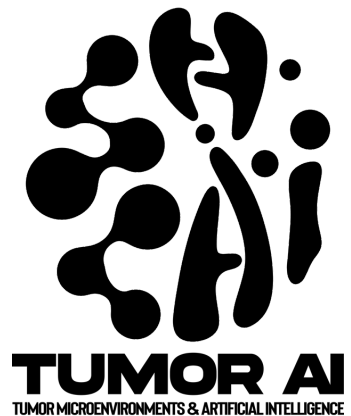


Gene attribute predictions through interpretable AI

Avi Sahu

TumorAI.org



Danger of uninterpretable AI

When AI Fails, Lives Are at Risk

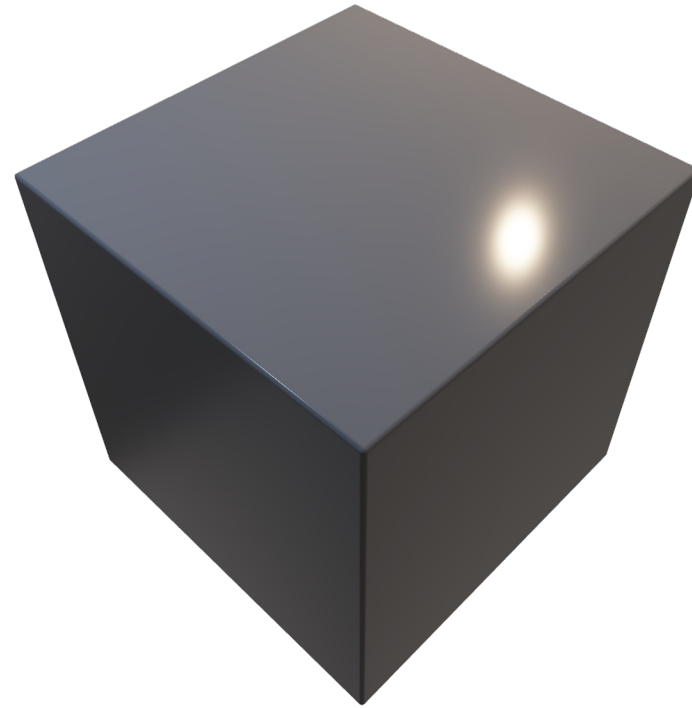
A significant number of AI failures are attributed to **interpretability** and **bias**.



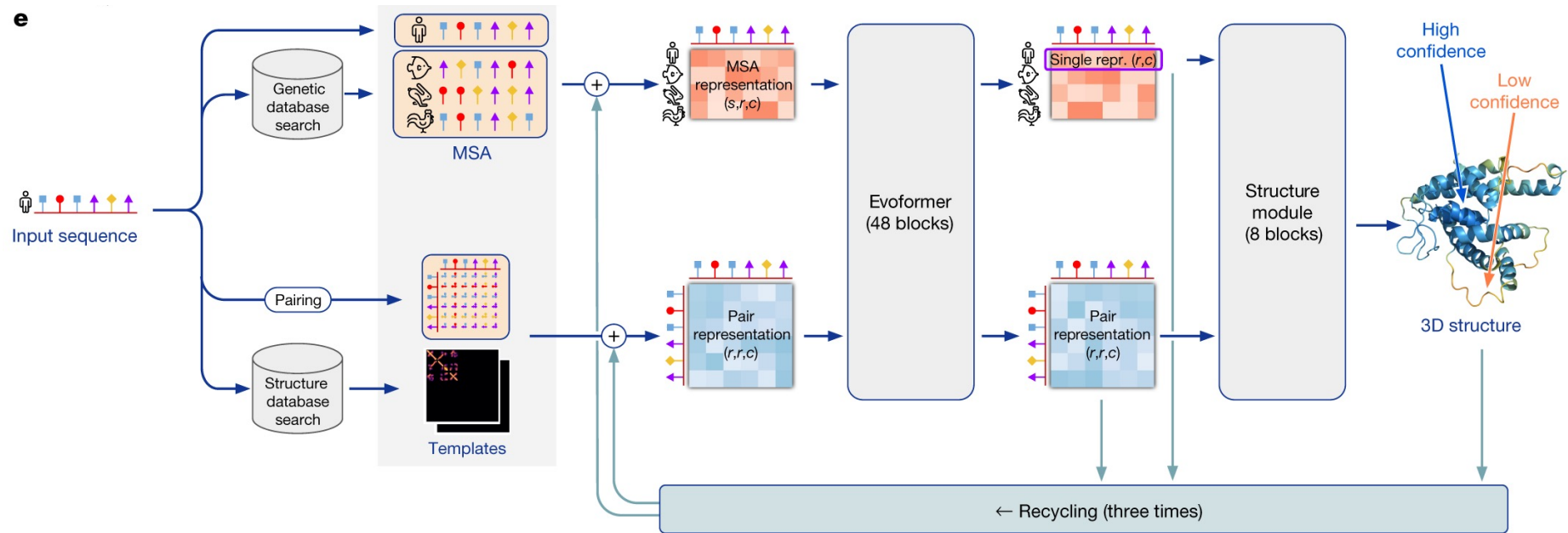
Transparency and
fairness are crucial for AI
adoption in genomics

Interpretability

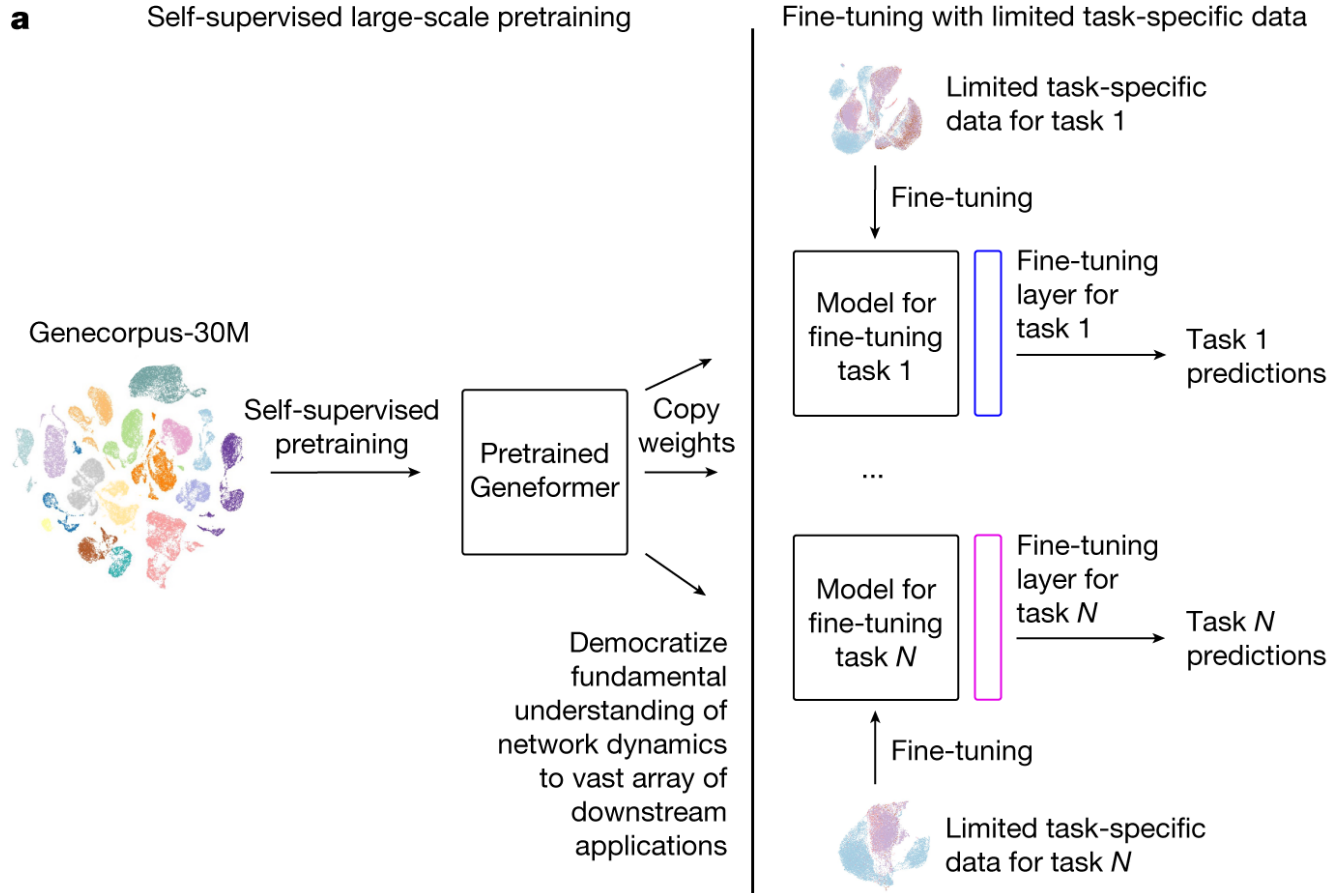
Bias



Predicting structure attributes from info amino acid sequence



Predicting cell or gene attributes from expression



Problem

Use only expression or sequence data but do not incorporate **text** – the extensive literature knowledge?



