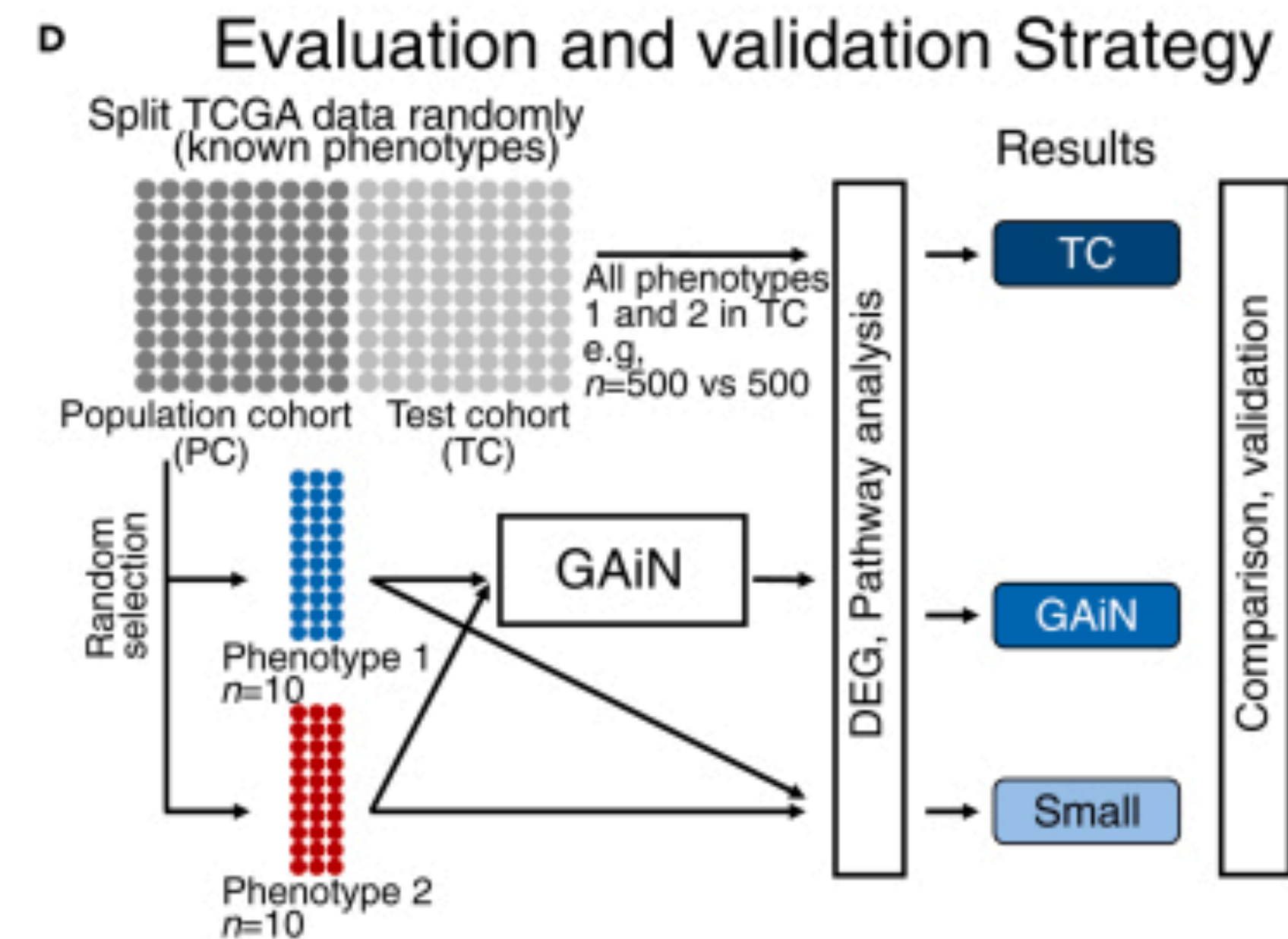
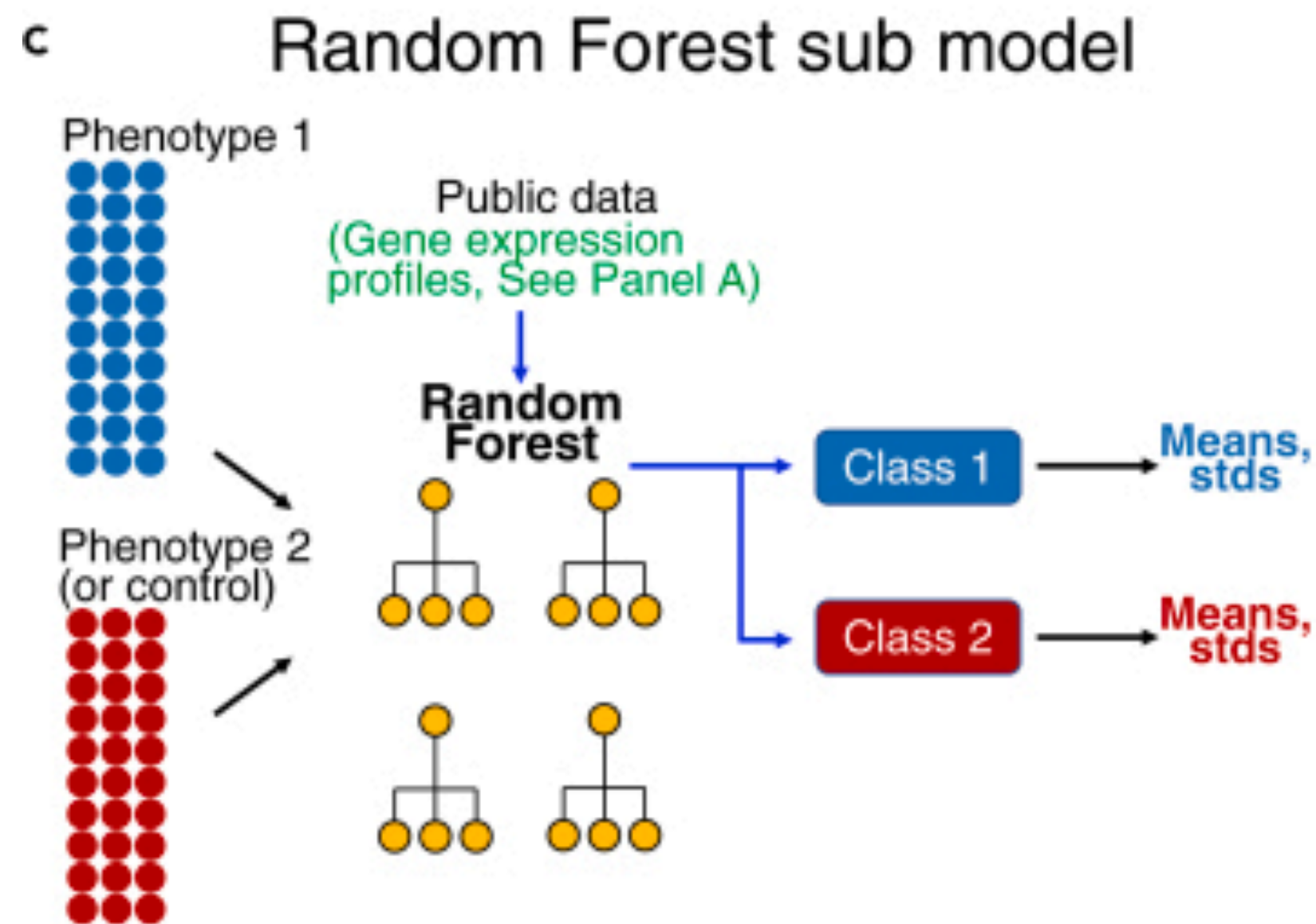**Hit overlaps with known psoriasis hits**
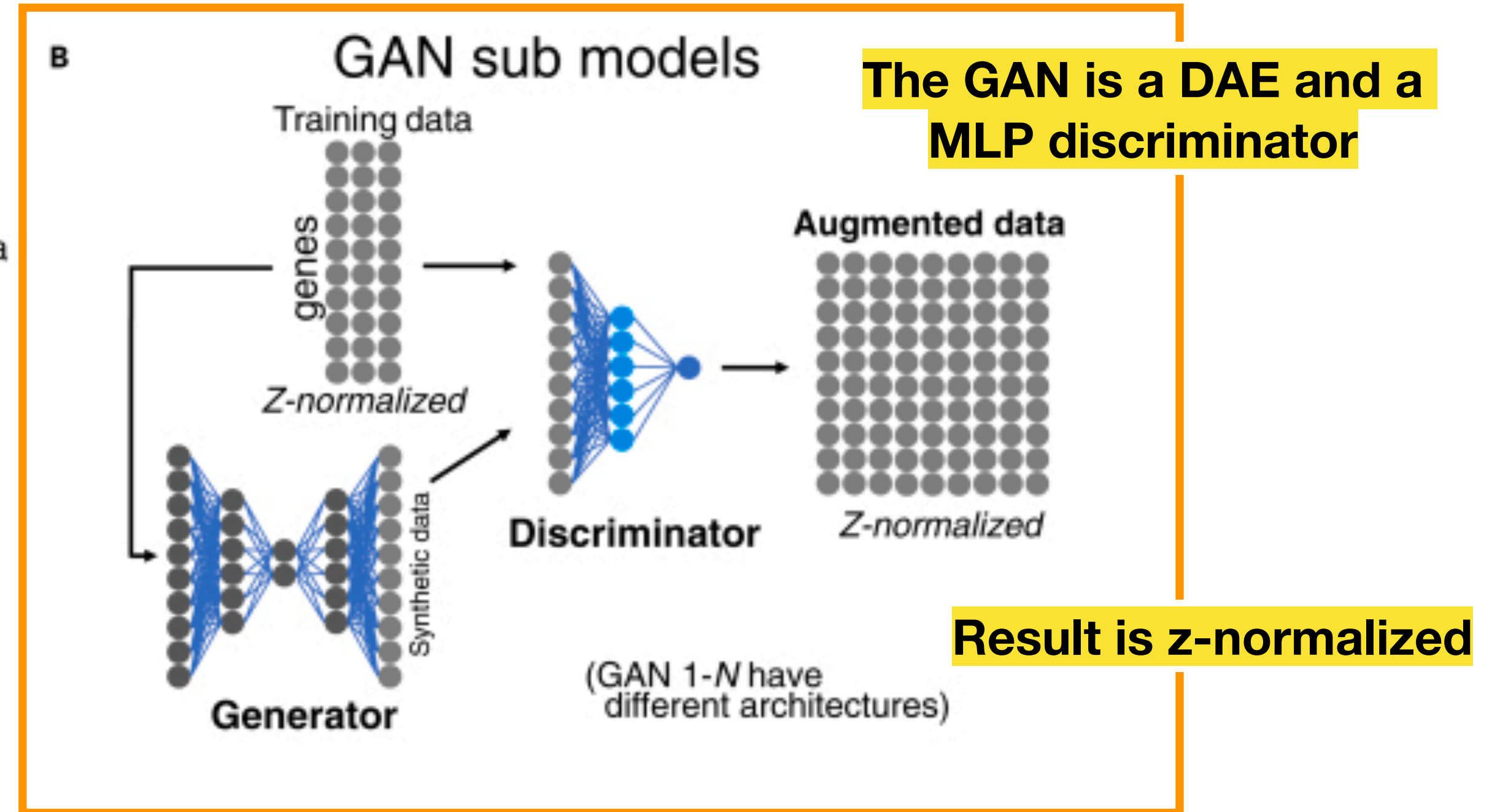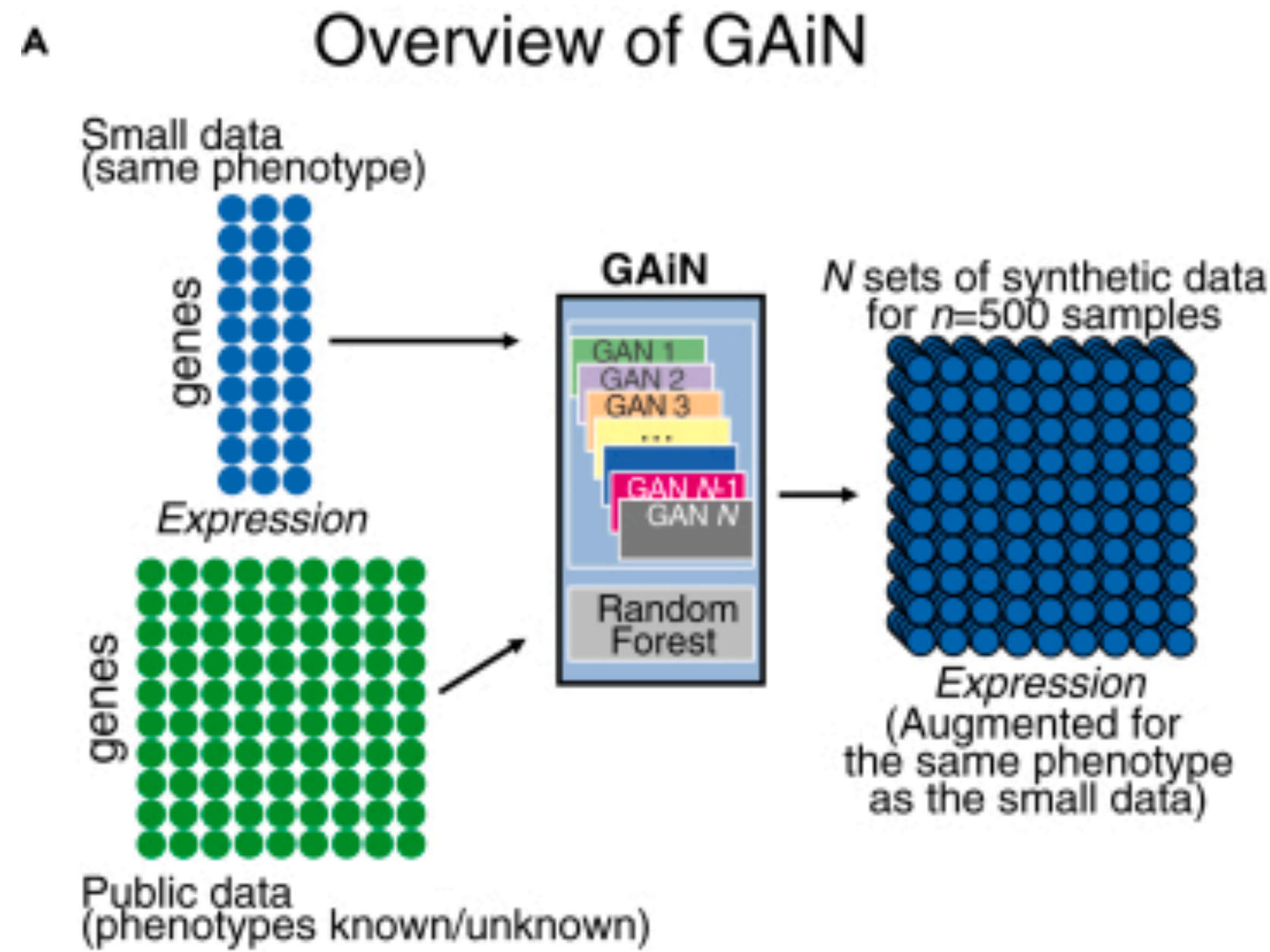
# "Good Luck, Babe"

