

#### Functional Gene Embeddings (Group A)

Niklas Bühler, Keming Zhang, Marco Basmaji

Munich, February 2024





#### Outline

- 1. Motivation & Goals
- 2. Data Description
- 3. Methods
- 4. Results
- 5. Further Research



#### Motivation & Goals



#### **Motivation & Goals**

- GWAS: Investigate the effect of common variants on traits and diseases
- Functional Gene Embeddings: Numerical vectors capturing gene functions
- **Goal**: Generate useful representations of genes for downstream tasks
- Evaluate quality of embeddings on GWAS signal prediction





- 1. GTEx
- 2. Tabula Sapiens
- 3. DepMap Gene Effect
- 4. RNA Isoform



#### GTEx, The Genotype-Tissue Expression:

• Collected samples from non-diseased tissue across many individuals (50k x 11k dense matrix)

#### Tabula sapiens:

• Contains single cell transcriptomics data of 483,152 cells across 58,870 genes.

1. G. Consortium, "The gtex consortium atlas of genetic regulatory effects across human tissues," Science, vol. 369, no. 6509, pp. 1318–1330, 2020. 2. T. S. Consortium\*, R. C. Jones, J. Karkanias, M. A. Krasnow, A. O. Pisco, S. R. Quake, J. Salzman, N. Yosef, B. Bulthaup, P. Brown, et al., "The tabula sapiens: A multiple-organ, single-cell transcriptomic atlas of humans," Science , vol. 376, no. 6594, p. eabl4896, 2022.



DepMap Gene Effect dataset:

- Scores that measure the effect size of knocking out a gene.
- Measure the dependency between genes and cell lines of around 17k genes and 500 cell lines.

#### **RNA** Isoform Dataset:

• It contains RNA levels in 32 tissues, based on RNA-seq.

2. M.Uhl en,L.Fagerberg,B.M.Hallstr om, C.Lindskog, Oksvold,A. Mardinoglu, A. Sivertsson, C. Kampf, E. Sj ostedt, A. Asplund, et al.,

<sup>1.</sup> Tsherniak, F. Vazquez, P. G. Montgomery, B. A. Weir, G. Kryukov, G. S. Cowley, S. Gill, W. F. Harrington, S. Pantel, J. M. Krill-Burger, et al. "Defining a cancer dependency map," Cell vol. 170, no. 3, pp. 564–576, 2017.

<sup>&</sup>quot;Tissue-based map of the human proteome," Science vol. 347, no. 6220, p. 1260419, 2015.



#### Data Approach:

• All of the datasets were aggregated based on unique Ensembl gene IDs and subsetted according to the protein-coding genes (Around 19k genes).

• GTEx was cut over samples and the Autoencoder was trained over only 5k samples.

• Tabula Sapiens was aggregated by cell ontology classes (Pseudo bulks).

#### Methods

- 1. Generating functional gene embeddings
  - a. Principal Components
  - b. Autoencoder
  - c. Variational Deep Tensor Factorization Model
- 2. Evaluating the utility of functional gene embeddings

# Method – Generating Embeddings – Principal Components

• Principal components that could explain 80-95 percent of variance





#### Method – Generating Embeddings – Autoencoder

• Basic structure of an autoencoder



• Latent Space: the compressed representation is extracted and serve as information-dense embeddings of the genes



## Method – Generating Embeddings – Variational Deep Tensor Factorization Model

• Basic structure of Variational Deep Tensor Factorization Model<sup>1</sup>



• Training the gene and sample embeddings to reconstruct the data matrix

1.F. Brechtmann, T. Bechtler, S. Londhe, C. Mertes, and J. Gagneur, "Evaluation of input data modality choices on functional gene embeddings," NAR Genomics and Bioinformatics, vol. 5, no. 4, p. Iqad095, 2023.



# Method – Evaluating the Embeddings – GWAS Signal Prediction

- Evaluate the utility of gene embeddings by using them as gene features to predict genome-wide association study (GWAS) summary statistics, and investigating how they could improve the accuracy scores.
- General process of the evaluation on **30 traits**.