Ala
Jaraweh



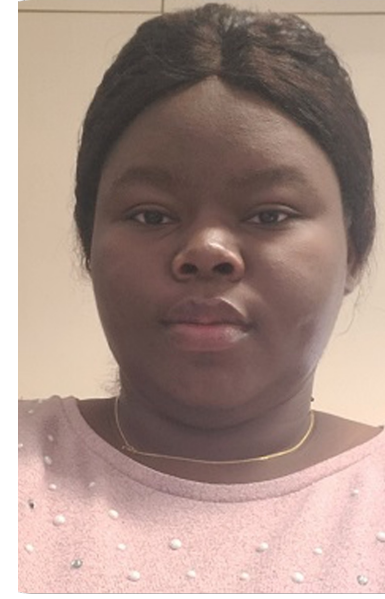
Oladimeji
Macaulay



David
Arredondo



Luis
Tofoya



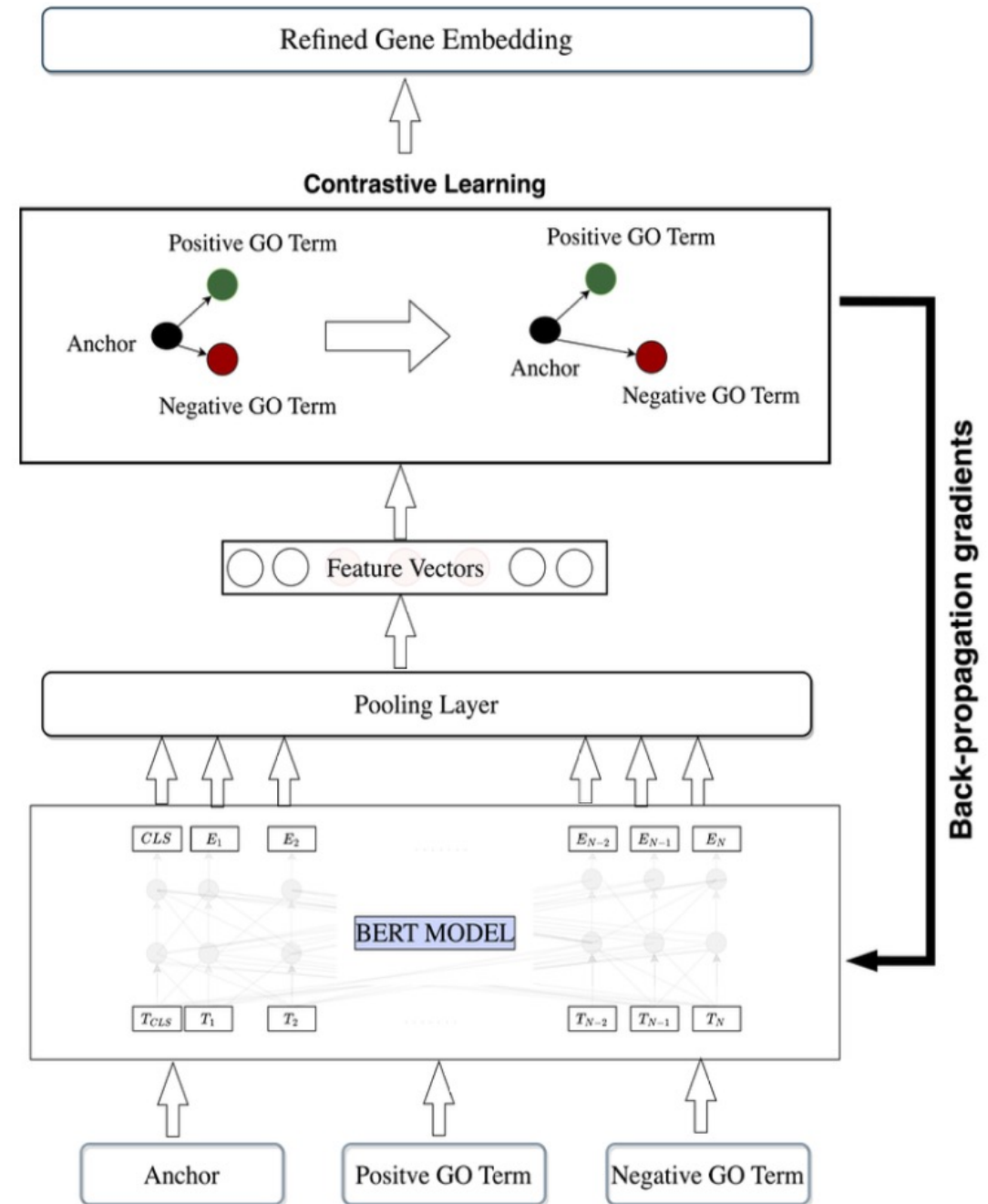
Funmi
Oyebamiji



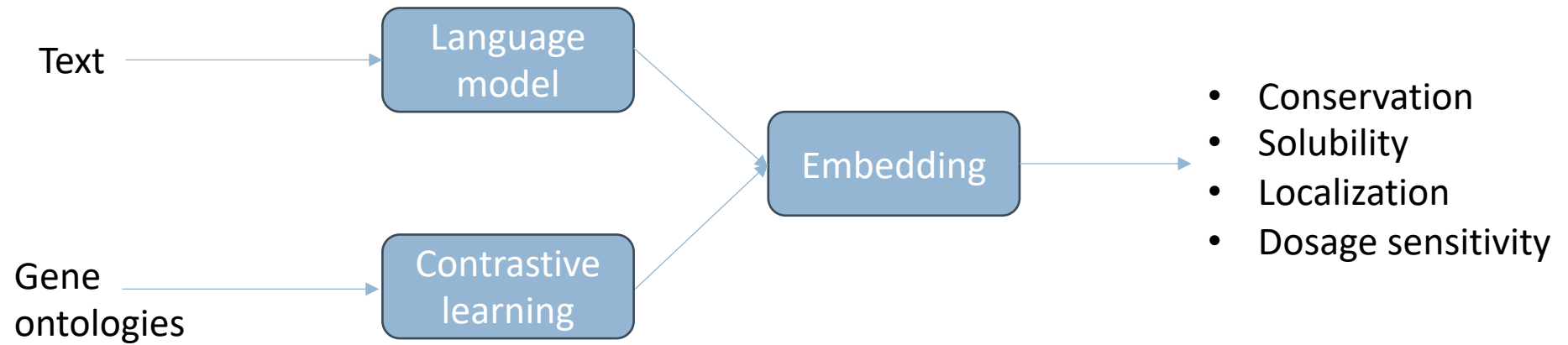
Kushal
Virupakshappa

Infer gene attributes using text

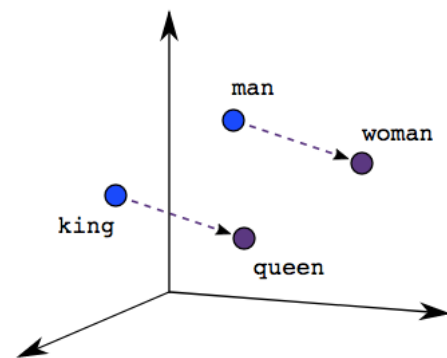
GeneLLM: A novel framework for deriving insights from text and ontologies