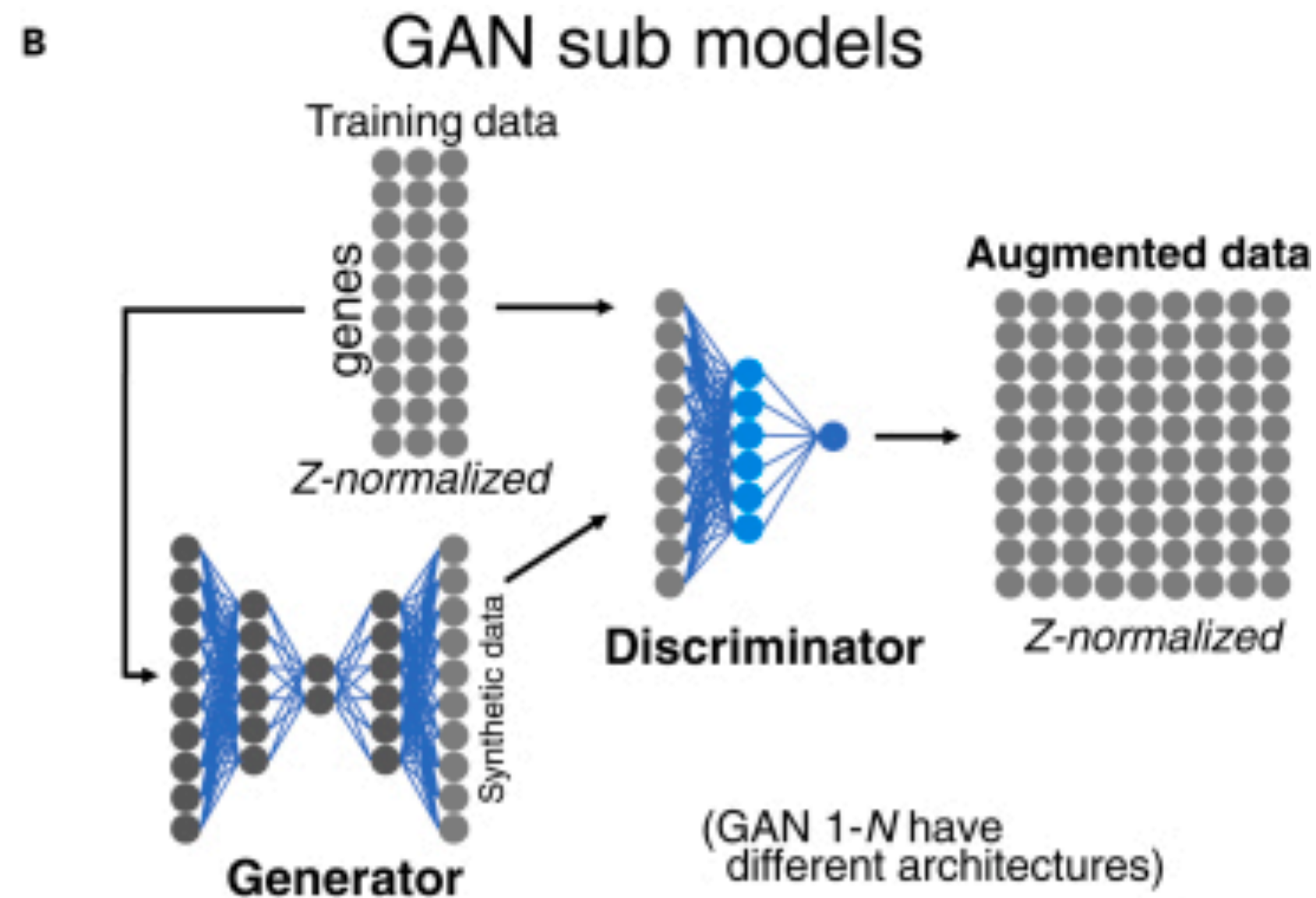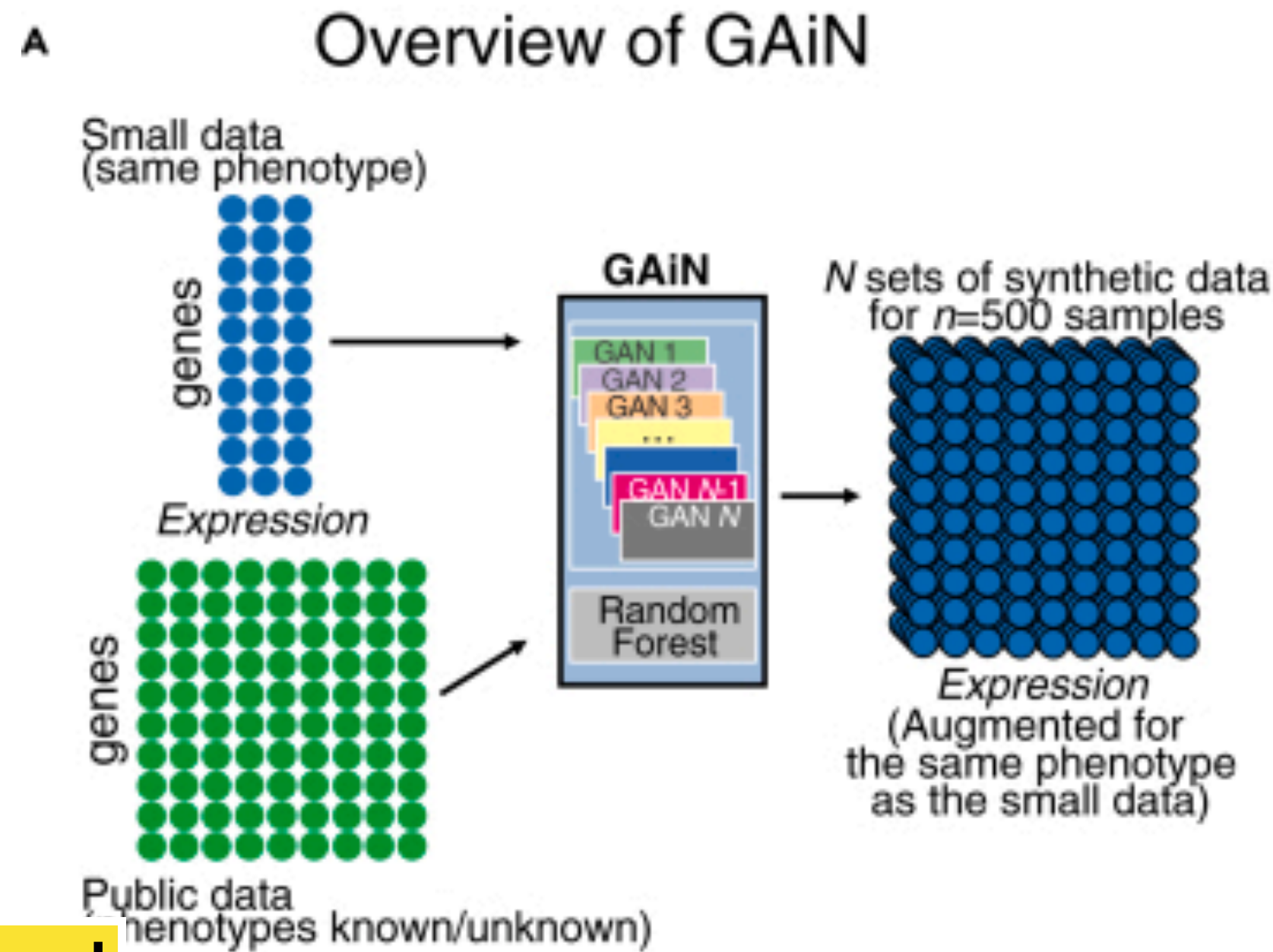## Bio Euphoria — Integrating Clinical & Molecular Data

# GAiN: An integrative tool utilizing generative adversarial neural networks for augmented gene expression analysis (Waters et al, *Patterns*)

- Goal: Boost the power of expression experiments with small sample sizes

- Method:

    - Use a Generative Adversarial Network (GAN) comprised of a Denoising Autoencoder (DAE) and an MLP