## Method – Evaluating the Embedding – GWAS Signal **TI** Prediction – Details

- Prediction model: linear regression and XGBoost regression
- Traits considered: 30 traits, mostly maximally independent (see the Supplement)
- Covariates: gene density, effective gene size, inverse of the mean minor allele count of SNPs in the genes, as well as logarithmic values of them
- Cross validation strategy: Leave-One-Chromosome-Out
- GWAS data sources:
  - The Pan UKB study<sup>1</sup>
  - The UKBiobank<sup>2</sup>
  - The publication results of Brechtmann et al<sup>3</sup>

<sup>1.</sup>Pan-UKB team, "Pan-uk biobank," 2020.

<sup>2.</sup>C. Bycroft, C. Freeman, D. Petkova, G. Band, L. T. Elliott, K. Sharp, A. Motyer, D. Vukcevic, O. Delaneau, J.O'Connell, et al., "The uk biobank resource with deep phenotyping and genomic data," Nature, vol. 562, no. 7726, pp. 203209, 2018.

<sup>3.</sup>F. Brechtmann, T. Bechtler, S. Londhe, C. Mertes, and J. Gagneur, "Evaluation of input data modality choices on functional gene embeddings," NAR Genomics and Bioinformatics, vol. 5, no. 4, p. Igad095, 2023.



#### Results

- 1. Embedding Results
- 2. Significant Differences between Full and Null Models
- 3. Directly Examining  $\Delta R^2$
- 4. Well-Predictable Traits Improve More with Additional Data
- 5. Which Embeddings work Well on which Traits?
- 6. Evaluating with Linear Regression



#### **Embedding Results**

- Eleven embeddings from three architectures and four data sources.
- **XGBoost** Regression Model: Better general results (R<sup>2</sup>) and more interesting analyses

Embedding Source	Embedding-Generating Model or Method	Dimension
DepMap Dataset	Principal Components	81
DepMap Dataset	Autoencoder	128
DepMap Dataset	VDTFM	128
GTEx Dataset	Principal Components	64
GTEx Dataset	Autoencoder	128
GTEx Dataset	Autoencoder	256
RNA Isoform Dataset	Principal Components	4
RNA Isoform Dataset	Autoencoder	32
Tabula Sapiens Dataset	Principal Components	5
Tabula Sapiens Dataset	Autoencoder	64
Tabula Sapiens Dataset	VDTFM	256

Table: Eleven embeddings from three different architectures and four different data sources.

## Significant Differences between Full and Null Models

Comparatively strong improvements with **Omics** and **GTEx** embeddings.

Wilcoxon Signed-Rank Test:

- Significant differences (α = 0.05)
- Bonferroni Correction (α ≈ 0.0042): still significant



#### Directly Examining $\Delta R^2$

- Omics and GTEx strongest
- RNA and TS embeddings also increased R<sup>2</sup> for most traits
- DepMap and VDTFM embeddings of TS didn't improve R<sup>2</sup> on most traits