GeneLLM: A novel framework for deriving insights from text and ontologies

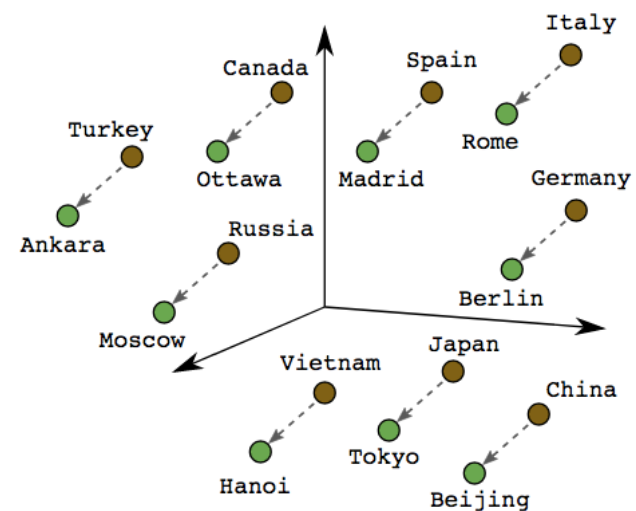


Learning relationships between objects with Embedding



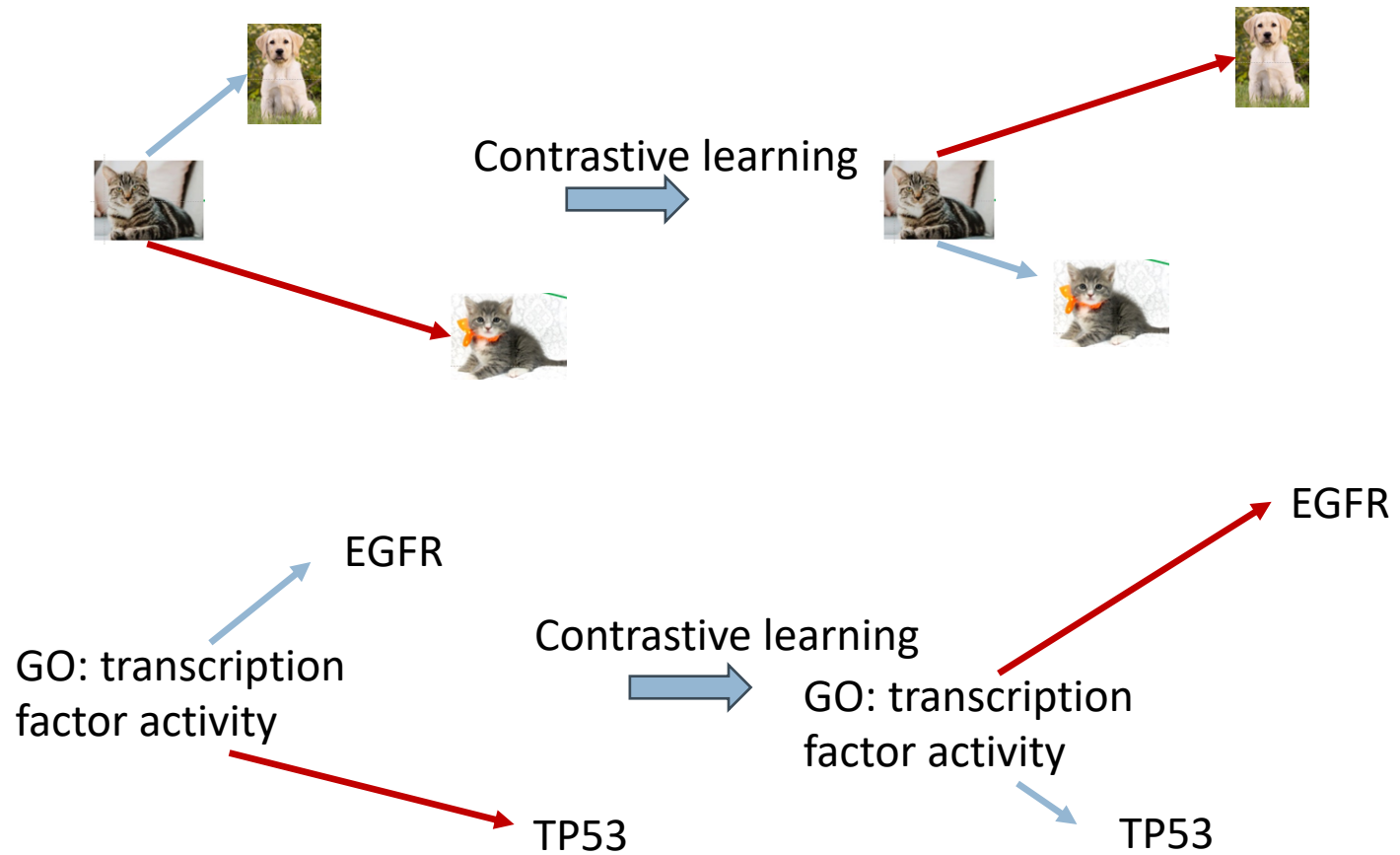
Closer points are more similar

Male-Female

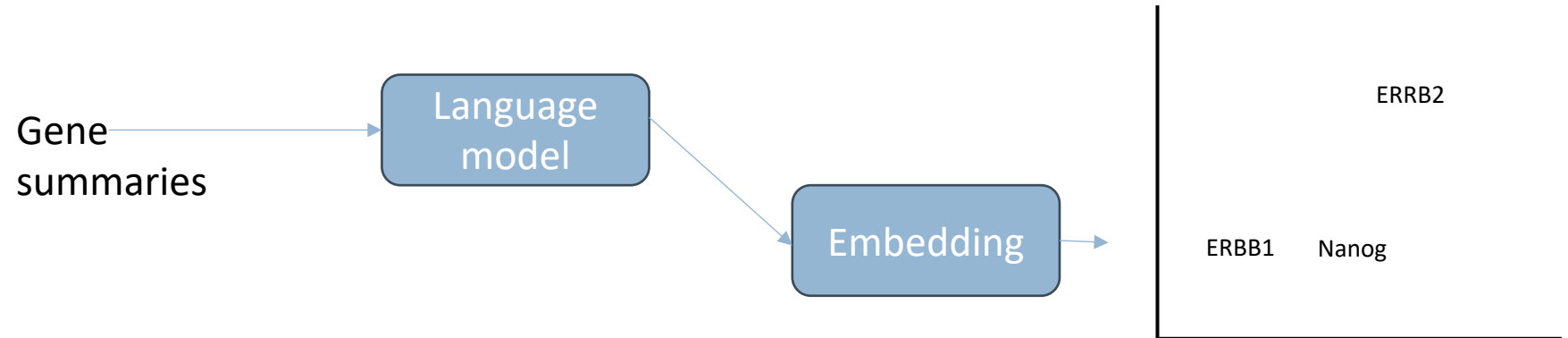


Country-Capital

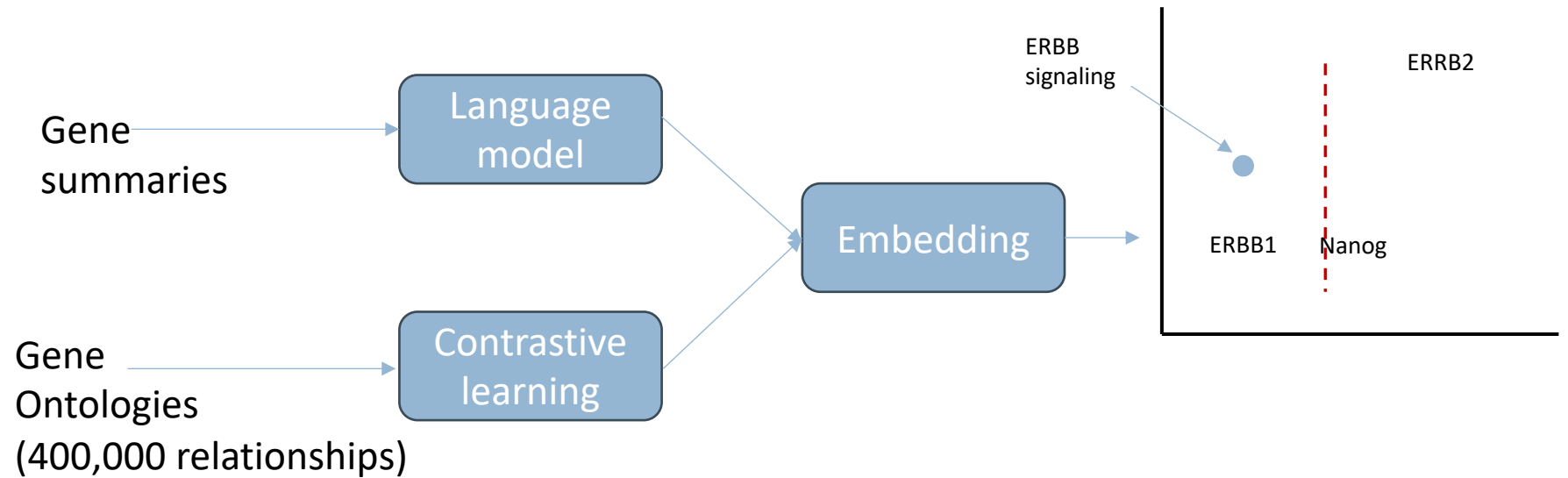
Contrasting similar and dissimilar Ontology-Gene relationships