    - H: Removing the noise in small sample experiments will boost their power

- Result: Able to recapitulate an experiment of N=533 with just 10 samples!
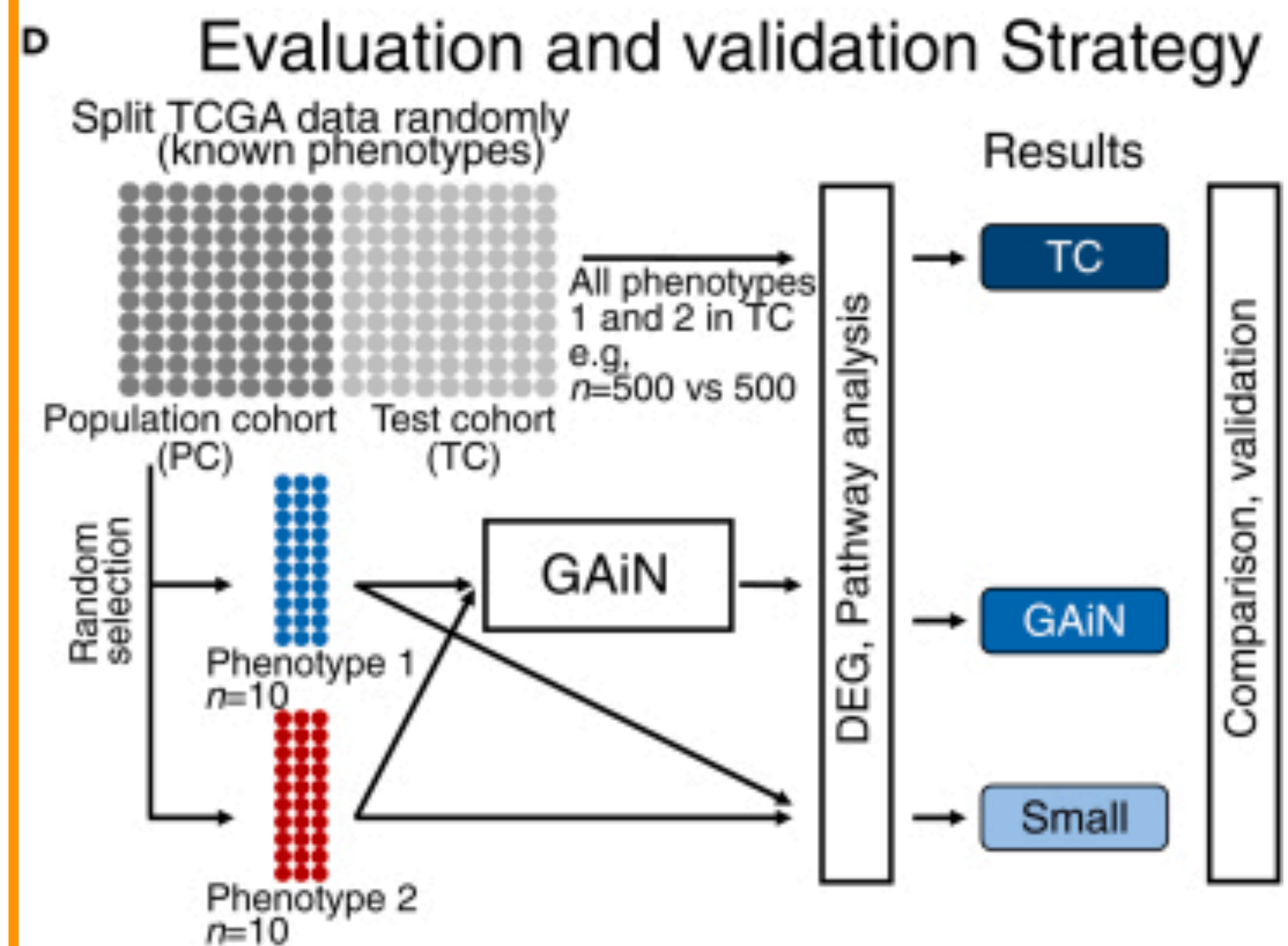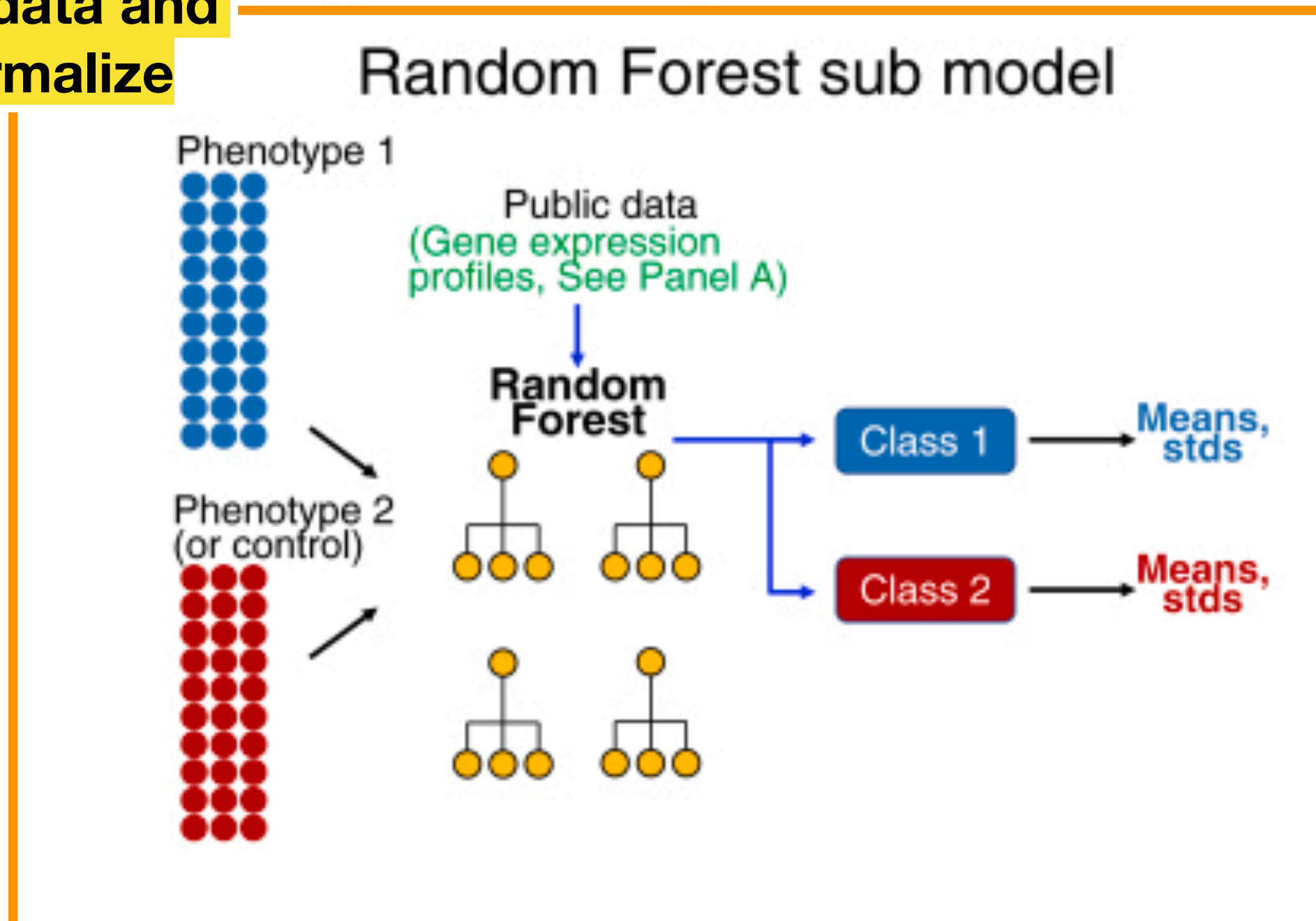
- Conclusion: This is either very clever or they are violating the laws of information theory

Uses a GAN and an Random Forest

**A** Overview of GAiN

**B** GAN sub models

The GAN is a DAE and a MLP discriminator

Result is z-normalized

**C** Random Forest sub model

**D** Evaluation and validation Strategy

Use population data and RF to un-z-normalize

**A** Overview of GAiN

**B** GAN sub models

**C** Random Forest sub model

**D** Evaluation and validation Strategy

Downsample from a larger dataset to validate

## DE genes from using TC, GAiN, and Small data



**Better Overlap with real data than equal amount of small data**

**Better prediction of real data than small data even if duplicated**

Figure S6. Rank comparison of DE genes for GAiN in phenotypic comparisons, related to Figure 3.



**Ranks are pretty darn close to real data ranks**

# Trainee Spotlight



**PRATISTHA GUCKHOOL**
DREXEL UNIVERSITY

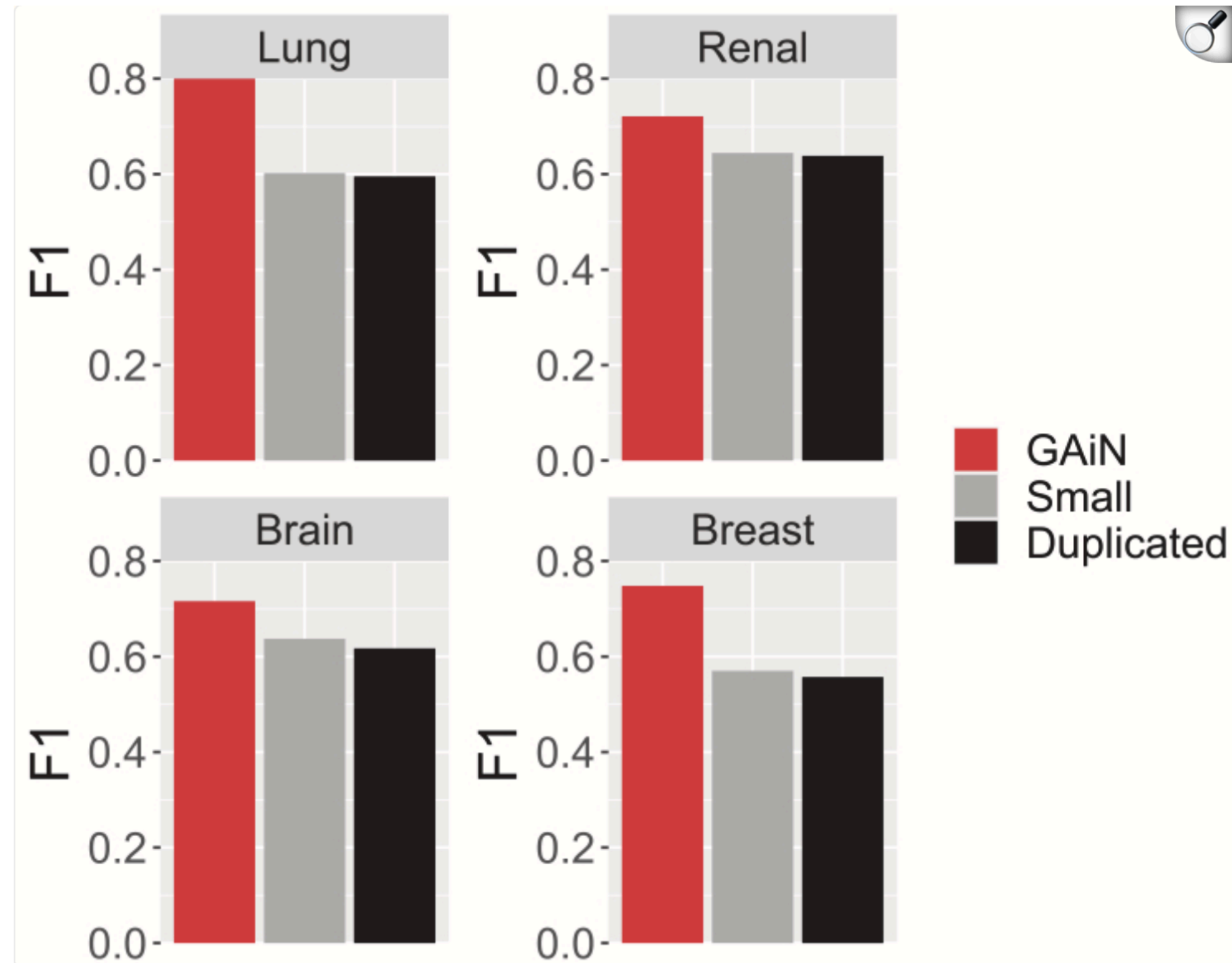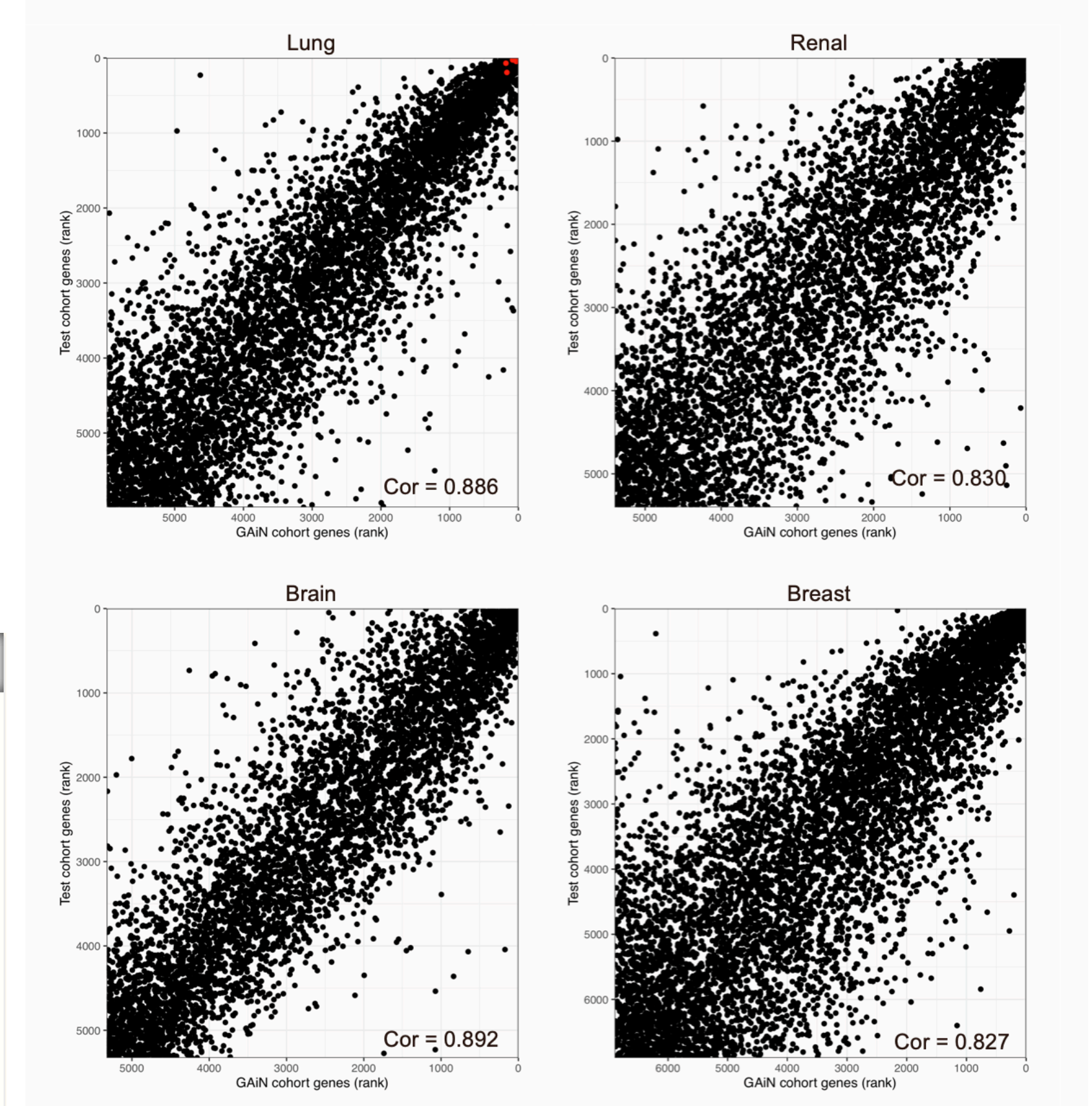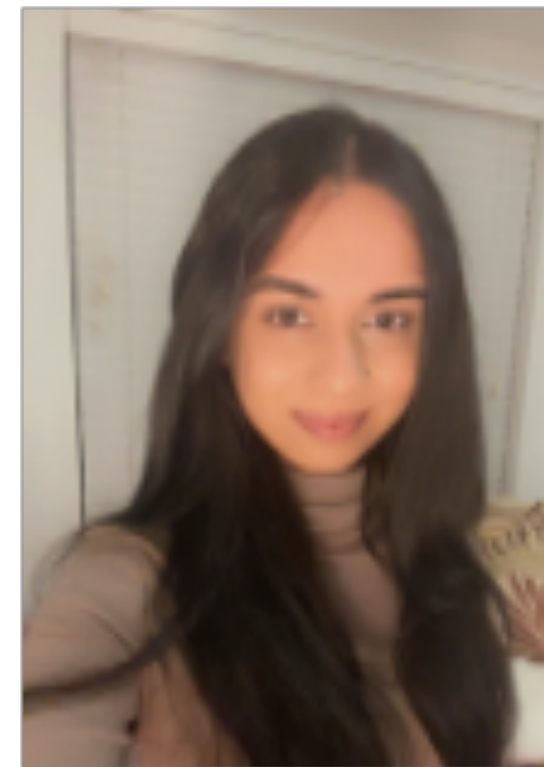# Scalable and unbiased sequence-informed embedding of single-cell ATAC-seq data with CellSpace (Tayyebi et al, *Nature Methods*)