GeneLLM: framework to extract insights from text and ontologies

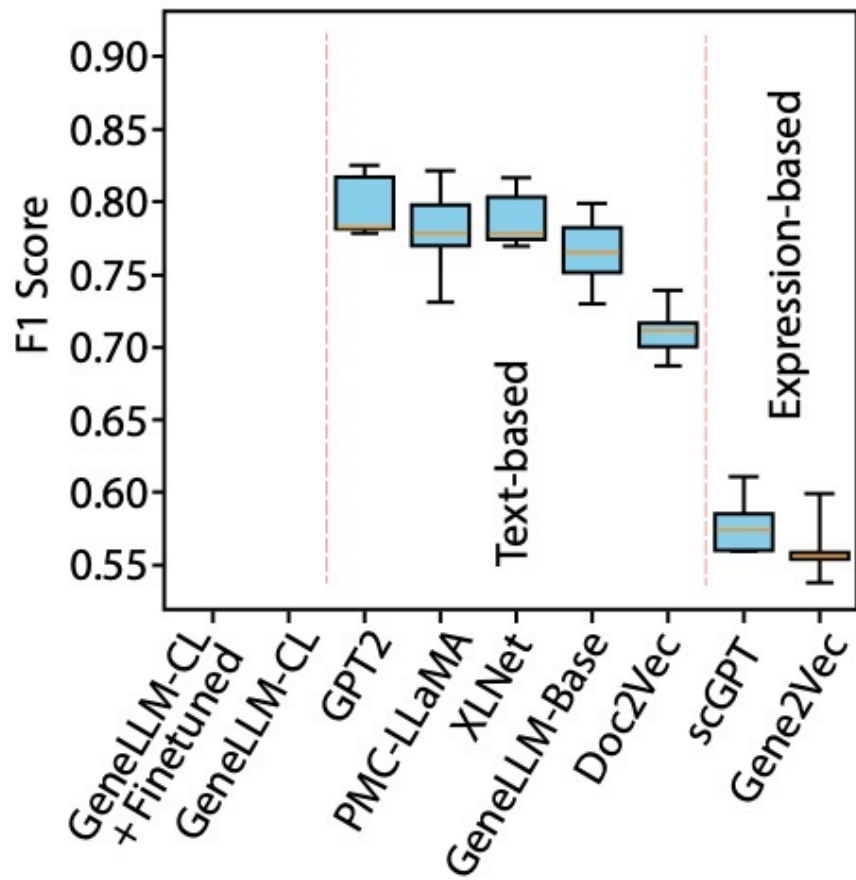


GeneLLM: framework to extract insights from text and ontologies



- Problem
- Solution
- Application

Text is more effective for protein solubility prediction

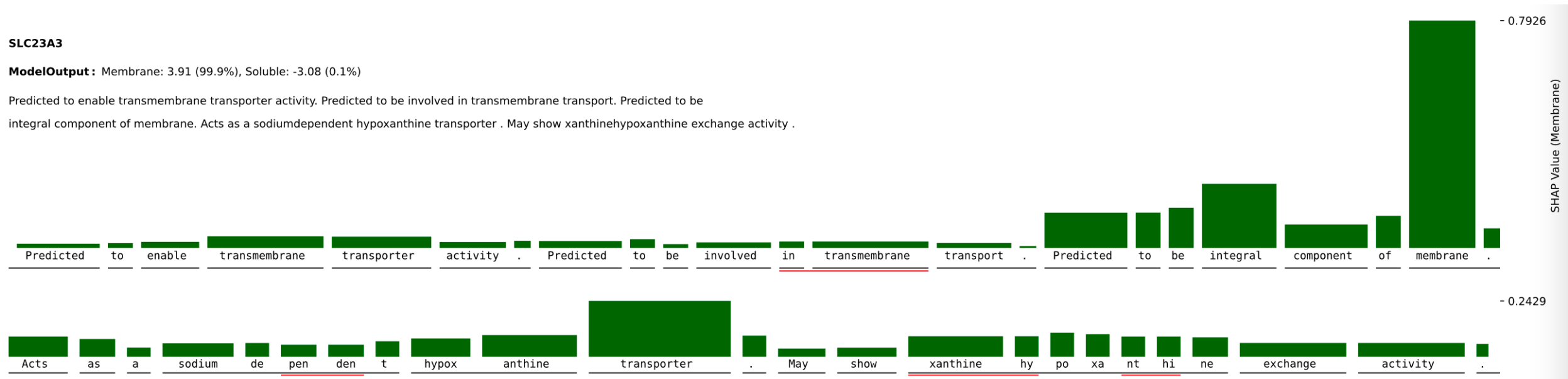


Membrane important word for predicting solubility: Interpreting Predictions with SHAP

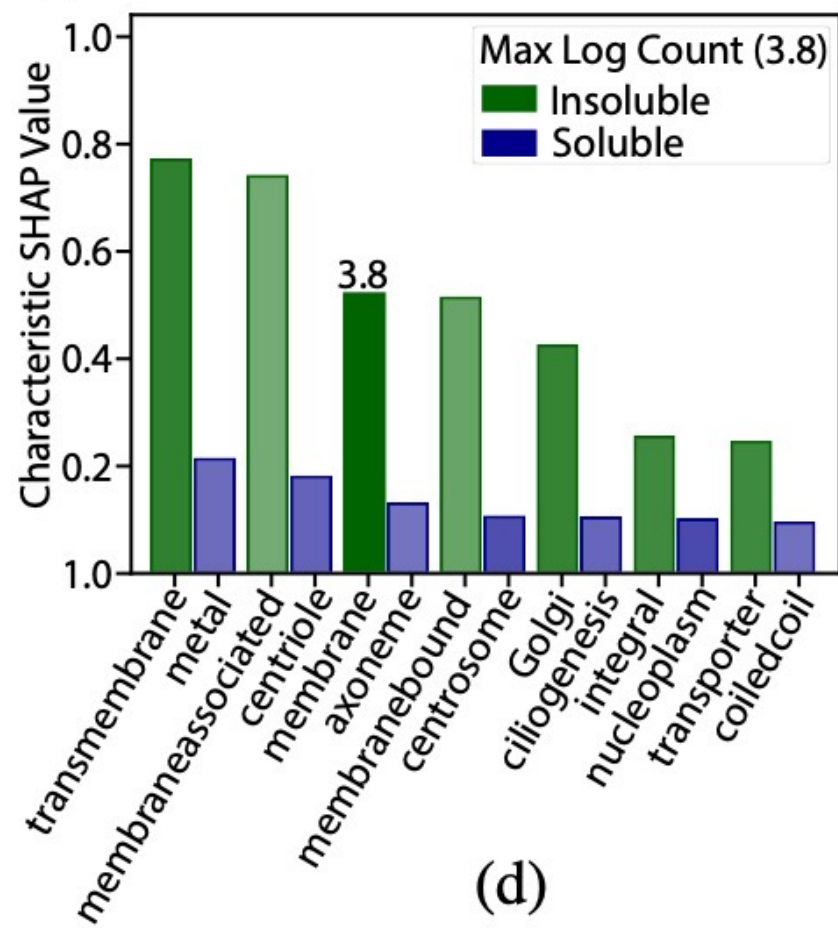
SLC23A3

ModelOutput: Membrane: 3.91 (99.9%), Soluble: -3.08 (0.1%)

Predicted to enable transmembrane transporter activity. Predicted to be involved in transmembrane transport. Predicted to be integral component of membrane. Acts as a sodiumdependent hypoxanthine transporter . May show xanthinehypoxanthine exchange activity .

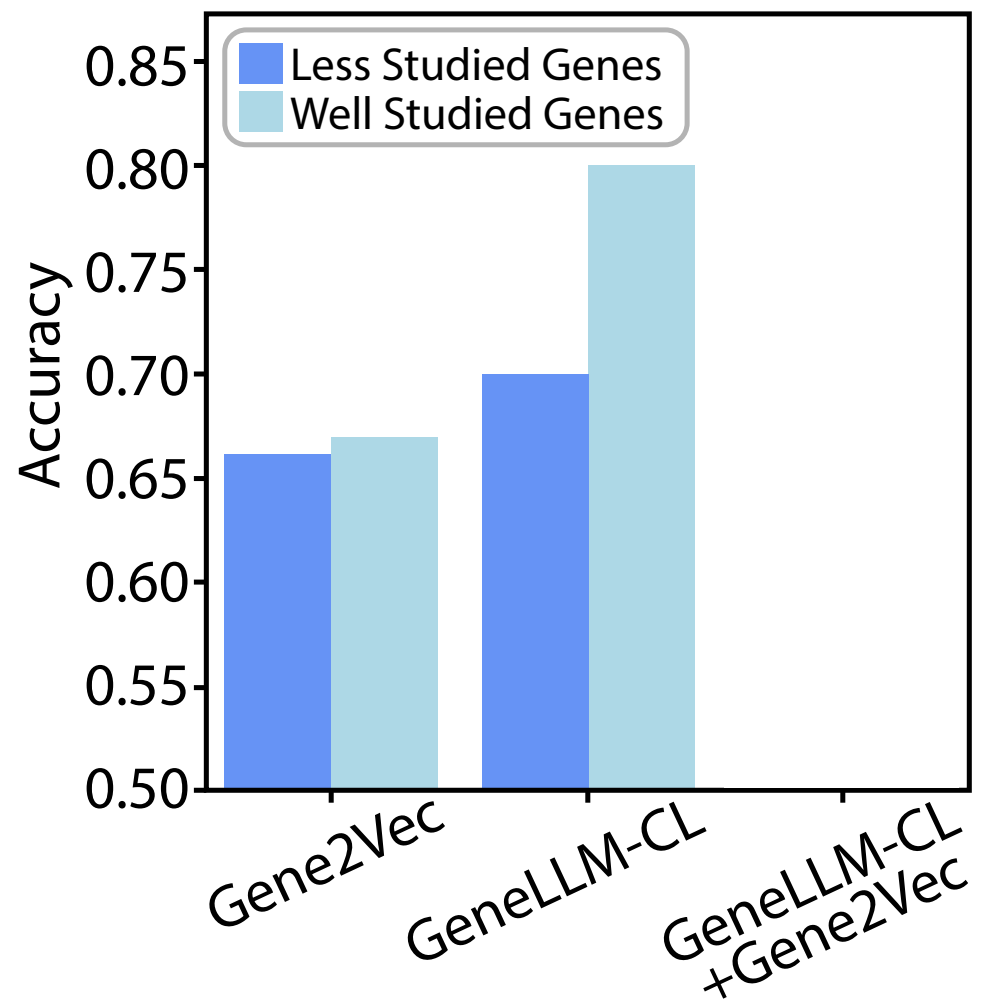


Important word for predicting solubility

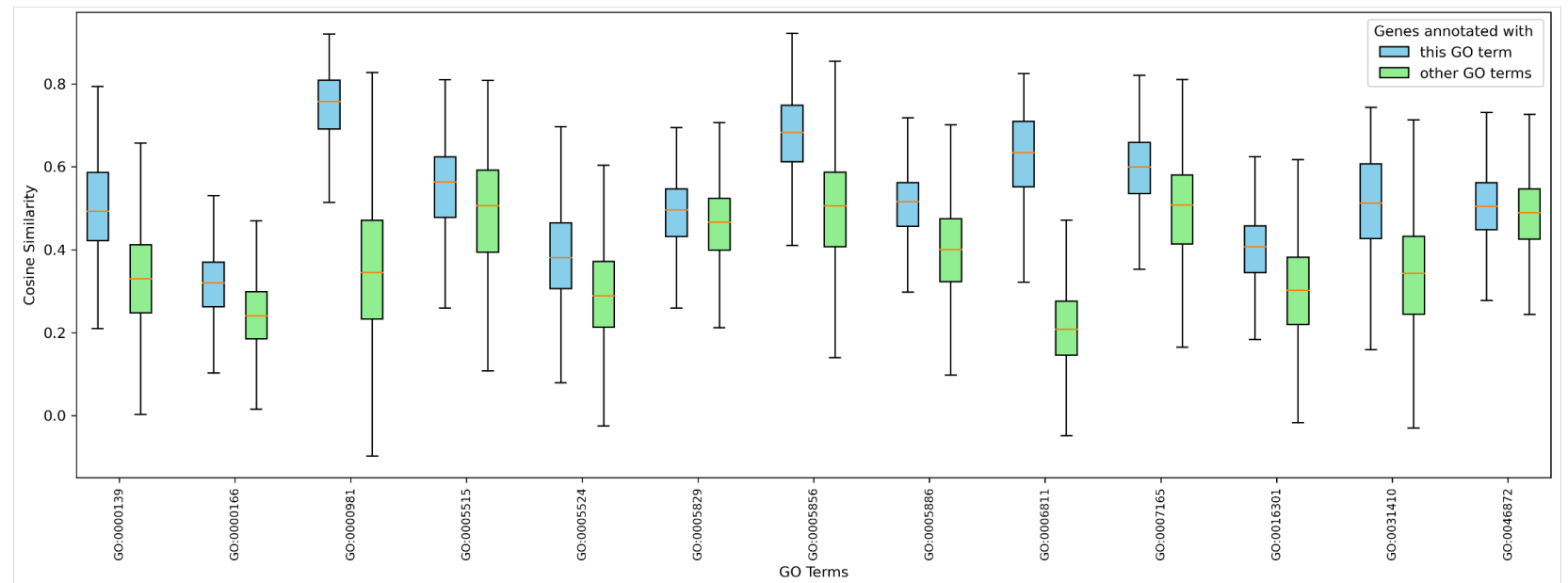
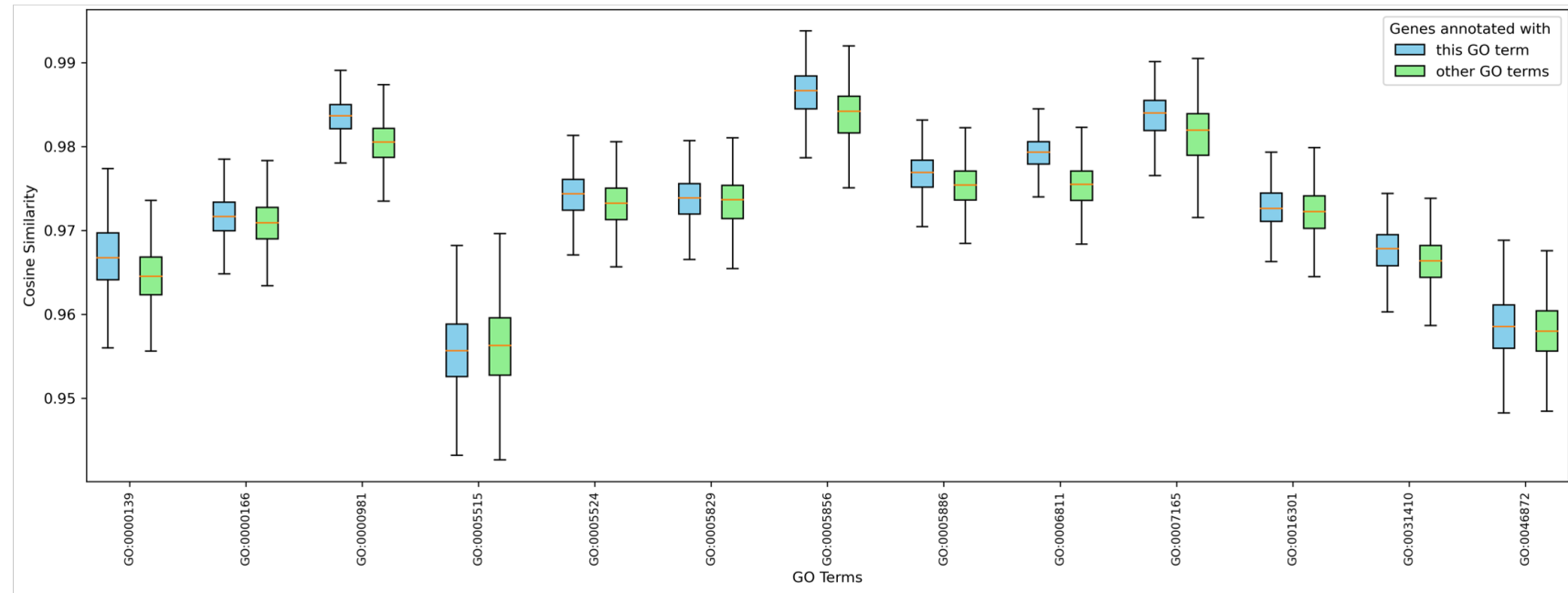


(d)