- Goal: Improve ATAC-seq analysis by integrating a genome sequence context

- Method:

  - Train a new embedding model that learns from the k-mers of "open" events (i.e. regions of the genome with enough reads)

  - Jointly learn with the cells to produce cell-sequence embeddings

- Result:

  - Embeddings naturally control for batch/donor effects

  - Can predict motif-activity on cell-by-cell basis

  - Creates a denser cell-by-event matrix for use with other analytical methods

- Conclusion: Sequence context is information that can be captured and integrated into our datasets

**a**

Cells

Accessibility events
(peaks/tiles)

Genome ⟶

**b**

Accessibility across all cells

Negative sampling
for cell labels

Sampling genomic sequence

...ACTTCCGGGTC...

Push towards
'positive' cells

Pull 'negative'
cells away

ACTTCCGG

CTTCCGGG

TTCCGGGT

TCCGGGTC

N-grams of k-mers

Training the embedding of
cells and k-mers with CellSpace

**c**

CCACTTCCTGT

CCACTTCC
CACTTCCT

**Each k-mer that is
"open" is used to
generate a training set**

**b**

Accessibility across all cells

Negative sampling for cell labels

Sampling genomic sequence

...ACTTCCGGGTC...

Push towards 'positive' cells

Pull 'negative' cells away

ACTTCCGG

CTTCCGGG

TTCCGGGT

TCCGGGTC

N-grams of k-mers

Training the embedding of cells and k-mers with CellSpace

**c**

CCACTTCCTGT

CCACTTCC
CACTTCCT
ACTTCCTG
CTTCCTGT

Direct embedding
● k-mers
■ Cells

Induced embedding
● Accessibility events (peaks/tiles)
▲ TF motifs

**Producing a sequence-based embedding space**

Cells

Accessibility events (peaks/tiles)

Genome

**This embedding space recovers known biology**



**… like TF motifs**

**… and cell types**

Correlated by not same as other methods

Significantly outperforms other methods

# Accelerating histopathology workflows with generative AI-based virtually multiplexed tumour profiling (Pati et al, *Nature Machine Intelligence*)

- Goal: Synthezsize multiplexed immunohistochemistry (IHC) images from standard H&E slides

- Method:

  - VirtualMultiplexer trained on *unpaired* H&E and IHC data to predict IHC markers (AR, NKX3.1, CD44, CD146, p53, ERG) from H&E alone

  - Graph Transformer to predict clinically relevant endpoints

- Result:

  - model-generated IHC images were indistinguishable from real ones

  - Improve time to get IHC from weeks to seconds

- Conclusion: Pretty incredible results and love their "Visual Turing Test" evaluation metric

**Unpaired H&E and IHC as input**

**Multiple loss functions**
- **neighborhood consistency loss (path level)**
- **Global consistency loss**
- **Local consistency loss (cell type priors)**

**Real Input**

**Real Reference**

**Virtual Staining**

**Indistinguishable by experts**



**a** Real H&E staining

**b** Real IHC staining

**c** Virtual IHC staining

**d** FID

**e** Visual Turing test

**f** Staining quality

Unpaired S2S translation methods
- CycleGAN
- CUT
- CUT + KIN
- AI-FFPE
- VirtualMultiplexer

Staining quality
- Acceptable
- Background
- Border
- Not acceptable

**Molecular morphology that the model gets right**

**… and some it doesn't**

**Visual ablation study**

**Evaluating each of those losses and their effects on the virtual stain**

# A cell atlas foundation model for scalable search of similar human cells (Heimberg et al, *Nature*)

- Goal: Build a scRNA-Seq foundation model that enables super fast searches

- Method:

  - Introduce SCimilarity which embeds cells into low-dimensional space while preserving overall representation of the profile

  - Trained on 23.4 million cells from 412 scRNA-seq studies

    - Use **triplet loss** to force similar cells together and **reconstruction loss** to maintain subtle expression variation

  - Precompute an Approximate Nearest Neighbor Index using hnswlib

- Result:

  - Performs great on benchmarks and is super fast:

    - 10,000 similar cells retrieved in 0.05 seconds from a 23.4-million-cell reference dataset

- Conclusion: Allows for rapid scalable searches through massive single cell datasets

**Triplet loss not only pushes similar samples together, it pushes dissimilar ones apart**

**Accurately annotates cells (and does it in 0.02s)**



**b** Author labels

**c** SCimilarity

B cell
CD8+ αβ T cell
cDC
Endothelial cell
Epithelial cell of proximal tubule
Collecting duct intercalated
Collecting duct principal
Connecting tubule epithelial
Distal convoluted tubule epithelial
Kidney interstitial fibroblast
LoH TAL
LoH tDL
Macrophage
Mast cell
NK T cell
Monocyte
Myofibroblast
NK cell
Non-classical monocyte
Parietal epithelial
Plasma cell
pDC

**d**

Predicted annotation

Author-annotated cells (%)
0  20  40  60  80  100

Cell annotations better than (or as good as) others

# Supervised discovery of interpretable gene programs from single-cell data (Kunes, Walle et al, *Nature Biotechnology*)

- Goal: Use prior knowledge to extract biologically meaningful gene programs from single-cell RNA-seq data

- Method: Introduce the Spectra algorithm, a matrix factorization method constrained by existing biological knowledge

- Result:

  - Outperforms existing methods (expiMap, Slalom, NMF)

    - Finds immune checkpoint therapy (ICT) response factors in CD8+ T cells

    - Identifies tumor-reactive vs. exhausted T cell programs (prev. methods struggled here)

    - Predicts patient response to anti-PD-1 therapy

- Conclusion: Biological constraints on an unsupervised approach can get you a good balance between discovery and realism

Penalty function forces the matrix factorization to align the graph

Build gene-gene graphs from existing knowledge

Ready for downstream analysis

Eval on breast cancer dataset

Better overlap with expected genes

Better reconstruction of held out genes

Faster!

**Differentiating tumor-reactive CD8 cells from exhausted CD8 cells is difficult**



**Cells that expand are still active**

**Previous methods can't separate these cells**

**Spectra can separate these cells**

# "What Was I Made For?"

## Biomarker Discovery & Validation

# Identifying the joint signature of brain atrophy and gene variant scores in Alzheimer's Disease (Cruciani et al, *JBI*)

- Goal: Move beyond univariate modeling in AD

- Method:

  - Derive individual matrices to represent imaging and genetics separately

  - Using partial least squares to identify joint latent space

  - Validate using permutation testing and subsequent transcriptional analysis

- Result:

  - EPHX1 (Biological oxidation pathway) → Linked to subcortical atrophy

  - BCAS1 (Myelination process) → Temporal lobe atrophy (especially dentate gyrus)

- Conclusion: Statistical (not AI) based approaches are still relevant

Latent spaces separate patients fro controls well

Partial least squares components map to relevant brain structures

# A structurally informed human protein–protein interactome reveals proteome-wide perturbations caused by disease mutations (Xiong et al, *Nature Biotechnology*)

- Goal: To predict protein-protein interaction interfaces without direct structural information

- Method:

  - An ensemble of four deep learning methods (structure-structure, sequence-sequence, structure-sequence, and sequence-structure)

  - Graph Convolutional Networks and RNNs

- Result:

  - Predicts alleles enriched for disease-associated mutations

  - Applied to 11k cancer genomes and found 586 "oncoPPIs"

- Conclusion: Adds a critical level of understanding to PPIs; goes way beyond what we've seen from computational methods in this space before.