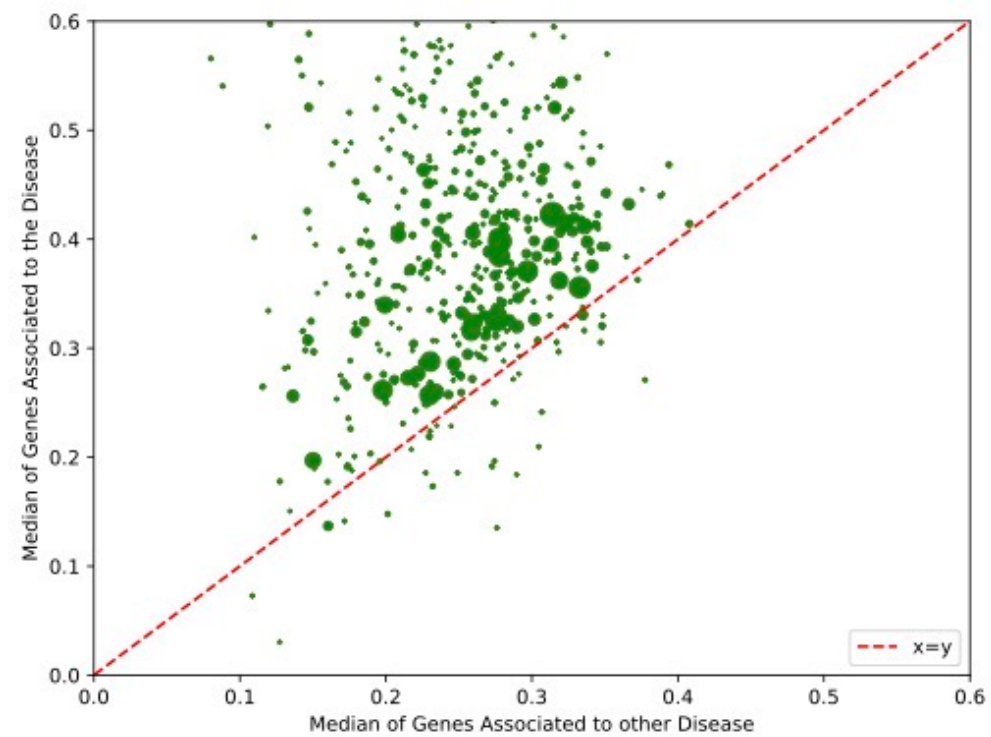
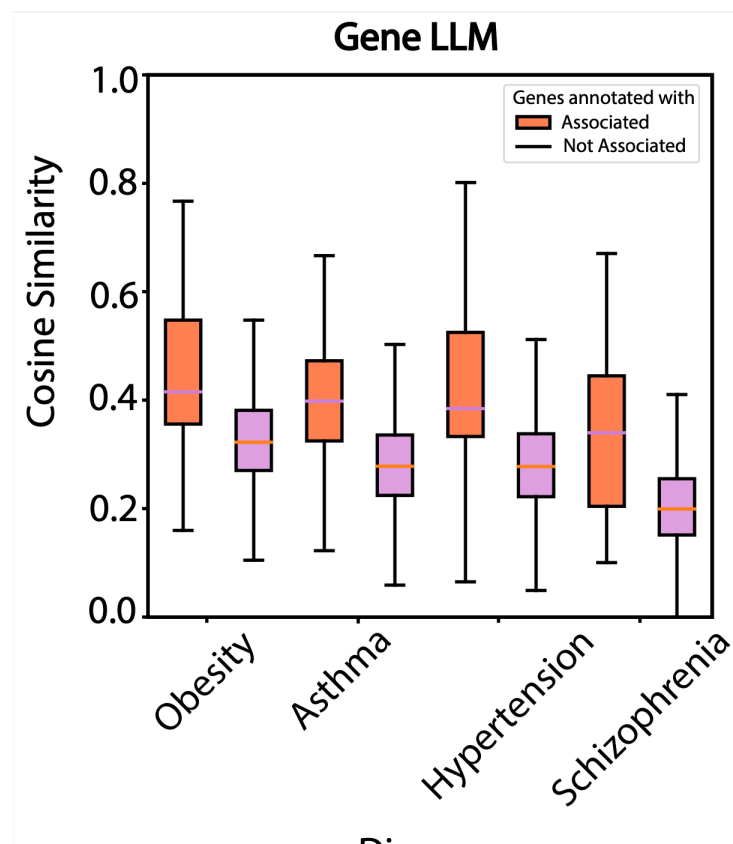
Structured data can
mitigate knowledge bias in
text



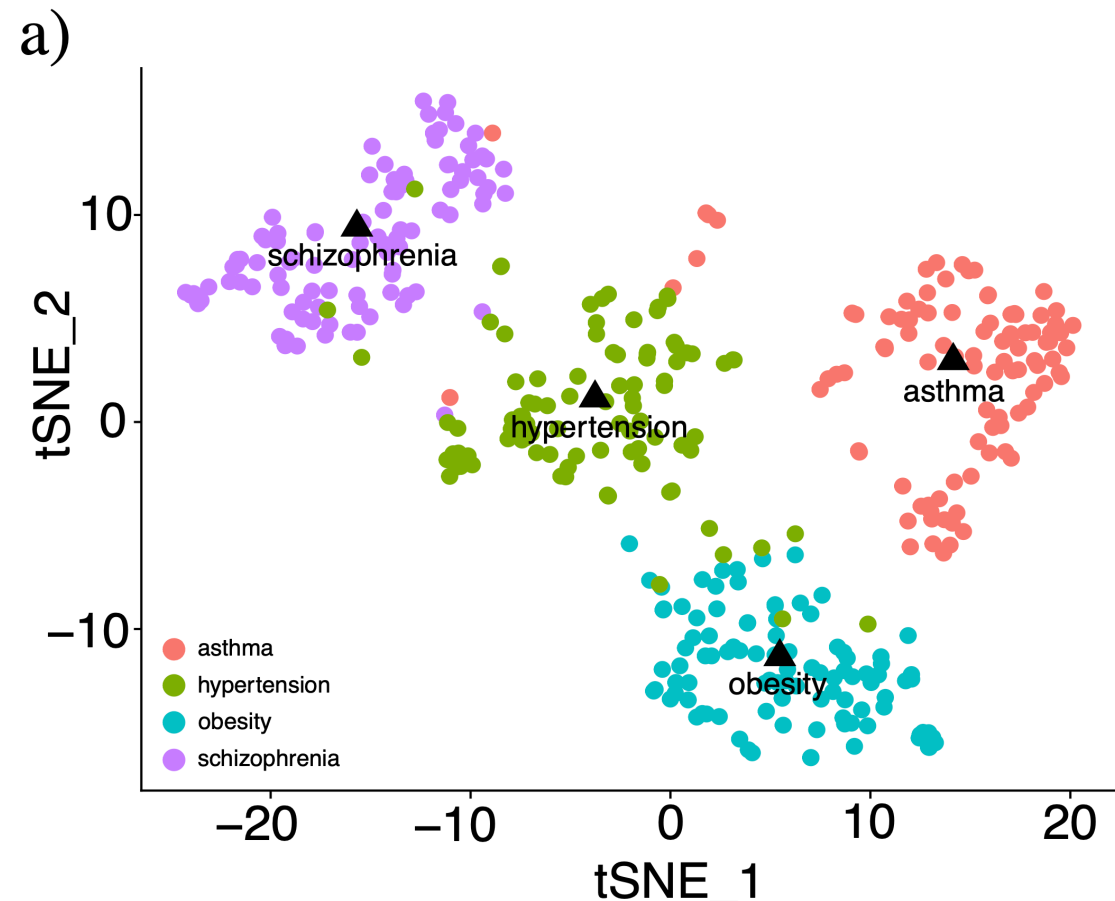
Contrastive learning
enables zero-shot
learnability