**The training set**

**A small amount of PPIs are resolved structurally**

**Using Proteins with >1 partner allows for predicting "partner-specific" interfaces**

**a**

90.92%    6.24%    2.83%
*H. sapiens*
146,138

94.69%    1.27%    4.03%
*A. thaliana*
44,630

89.61%    8.77%    1.62%
*S. cerevisiae*
35,193

97.12%    1.13%    1.75%
*D. melanogaster*
30,682

95.10%    1.14%    3.76%
*C. elegans*
12,201

68.56%    23.47%    7.96%
*M. musculus*
6,944

82.15%    11.33%    6.52%
*S. pombe*
3,557

43.16%    2.18%    54.65%
*E. coli*
2,750

■ Co-crystalized
■ Homology modeled
■ Unresolved

**b**

Training set

1 partner
1,826 (65.45%)

2 partners
777 (27.85%)

3 partners
125 (4.48%)

> 3 partners
62 (2.22%)

**c**    Structure–Structure model

Target protein structure and sequence          Partner protein structure and sequence

BIOP          BIOP

**d**    Sequence–Sequence model

Target protein sequence          Partner protein sequence

**The ensemble**

**Outperforms both state-of-the-art Structure-based and Sequence-based methods**

**Enrichment of mutation burden**

**Example of predicted PPI interfaces**

**The cancer protein-protein-interactome**

**… is predictive of survival in TCGA**

I'm told this is experimental  validation of the predicted interface 🤷



c

CDK4     CCND1     TSC2

d

| | | | | | | |
|---|---|---|---|---|---|---|
| − | + | − | + | − | + | TSC2-3× MYC |
| + | − | + | − | + | − | CDK4-3× MYC |
| − | − | − | − | + | + | CCND1 (Glu162Lys)-3× FLAG |
| − | − | + | + | − | − | CCND1 (Lys114Arg)-3× FLAG |
| + | + | − | − | − | − | CCND1-3× FLAG |

Input

250 —     ← TSC2    Anti-MYC WB

37 —     ← CDK4

37 —     ← CCND1 or CCND mutant    Anti-FLAG WB

IP

250 —     ← TSC2    Anti-MYC WB

37 —     ← CDK4

37 —     ← CCND1 or CCND mutant    Anti-FLAG WB

kDa

# "I Remember Everything"

## EHR, Real-World Evidence, & Epidemiology

# Genetic factors associated with reasons for clinical trial stoppage (Razuvayevskaya et al, *Nature Genetics*)

- Goal: Investigate if genetic factors play a role in clinical trial failures

- Method:

  - Fine-tuned a BERT model to analyze free-text reasons for stoppage in 28k clinical trials (CT.gov)

  - Link trials to target genetics

  - Stratify trials by the amount of genetic evidence

- Result:

  - Trials with weak genetic evidence are more likely to stop due to lack of efficacy (OR=0.61, p<0.0000…001)

  - Highly constrained genes (i.e. loss intolerant) are more likely to result in safety issues if targeted

- Conclusion: SOTR NLP/LLMs continue to make unstructured data available for research!

Reasons for stopping trials

**Studies that have stopped are less likely to have genetic evidence**



**a** Exposure: human genetic evidence

Odds ratio

| | | n | OR (95% CI) | P value |
|---|---|---|---|---|
| **All studies** | | | | |
| Phase IV | | 6,189 | 1.31 (1.27–1.34) | $<1 \times 10^{-40}$ |
| Phase III+ | | 14,616 | 1.39 (1.35–1.42) | $<1 \times 10^{-40}$ |
| Phase II+ | | 22,708 | 0.7 (0.68–0.72) | $<1 \times 10^{-40}$ |
| **Stopped studies** | | | | |
| Terminated | | 2,789 | 0.76 (0.73–0.79) | $<1 \times 10^{-40}$ |
| Withdrawn | | 810 | 0.72 (0.67–0.78) | $2.4 \times 10^{-20}$ |
| Suspended | | 119 | 0.73 (0.61–0.88) | 0.00076 |
| **Stop reason** | | | | |
| COVID-19 | | 37 | 0.91 (0.65–1.27) | 0.62298 |
| Safety | | 192 | 0.86 (0.75–1) | 0.05637 |
| Business or administrative | | 756 | 0.79 (0.73–0.85) | $2.6 \times 10^{-10}$ |
| Study design | | 178 | 0.68 (0.58–0.79) | $1.7 \times 10^{-7}$ |
| Insufficient enrolment | | 947 | 0.61 (0.57–0.65) | $<1 \times 10^{-40}$ |
| Negative | | 284 | 0.61 (0.54–0.69) | $6.0 \times 10^{-18}$ |

All indications

If the Target is constrained, more side effects

Low tissue-specific —> more side effects

More interacting partners —> more side effects

# An open-source framework for end-to-end analysis of electronic health record data (Huemos et al, *Nature Medicine*)

- Goal: Democratize EHR data analysis

- Method:

    - Develop open source modular tools based on existing data analysis standards

    - Implement a suite of commonly used analytical methods and enable multi-modal data integration

- Result:

    - Demonstrate on 4 different datasets with EHR data (e.g. UKB)

        - Phenotype stratification; biomarker discovery; causal inference; risk prediction

- Conclusion: EHRs are quickly becoming one of our most valuable research assets; the more people with access, the better

**Use AnnData format**

**Integrated data pre-processing tools**

**Various automatic modeling capabilities**

Other acute lower respiratory infections
Lung diseases due to external agents
Chronic lower respiratory diseases
Other diseases of upper respiratory tract
Other diseases of the pleura
Other diseases of the respiratory system
Influenza and pneumonia

Count
$10^0$   $10^1$   $10^2$   $10^3$

UMAP2

seasonal influenza virus identified ($n = 1$)
Influenza with other respiratory manifestations, virus not identified ($n = 1$)
Other disorders of lung ($n = 1$)
Other specified pleural conditions ($n = 4$)
Pneumonia due to streptococcus pneumoniae ($n = 1$)
Pneumonia due to staphylococcus ($n = 1$)
Pneumonia, unspecified ($n = 277$)
Respiratory failure, unspecified ($n = 1$)
Viral pneumonia, unspecified ($n = 1$)

UMAP1

**Found ICD10 code for "unspecified pneumonia"**

**e**

Pneumonia unspecified - annotated

UMAP2

UMAP1

**Used ehrapy to annotate them by their clinical features**

- Mild bacterial pneumonia ($n = 78$)
- Sepsis-like pneumonia ($n = 28$)
- Severe pneumonia with co-infection ($n = 74$)
- Viral pneumonia ($n = 97$)

**f**

Death
Length of stay
Sputum bacteria
Sputum fungi
Lymphocytes
Neutrophils
PCT
CRP
ALT
AST
GGT
Platelets
Age in months

Viral pneumonia

Mild bacterial pneumonia

Severe pneumonia with co-infection

Sepsis-like pneumonia

z-score
7.5
5.0
2.5
0
-2.5
-5.0
-7.5
-10.0

**g**

**Built timelines of disease**

Drug class

— Antibiotics (mg)
— Antivirals (mg)
— Catecholamines (mg ml⁻¹)
— Corticosteroids (mg ml⁻¹)
— Anticoagulants (mg)
— Electrolytes (ml)

Cultures negative        Cultures negative          *A. baumannii*

Medication categories
- Administered
- Not administered

Laboratory measurements
- C-reactive protein (mg dl⁻¹)
- Monocytes ($10^9$/L)
- Neutrophils ($10^9$/L)
- PCT (ng ml⁻¹)
- ALT (U L⁻¹)

Unit
60
40
20
0

Other acute lower respiratory infections
Lung diseases due to external agents
Chronic lower respiratory diseases
Other diseases of upper respiratory tract
Other diseases of the pleura
Other diseases of the respiratory system
Influenza and pneumonia

UMAP2

seasonal influenza virus identified (n = 1)
Influenza with other respiratory manifestations, virus not identified (n = 1)
Other disorders of lung (n = 1)
Other specified pleural conditions (n = 4)
Pneumonia due to streptococcus pneumoniae (n = 1)
Pneumonia due to staphylococcus (n = 1)
Pneumonia, unspecified (n = 277)
Respiratory failure, unspecified (n = 1)
Viral pneumonia, unspecified (n = 1)

UMAP1

**Found ICD10 code for "unspecified pneumonia"**

**b**

Kaplan–Meier survival curves for all pneumonia groups

Rate of survival

**f**

**nnotate features**

Death
Length of stay
Sputum bacteria
Sputum fungi
Lymphocytes
Neutrophils
PCT
CRP
ALT
AST
GGT
Platelets
Age in months

Viral pneumonia

Mild bacterial pneumonia

Severe pneumonia with co-infection

Sepsis-like pneumonia

z-score
7.5
5.0
2.5
0
−2.5
−5.0
−7.5
−10.0

Groups
Viral pneumonia (n = 94)
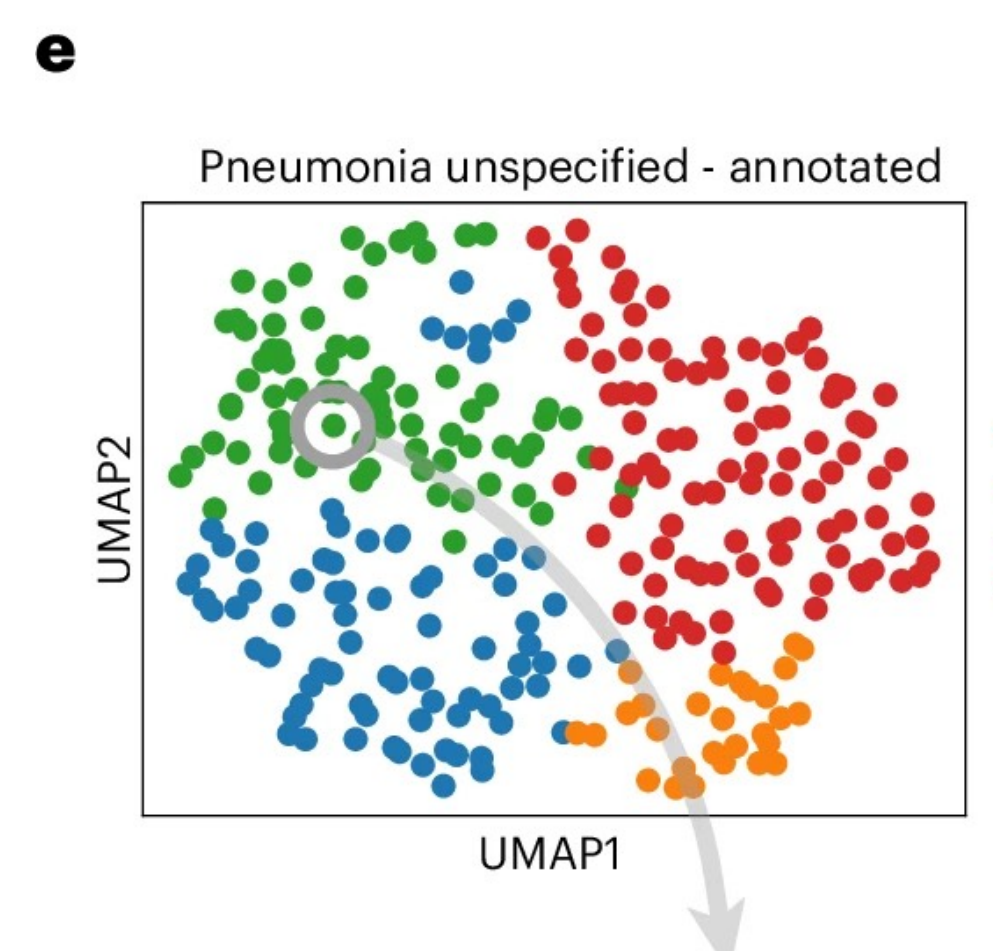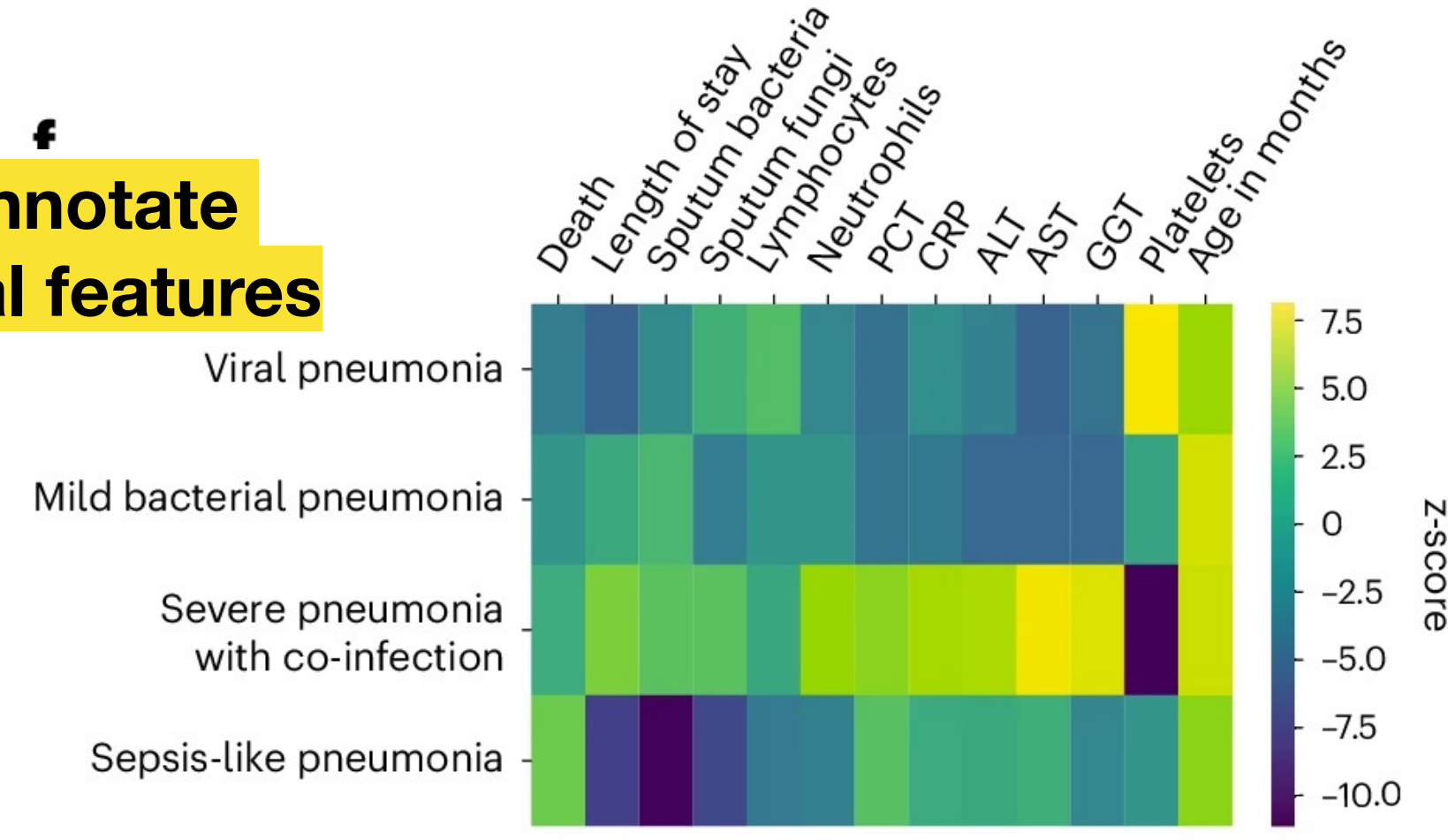Mild bacterial pneumonia (n = 58)
Sepsis-like pneumonia (n = 28)
Severe pneumonia with co-infection (n = 25)