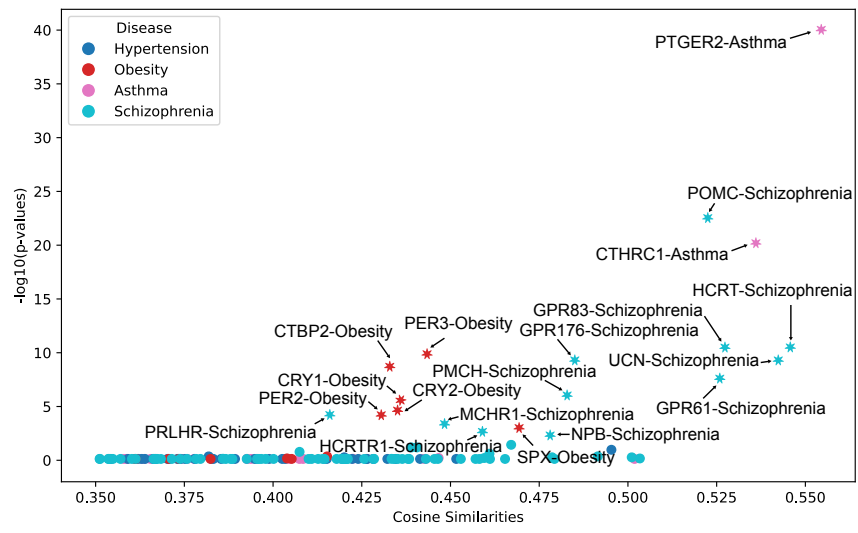
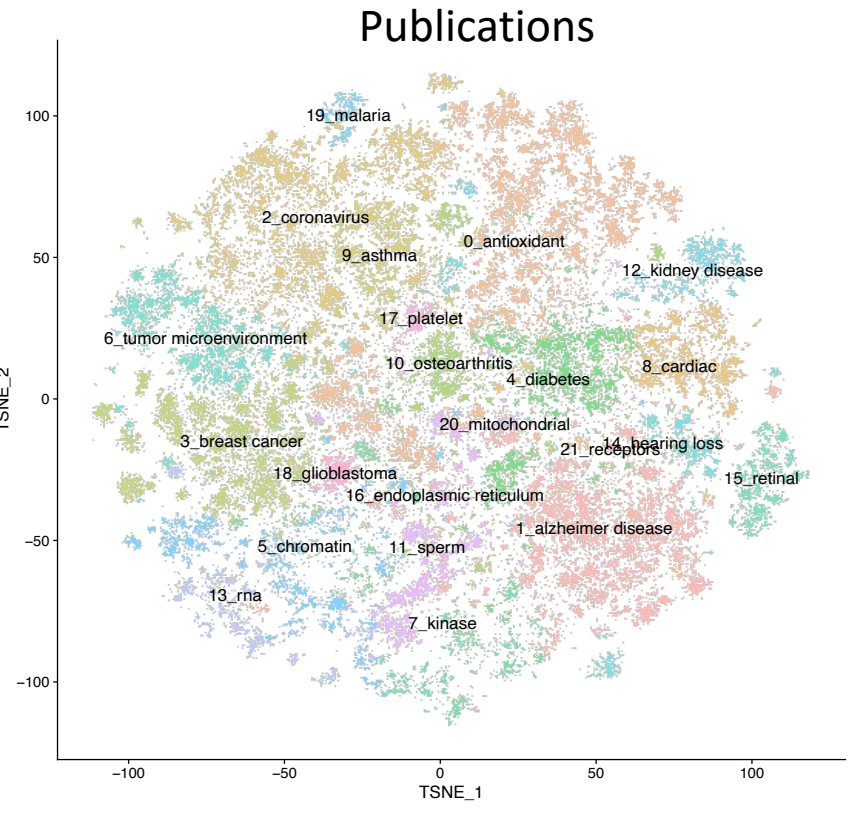
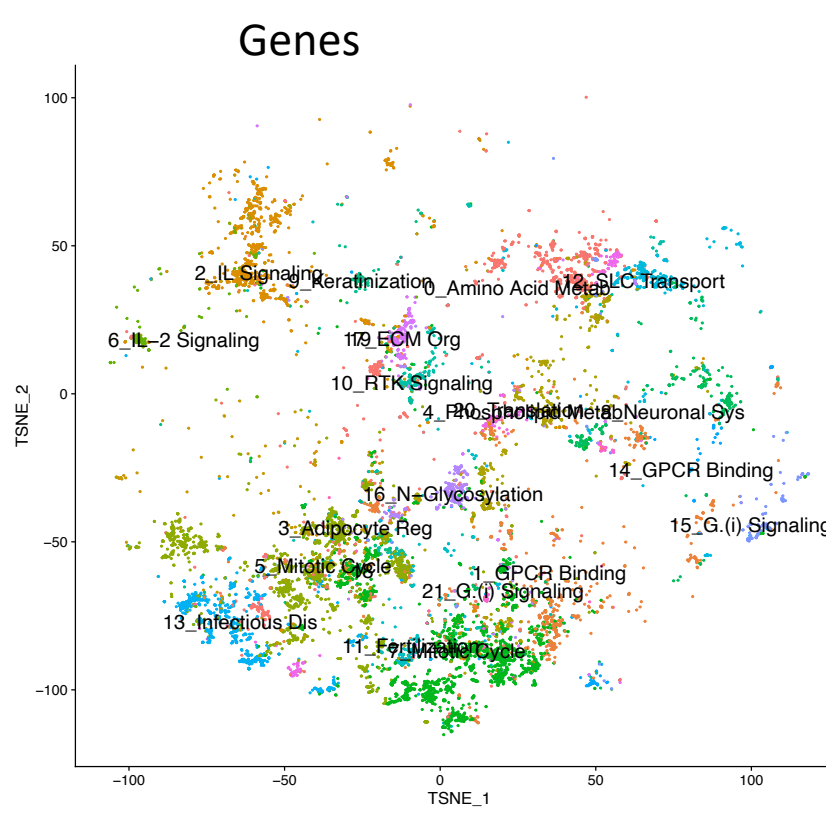
Zero-shot prediction of risk genes



Novel risk genes identification



Can GeneLLM ground disease-gene predictions in scientific literature?



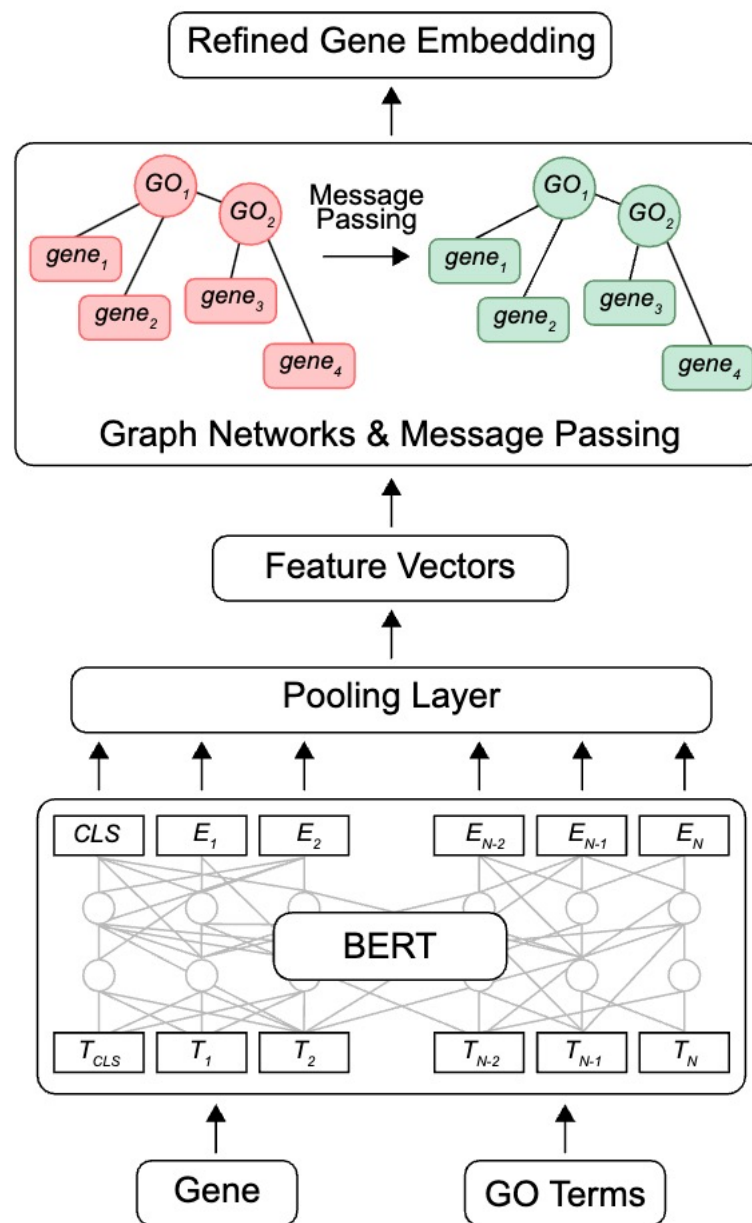
GeneLLM is a foundation model

Model	Dosage Sensitivity	BivalentVs Lys4 Methylated	BivalentVs Non Methylated	Tf range	Tf target type	Solubility	Subcellular localization	Conservation (Pearson Corr.)
Majority Classifier	0.73 ± —	0.58 ± —	0.75 ± —	0.73 ± —	0.41 ± —	0.52 ± —	0.39 ± —	— ± —
GPT2	0.74 ± 0.04	0.86 ± 0.04	0.80 ± 0.11	0.71 ± 0.03	0.18 ± 0.02	0.80 ± 0.02	0.77 ± 0.01	0.31 ± 0.02
Doc2Vec	0.74 ± 0.04	0.84 ± 0.06	0.78 ± 0.05	0.66 ± 0.07	0.26 ± 0.01	0.71 ± 0.03	0.69 ± 0.02	0.34 ± 0.01
PMC-LLaMA	0.86 ± 0.05	0.77 ± 0.04	0.84 ± 0.07	0.64 ± 0.08	0.08 ± 0.01	0.78 ± 0.03	0.69 ± 0.01	0.55 ± 0.01
XLNet	0.74 ± 0.06	0.84 ± 0.06	0.83 ± 0.08	0.69 ± 0.05	0.12 ± 0.01	0.79 ± 0.02	0.76 ± 0.01	0.40 ± 0.01
Gene2Vec	0.84 ± 0.04	0.84 ± 0.06	0.75 ± 0.06	0.75 ± 0.08	0.21 ± 0.01	0.56 ± 0.02	0.54 ± 0.02	0.50 ± 0.02
BERT-Base	0.76 ± 0.09	0.83 ± 0.06	0.77 ± 0.10	0.68 ± 0.04	0.17 ± 0.01	0.77 ± 0.02	0.76 ± 0.01	0.43 ± 0.01
GeneLLM	0.87 ± 0.06	0.86 ± 0.09	0.82 ± 0.08	0.74 ± 0.07	0.49 ± 0.04	0.89 ± 0.01	0.83 ± 0.01	0.53 ± 0.01

Genen outperforms GeneLLM

Model	Dosage Sensitivity	BivalentVs Lys4 Methylated	BivalentVs Non Methylated	Tf range	Tf target type	Solubility	Subcellular localization	Conservation (Pearson Corr.)
GPT2	0.83 ± 0.06	0.91 ± 0.06	0.83 ± 0.09	0.67 ± 0.13	0.52 ± 0.04	0.88 ± 0.01	0.91 ± 0.00	0.31 ± 0.02
Doc2Vec	0.78 ± 0.09	0.90 ± 0.08	0.79 ± 0.11	0.47 ± 0.09	0.54 ± 0.04	0.85 ± 0.02	0.85 ± 0.01	0.34 ± 0.01
PMC-LLaMA	0.89 ± 0.04	0.87 ± 0.03	0.90 ± 0.09	0.52 ± 0.45	0.06 ± 0.01	0.78 ± 0.03	0.83 ± 0.01	0.55 ± 0.01
XLNet	0.81 ± 0.08	0.90 ± 0.06	0.81 ± 0.06	0.61 ± 0.08	0.52 ± 0.01	0.86 ± 0.02	0.89 ± 0.01	0.40 ± 0.01
Gene2Vec	0.88 ± 0.04	0.88 ± 0.07	0.75 ± 0.06	0.56 ± 0.12	0.58 ± 0.01	0.60 ± 0.02	0.73 ± 0.01	0.50 ± 0.02
BERT-Base	0.85 ± 0.06	0.87 ± 0.05	0.85 ± 0.07	0.49 ± 0.11	0.53 ± 0.02	0.84 ± 0.02	0.90 ± 0.01	0.43 ± 0.01
GeneLLM	0.89 ± 0.06	0.87 ± 0.08	0.79 ± 0.10	0.47 ± 0.04	— ± —	0.89 ± 0.01	0.94 ± 0.03	0.53 ± 0.01
GCN	0.89 ± 0.03	0.92 ± 0.01	0.90 ± 0.03	0.59 ± 0.08	0.52 ± 0.03	0.99 ± 0.02	0.95 ± 0.00	0.55 ± 0.01

GNN are more effective in leveraging relational data

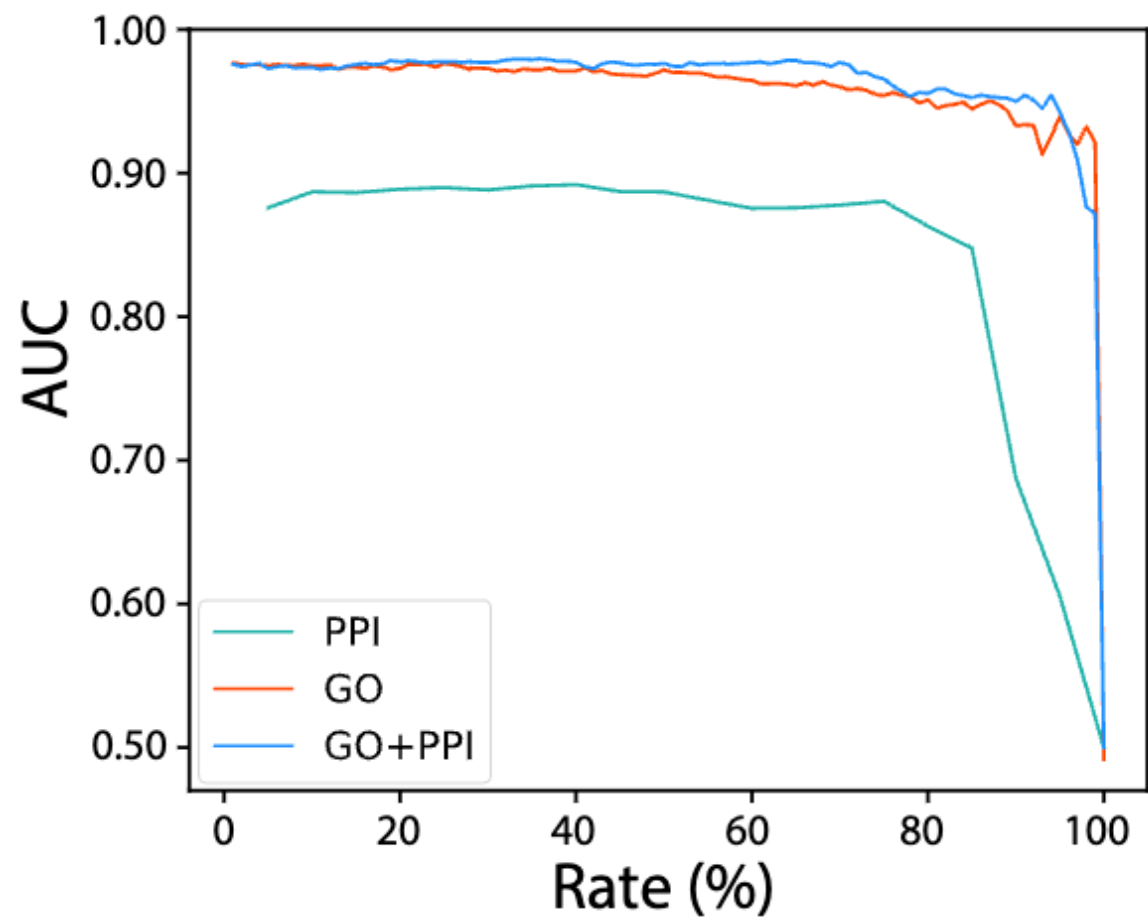


Yue Hu



Yanfu Zhang

Information propagation
infers gene information
with **75%** data hidden



GeneLLM website
litgene.tumorai.org



Gene-Disease-Drug Association

[Home](#)[About](#)[Feedback](#)[Contact Us](#)[Download](#)

Step 1: Create Your Prompt

Do you want to use a Gene, a Disease, or a Drug summary to create the prompt?

Gene

Disease

Drug

Summary

- Text data complements structured data complements and could mitigate knowledge bias
- Increasing interpretability and mitigating bias of AI models enhance their utility
- Scientific literature could be leveraged to increase AI reliability

Acknowledgements



Ala Jararweh
Macaulay Oladimeji
David Arradendo
Kushal Virupakshappa
Olufunmilola Oyebamiji
Luis Tafoya
Michael Servilla
Mikaela Dicome
Yue Hu



Mara Steinkamp
Marianne Berwick
Kim Leslie
Sarah Adams
Visu Palanimsamy
Eric Bartee
Olga Ponomarova



Yanfu Zhang



Patrick Finley