Hours spent in ICU

**Found big survival differences**

**timelines of disease**

Drug c

Antibiotics (mg)
Antivirals (mg)
Catecholamines (mg ml⁻¹)
Corticosteroids (mg ml⁻¹)
Anticoagulants (mg)
Electrolytes (ml)

Medication categories
Administered
Not administered

60

40

Unit

20

0

Cultures negative
Cultures negative
*A. baumannii*

Laboratory measurements
C-reactive protein (mg dl⁻¹)
Monocytes (10⁹/L)
Neutrophils (10⁹/L)
PCT (ng ml⁻¹)
ALT (U L⁻¹)

# The biomedical knowledge graph of symptom phenotype in coronary artery plaque: machine learning-based analysis of real-world clinical data (Huan et al, *BioData Mining*)

- Goal: Better understand the precursors of heart disease and you'll better understand heart disease

  - i.e. Study the multifactorial nature of coronary artery plaque

- Method:

  - Use ~1,500 patients EHRs to get a diverse presentation of plaques and use for training/evaluation

  - Integrate genetic and pathway knowledge (e.g. STRING) to add mechanistic layer and to identify hub genes which could explain differences

- Result:

  - Identified 23 symptom phenotypes, 41 association rules, and 61 hub genes

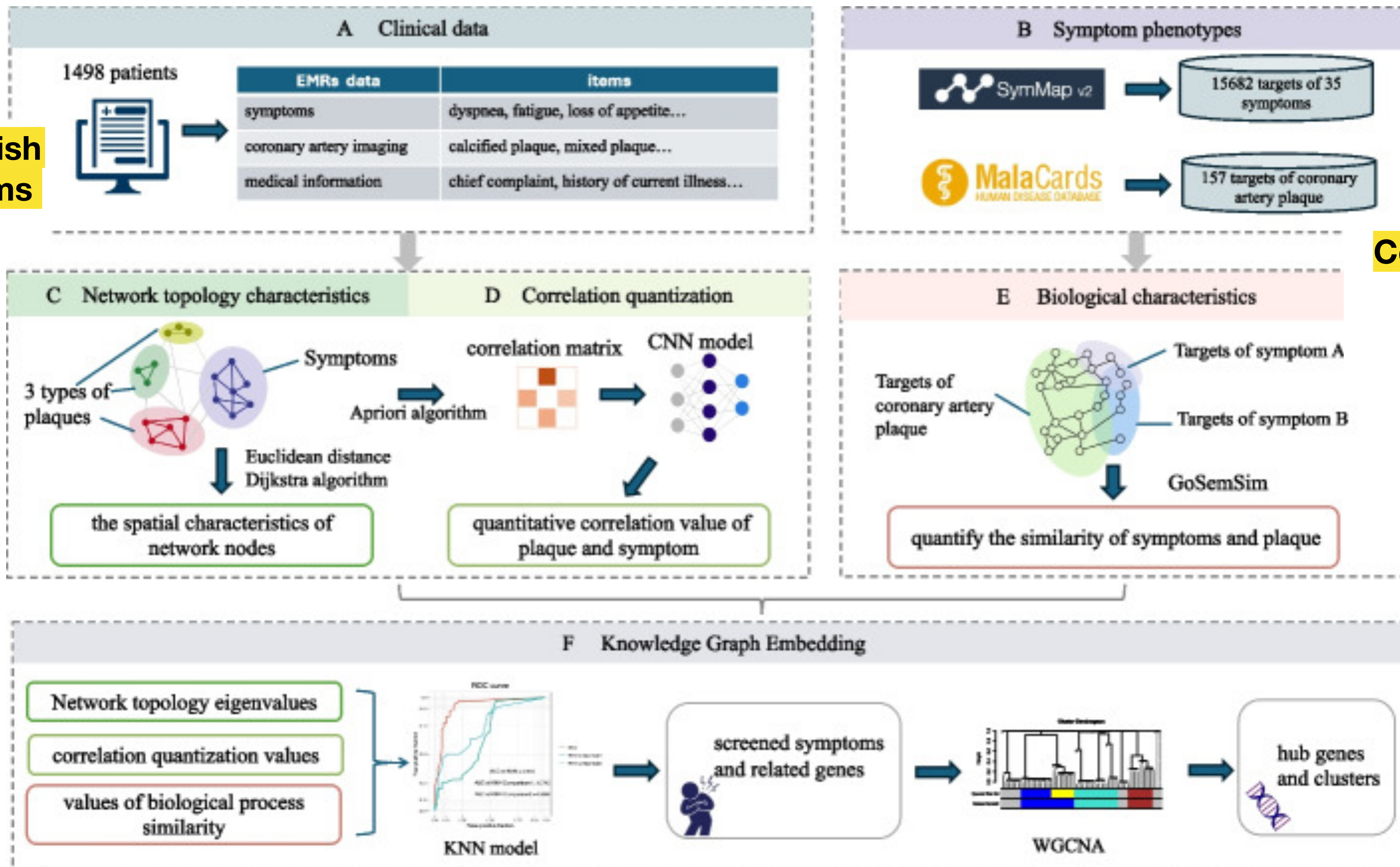  - Lipid metabolism and inflammation pathways are key drivers of symptom differences

- Conclusion: The EHR is a rich resource that's underutilized. Leveraging the symptomatic variances of how disease present is a great new avenue to study human disease.

Use the EHR to establish link between symptoms and plaque

Connect symptoms to biological network through annotated targets

Use classic ML approaches to build predictive models

**Connect symptoms to plaques to genes**

**Hub genes are enriched for lipid metabolism**

# "Houdini"

Emerging Therapeutics & Technologies

# Machine-guided design of cell-type-targeting cis-regulatory elements (Gosai, Castro et al, *Nature*)

- Goal: Design and validate synthetic cis-regulatory elements (CREs) that drive cell-type-specific gene expression

- Method:

  - A CNN trained on massively parallel reporter assay (MPRA) data to predict effect of a given sequence on cell-type-specific expression (~776 million sequences assayed!)
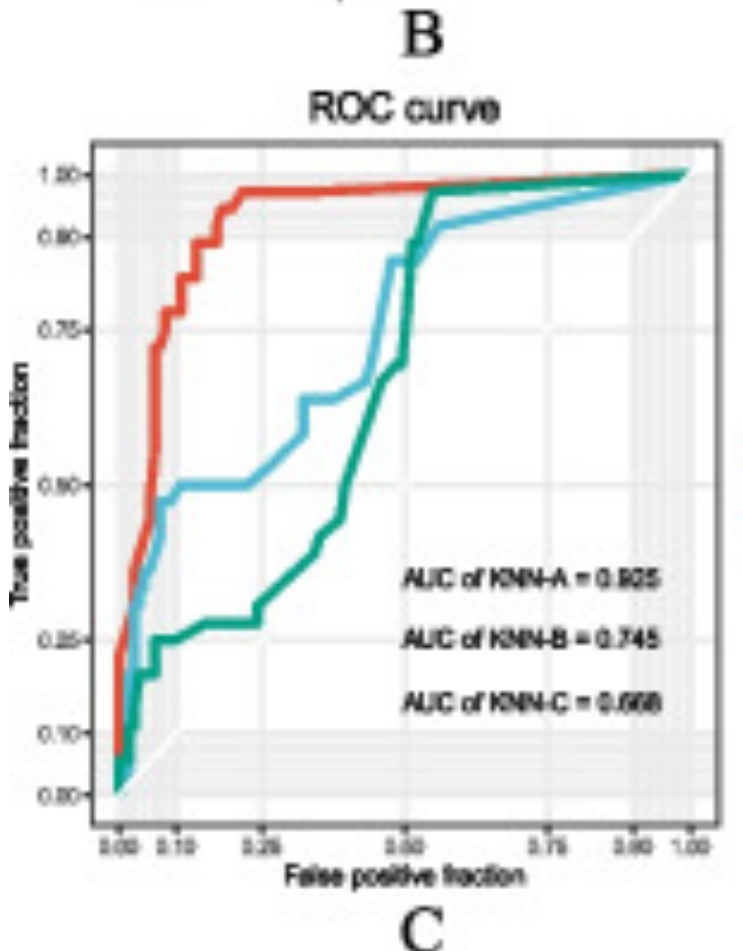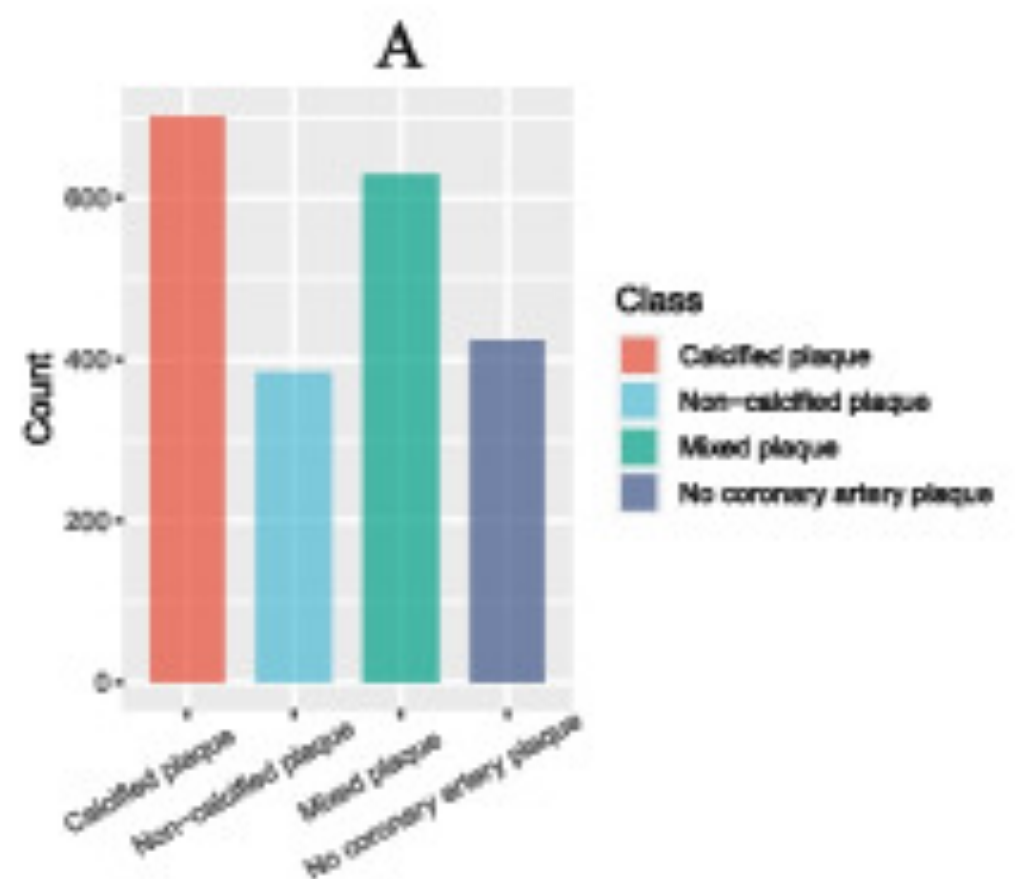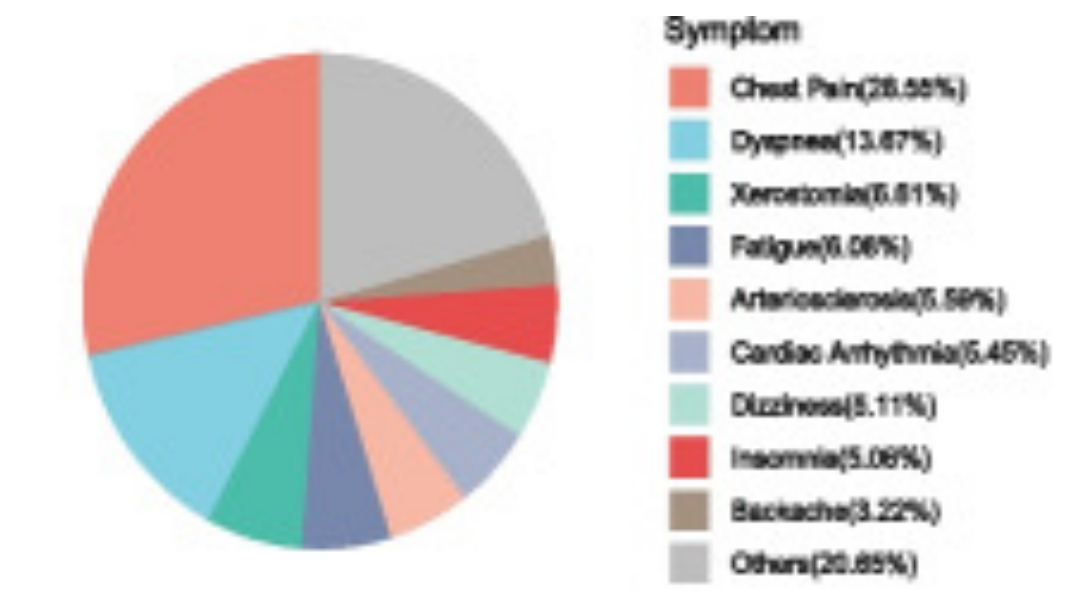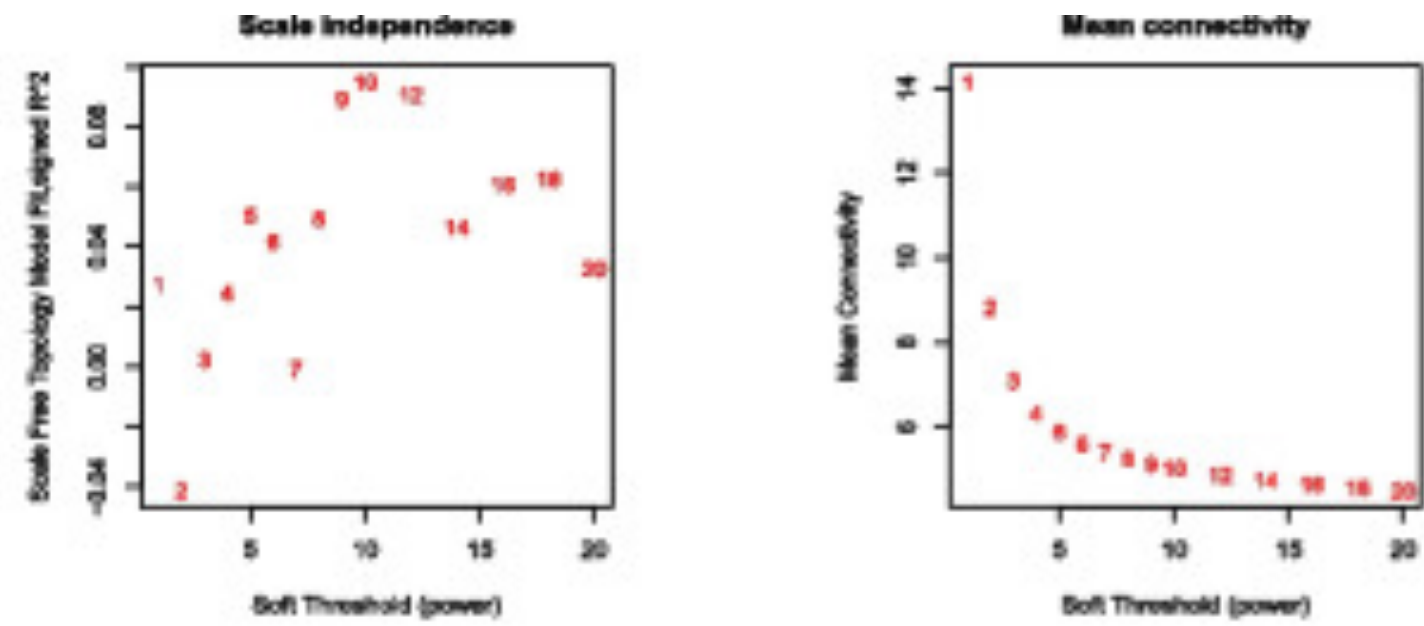
  - Since there are >2E120 possible sequences, need a strategy for selecting them — introduce CODA which is an optimization strategy

- Result:

  - Synthetic CREs outperformed natural sequences in driving cell-type-specific expression; confirmed with extensive *in vitro* testing

  - Regulatory grammar of synthetic CREs showed a distinct motif vocabulary

- Conclusion: Award for most mind-blowing paper this year. 🤯

**Architecture of Manilois**

**Data comes from these massively parallel reporter assays (MPRA)**

**a** Empirical reporter assays

200 bp

Natural and synthetic CRE oligos

Reporter transfection across cell types

$$\text{Activity} = \frac{\text{mRNA}}{\text{plasmid}}$$

Sequencing

CRE activity measurement

MPRA model training

**b** In silico CRE modelling

T A A G A T G T C

One-hot sequence encoding

Machine learning model

CRE activity prediction

CAGATAAGTGCAG

Sequence contribution scores

**Consistent with experimental measures**

**e** Malinois    STARR    DHS    H3K27ac

kb    kb    kb    kb

2,413 cCREs

Signal intensity

**Predicted v. Actual shows high correlation**

**c**

Empirical activity

Predicted activity

K562 r = 0.88    HepG2 r = 0.89    SK-N-SH r = 0.88

Experimental results show synthetic CREs are stronger than naturals

**Synthetic and naturals are quite a bit different**



**Synthetic and naturals are quite a bit different**

**Show striking specificity in both zebrafish and mice**

# mRNA-LM: full-length integrated SLM for mRNA analysis (Li et al, *Nucleic Acids Research*)

- Goal: Develop a small language model for mRNA that includes the coding region PLUS the 5' and 3' untranslated regions (UTR)

- Method:

  - Use contrastive learning (CLIP) but instead of text + image, it's 3 different mRNA sequences to construct a joint language model for all three regions

  - Pretrained BERT-based models on each of the 3, then combined using CLIP to create mRNA-LM

  - Fine-tuned mRNA-LM on mRNA half-life, translation rate, transcript expression, and protein expression

- Result: Significantly outperformed other methods on the trained tasks; Performed well at zero shot tasks

- Conclusion: Nice clear methods made this a joy to read. Recommended for your JCs!

# "Feather"

A Feather in Your Cap — Shout-Outs & Honorable Mentions

# A comparative study of large language model-based zero-shot inference and task-specific supervised classification of breast cancer pathology reports (Sushil et al, *JAMIA*)

- Goal: Compare state-of-the-art LLMs to supervised ML techniques for extracting important information from pathology reports

- Method:

  - 769 manually annotated breast cancer pathology reports from UCSF

  - LLMs (zero-shot inference): GPT-4, GPT-3.5, Starling-7B-beta, ClinicalCamel-70B

  - Supervised ML models: Random Forests, LSTM-Att, UCSF-BERT

- Result:

  - GPT-4 performed better or as well as the best supervised method, in particular on tasks with high label imbalance

  - Supervised models struggled with generalizability

- Conclusion: For information extraction, LLMs are great

~15k patients and ~50k reports available

Example of a report annotation

**Label distributions**



**Performances**

**Macro F1-Score**

| | Specimen type | Estrogen receptor status | Progesterone receptor status | HER-2 receptor status | Tumor grade | Lymph node involvement | Lymphovascular invasion | Margin status | DCIS Margin status | Tumor histology | Sites examined | Sites of disease |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **GPT-4** | 0.88 | 0.94 | 0.93 | 0.85 | 0.88 | 0.89 | 0.98 | 0.86 | 0.88 | 0.79 | 0.77 | 0.70 |
| **LSTM-ATT** | 0.83 | 0.77 | 0.83 | 0.73 | 0.84 | 0.49 | 0.97 | 0.41 | 0.54 | 0.80 | 0.83 | 0.78 |
| **RANDOM FORESTS** | 0.75 | 0.76 | 0.79 | 0.69 | 0.53 | 0.44 | 0.70 | 0.36 | 0.30 | 0.57 | 0.69 | 0.49 |
| **UCSF-BERT** | 0.89 | 0.60 | 0.59 | 0.67 | 0.67 | 0.49 | 0.80 | 0.52 | 0.36 | 0.37 | 0.41 | 0.29 |
| **GPT-35-TURBO-16K** | 0.53 | 0.67 | 0.61 | 0.49 | 0.78 | 0.53 | 0.61 | 0.45 | 0.62 | 0.38 | 0.59 | 0.32 |
| **STARLING-7B-BETA** | 0.59 | 0.39 | 0.39 | 0.20 | 0.37 | 0.55 | 0.33 | 0.26 | 0.20 | 0.42 | 0.44 | 0.22 |
| **CLINICALCAMEL-70B** | 0.45 | 0.67 | 0.59 | 0.30 | 0.48 | 0.17 | 0.30 | 0.37 | 0.22 | 0.34 | 0.12 | 0.08 |

# Proteogenomic analysis reveals non-small cell lung cancer subtypes predicting chromosome instability, and tumor microenvironment (Song, Choi, Kim, Hwang et al, *Journal*)

- Goal: To redefine NSCLC subtypes based on **proteogenomic and multiomics analysis**; and link these subtypes to clinical outcomes

- Method:

  - Multiomics integration of genomic (WES), transcriptomic (RNA-seq), proteomic (TMT-labeling), phosphoproteomic, and acetylproteomic data

  - Non-negative matrix factorization (NMF) clustering to identify molecular subtypes

    - Good for heterogenous data distribution, scales, and sparseness; allows for "soft-clustering"

- Result:

  - Found 5 molecular subtypes

  - 1 of which is new (hypoxic subtype) and associated with poor outcomes

- Conclusion: Right method with the right data has the right result.

I'm a sucker for this style
of a Table 1

How these subtypes relate to those found by others

Driving factors can be identified

Novel subtype is associated with poor survival

# Novel machine learning model for predicting cancer drugs' susceptibilities and discovering novel treatments (Cao et al, *JBI*)

- Goal: Predict drug effect on cancer cells using genetics (a tale as old as time)

- Method:

  - Introduce Kernalized Residual Stacking (KRS)

    - Multi-task learning method (each drug's effect on the line is a task); tasks can share information through residual correction

    - Radial Basis Function (RBF) kernel reduces dimensionality

    - Allows for feature importance to be quantified to facilitate interpretability

- Result:

  - Outperforms benchmarks; Identified the PI3K-Akt pathway as a key cancer drug response regulator

  - Identify 8 novel cancer drug repurposing candidates

- Conclusion: I like the information sharing aspects of this, although wish they pushed it a bit further.

**Can identify drivers of performance (e.g. ERBB2)**

**Better at predicting IC50 values**

# Reverse Metabolomics for Discovery of Chemical Structures from Humans (Gentry et al, *Nature*)

- Goal: Only 10% of metabolites in metabolomics studies are identified; Identify them!

- Method:

  - Generated mass spec spectra for 2,430 molecules and do a reverse database lookup on 1.2 billion metabolomics results

  - Associate the molecule with phenotypes from the metabolimcs studies

- Result:

  - Identified 139 previously unreported bile acid conjugates

  - Some were significantly associated with disease (e.g. Crohn's disease)

- Conclusion: This is why we make our data public!

Generate new spectra

Figure out the metabolites

Match existing spectra

# Single-cell chemoproteomics identifies metastatic activity signatures in breast cancer (Pillai et al, *Science Advances*)

- Goal: Develop a single-cell chemoproteomics platform to measure protein activity rather than just protein/mRNA abundance

- Method:

  - single-cell activity-dependent proximity ligation (scADPL), a microfluidic-based chemoproteomic method to measure active proteins

  - Applied to breast cancer cell lines and patient-derived organoids (PDOs)

  - Focused on a six-enzyme panel (Ag-6) involved in cancer aggressiveness

- Result:

  - Identified increased enzyme activity in highly metastatic breast cancer cells

  - Showed that **enzyme activity**, not just abundance, correlates with tumor aggressiveness and metastatic potential

- Conclusion: A omics method that measures protein activity?? Sign me up!

**A**

1. Treatment with activity probe(s)

2. Single cell trapping and lysis

Whole proteome

3. Probe and target protein proximity barcoding

4. Activity barcode ligation and amplification

Active proteins-of-interest (POIs)

Inactive POIs and whole proteome

Valve

I  II  III  IV  V

**B**

**C**

Barcoded amplicons

Off-chip quantification (qPCR or ddPCR)

Inputs

Output/Retrieval ports

Bright field (BF)

Isolated single cell

Immunofluoresence

Isolated FITC+ cell

**Enzyme activity is more sensitive than mRNA abundance**

# A prognostic and predictive computational pathology immune signature for ductal carcinoma in situ: retrospective results from a cohort within the UK/ANZ DCIS trial (Li et al, *Lancet Digital Health*)

- Goal: Evaluate computational pathology biomarker (CPath TIL) to quantify tumor infiltrating lymphocytes (TIL) density in breast cancer

  - This is a critical measure, but time consuming to collect

- Method:

  - Perform a retrospective analysis of already completed Randomized Controlled Trial (n=755)

  - Computationally estimate TIL and perform survival and interaction analysis

- Result:

  - CPath TIL-high patients had a significantly higher risk of recurrence (HR 2.10, p = 0.0004)

  - Invasive progression risk was even higher (HR 3.09, p = 0.0013)

- Conclusion: Smart use of available RCT data to demonstrate computational methods

**CPath-TIL Low**    **CPath-TIL High**

**A**

HR 0·40 (95% CI 0·20–0·81); p=0·0079
10-year IBE-free survival with radiotherapy
92·1% (95% CI 85·8–95·7)
10-year IBE-free survival without
radiotherapy 82·6% (95% CI 77·0–87·0)

**B**

HR 0·32 (95% CI 0·19–0·54); p<0·0001
10-year IBE-free survival with radiotherapy
85·7% (95% CI 78·3–90·8)
10-year IBE-free survival without
radiotherapy 63·9% (95% CI 57·1–69·9)



Number at risk
(number of events)

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| No radiotherapy | 229 | (30) | 192 | (9) | 155 | (2) | 28 | 226 | (62) | 153 | (17) | 114 | (2) | 33 |
| Radiotherapy | 130 | (7) | 117 | (3) | 90 | (0) | 20 | 133 | (12) | 113 | (6) | 89 | (0) | 18 |

# Incorporation of emergent symptoms and genetic covariates improves prediction of aromatase inhibitor therapy discontinuation (Rattsev et al, *JAMIA Open*)

- Goal: Build predictive models breast cancer therapy discontinuation

- Method:

  - Use genetic, clinical features, and patient-reported outcomes with machine learning (CoxPH, Random Survival Forest, and GB)

  - 181 women from a *prospective cohort*

- Result:

  - AUROCs of about 70%

  - Identified *ESR1* variants and some symptoms as risk factors

- Conclusion: A recipe for translational work: good ML meets good study design

**Study setup**

**Can produce personalized survival curves**

# Developing and evaluating pediatric phecodes (Peds-Phecodes) for high-throughput phenotyping using electronic health records (Grabowska et al, *JAMIA*)

- Goal: Develop pediatric-specific EHR-phenotyping system

- Method:

  - Use EHRs (1M+ peds) + genetics (50k+ peds)

  - Adult phecodes with little peds data were removed or labeled as rare using clinical expertise

  - Adult phenols with many peds data were kept or split into new ones

  - Ran a PheWAS using these new codes to validate

- Result:

  - Found 2,051 Pediatric PheCodes, reclassified Pediatric conditions into 19 categories

  - Peds-Phecodes replicated more known genetic associations than phecodes (248 v. 192)

- Conclusion: Fantastic resource for EHR analysis of pediatric conditions and outcomes (which are woefully understudied)

**Stronger signals for what's expected**

# Inference of phylogenetic trees directly from raw sequencing reads using Read2Tree (Dylus et al, *Nature Biotechnology*)

- Goal: Develop a method (Read2Tree) to infer phylogenetic trees directly from raw sequencing reads, bypassing genome assembly and annotation

- Method: Aligns raw reads to orthologous genes, reconstructs sequences, and infers trees with improved speed and scalability

- Result: Up to 100× faster than assembly-based pipelines while maintaining high accuracy; successfully applied to yeast phylogenies and SARS-CoV-2 variant classification

- Conclusion: A useful tool for large-scale genomic comparisons — could have potential for TBI in the future!

**Some trees**

**Some calibration info on how it performs**

# "Too Sweet"

## Brain Candy — Tasty Tidbits & Intellectual Treats

# Poisoning Scientific Knowledge Using Large Language Models (Yang et al, *Nature Machine Intelligence*)

- Goal: Investigate how maliciously generated abstracts can manipulate scientific knowledge graphs (KGs), affecting downstream biomedical applications

- Method: Develop **Scorpius**, an AI model that injects false drug-disease relationships into KGs by generating fake abstracts

- Result: Even **one** false abstract can dramatically increase a drug's apparent relevance to a disease in KGs

- Conclusion: The inability of the LLM to "think" critically about the information it takes in is a MAJOR limitation

**Ranking is very sensitive to poisoning**



**d** Add one link

**e** Add two links

**f** Add three links

Hub node
- Leukocytopenia
- Agaricus catalepsy
- Methemoglobinemias

Disease node proportion
100%
80%

Disease-agnostic ranking
5
10
50
1000

Number of malicious links

CL285032, Allylglycine, Azaperone, BSH, Cervisol, Methsuximide, Carglumic acid, Lithium citrate, Astemizol, N-tritylmorpholine

# Coffee Drinking Timing and Mortality in US Adults (Wang et al, *European Heart Journal*)

- Goal: To determine whether the timing of coffee consumption influences all-cause and cause-specific mortality in US adults

- Method:

  - 40,725 adults from the National Health and Nutrition Examination Survey (NHANES, 1999–2018)

  - Validation cohort: 1,463 adults from the Women's and Men's Lifestyle Validation Study

  - Followed for 10 years tracking ~4k deaths; adjusting for confounders

- Result:

  - Morning coffee drinkers had significantly lower all-cause mortality compared to non-coffee drinkers (HR: 0.84; 95% CI: 0.74–0.95)

  - Lower CVD mortality was observed in morning coffee drinkers (HR: 0.69; 95% CI: 0.55–0.87)

  - All-day coffee drinkers did not show a significant mortality benefit compared to non-drinkers

  - Higher coffee intake was associated with lower all-cause mortality, but only for morning coffee drinkers (P-interaction = 0.031)

- Conclusion: ☕

# The Virtual Lab: AI Agents Design New SARS-CoV-2 Nanobodies with Experimental Validation (Swanson et al, *pre-print*)

- Goal: Develop an AI-driven virtual research team using large language models (LLMs) to automate and accelerate interdisciplinary scientific discovery

- Method:

  - Several LLM "agents": Principal Investigator, Biologist, Computational Scientist, Machine Learning Expert

  - hold virtual meetings, critique each other's work, and refine research strategies

  - Can use tools like ESM, AlphaFold, and Rosetta

- Result:

  - Applied to design nano bodies to bind SARS-CoV-2 variants

  - AI-generated 92 nanobody designs, with over 90% solubility and expression

- Conclusion: A step toward an closed loop AI scientific engine!

# Lookback at my predictions for 2024 ✨

❌ Temporal analysis of single cell sequencing data will emerge

✅ The start of a closed loop AI scientific method engine

❌ More biomechanics simulations used for pre-training models

✅ The rise of "foundation models"

RESULTS BY YEAR



2000                          2025

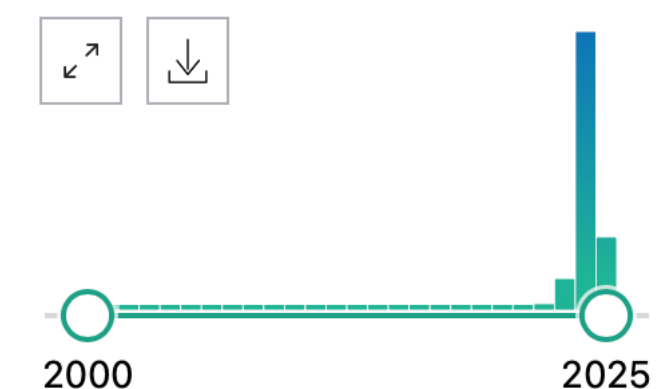✅ Multimodal deep learning will become the norm

✅ More examples of heterogeneous models being integrated together

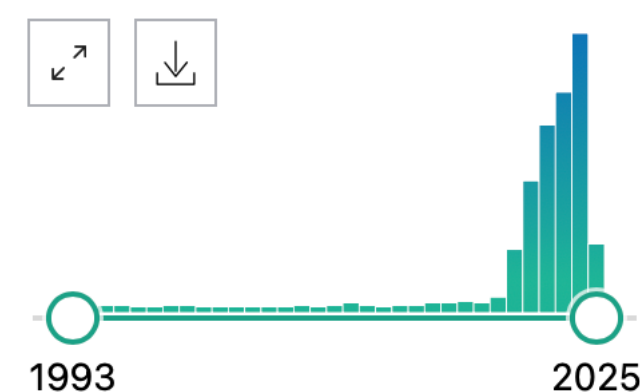❌ The year Mamba and Transformers battle it out

✅ More custom architectures for biology

RESULTS BY YEAR



❌ UMAP will FINALLY DIE! (wishful thinking)

1993                          2025

# Predictions for 2025 ✨

☐ Emergence of multimodal CLIP models for TBI, particularly with language

☐ Synthetic data will find some compelling use cases (I haven't seen one yet)

☐ Foundation models will have big impact on rare disease work

☐ Diffusion models for novel drug discovery coupled with experimental validation (and trials?)

☐ Uncertainty quantification in AI modeling will emerge

☐ AI efficiency boosts will lead to real time/streaming applications in TBI

☐ New explainable AI techniques (e.g. SAE) will begin to get used in TBI

☐ An initial biomedical application of quantum computing

☐ Lastly, new architectures that can leverage multimodal data (I don't think we've exhausted this at all)

# Thank you!