# Well-Predictable Traits Improve More with Additional Data



# Well-Predictable Traits Improve More with Additional Data



		ΔR <sup>2</sup>									TIM					
۱۸/I-	¦ala ⊑uas la astaliu		Alanine_aminotransferase	0.0010	0.0016	-0.0009	0.0157	0.0170	0.0152	0.0155	-0.0015	0.0053	0.0047	0.0088	0.0116	
vvn	ich Embeddir	מו	S Albumin	0.0032	0.0053	0.0045	0.0184	0.0189	0.0186	0.0206	-0.0045	0.0089	0.0077	0.0127	0.0124	0.025
		3	Apolipoprotein_A	-0.0074	-0.0029	-0.0035	0.0113	0.0099	0.0115	0.0119	-0.0105	0.0039	0.0038	0.0073	0.0072	- 0.025
wor	·k \N/all on \N/h	nic	C-reactive_protein	-0.0104	-0.0082	-0.0068	0.0098	0.0097		0.0112	-0.0055	0.0031	0.0042	0.0063	0.0067	
<b>WVUI</b>			Calcium_100024	0.0002	0.0011	0.0009	-0.0001	-0.0002	0.0002	-0.0000	-0.0002	0.0000	0.0002	-0.0000	0.0003	0.020
<b>T</b>	:1 - <b>O</b>		Calcium_30680	-0.0034	-0.0042	0.0013	0.0103	0.0095	0.0113	0.0106	-0.0015	0.0075	0.0067	0.0088	0.0091	- 0.020
Ira	ITS ?		Cholesterol	-0.0081	-0.0045	-0.0068	0.0059	0.0075	0.0088	0.0101	-0.0096	-0.0003	0.0018	0.0033	0.0045	
			Creatinine Direct bilirubin	-0.0099	-0.0041	-0.0048	0.0160	0.0155	0.0142	0.0243	-0.0068	-0.0033	0.0033	0.0086	-0.0072	0.015
			Glucose	-0.0041	-0.0008	-0.0014	0.0086	0.0083	0.0120	0.0090	-0.0024	0.0056	0.0071	0.0069	0.0078	- 0.015
•	VDTEM didn't halp		HDL	-0.0087	-0.0095	-0.0058	0.0116	0.0097	0.0129	0.0107	-0.0080	0.0046	0.0040	0.0076	0.0078	
•			HDL_cholesterol	-0.0086	-0.0085	-0.0049	0.0119	0.0094	0.0134	0.0084	-0.0078	0.0046	0.0041	0.0075	0.0079	
			IGF-1	-0.0083	-0.0063	-0.0026	0.0129	0.0126	0.0129	0.0105	-0.0027	0.0053	0.0059	0.0061	0.0084	- 0.010
ë	across datasets	aj	LDLC direct	-0.0077	-0.0023	-0.0032	0.0123	0.0160	0.0137	0.0192	-0.0061	-0.0014	-0.0010	0.0055	0.0105	
		Ē	LDL_direct_adjusted_by_medication	-0.0082	-0.0030	-0.0041	0.0119	0.0144	0.0128	0.0187	-0.0062	0.0031	0.0037	0.0053	0.0104	
•	Some traits didn't		Lipoprotein_A	-0.0040	-0.0018	-0.0010	-0.0018	-0.0017	-0.0013	-0.0025	-0.0038	-0.0007	-0.0005	-0.0005	-0.0010	- 0.005
			MCH	-0.0059	-0.0075	-0.0020	0.0136	0.0112	0.0100	0.0165	-0.0078	0.0128	0.0043	0.0007	0.0049	
	i <b>mprove</b> (e.g.		Mean_corpuscular_haemoglobin	0.0053	0.0024	0.0052	0.0202	0.0189	0.0139	0.0199	-0.0000	0.0161	0.0101	0.0079	0.0119	
			RBC	-0.0034	-0.0044	-0.0026	0.0049	0.0037	0.0008	0.0038	-0.0038	0.0002	0.0004	0.0080	0.0038	- 0.000
	Lipoprotein A)		Red_blood_cell_erythrocyte_count	-0.0021	-0.0041	0.0022	0.0203	0.0188	0.0177	0.0257	-0.0029	0.0141	0.0102	0.0123	0.0167	
	1 1 ,		SHBG	0.0002	-0.0010	0.0042	0.0195	0.0258	0.0232	0.0281	-0.0059	0.0074	0.0077	0.0117	0.0151	
•	Others greatly		Testosterone	-0.0014	-0.0017	-0.0020	0.0143	0.0135	0.0129	0.0159	-0.0051	0.0003	0.0030	0.0042	0.0049	0.005
_	9.00.1		Total_bilirubin	-0.0093	-0.0110	-0.0092	0.0002	0.0032	-0.0007	0.0001	-0.0043	0.0012	0.0015	0.0008	0.0014	
i	improved through		Triglycerides	-0.0121	-0.0092	-0.0128	0.0079	0.0112	0.0096	0.0040	-0.0115	-0.0002	-0.0002	0.0036	0.0031	
	improved through		Urate	-0.0026	-0.0016	-0.0019	0.0068	0.0073	0.0052	0.0085	-0.0045	0.0037	0.0045	0.0076	0.0072	0.010
	GTEx and Omics (A d		Vitamin_D	0.0025	0.0002	0.0005	0.0065	0.0046	0.0042	0.0079	-0.0021	0.0008	0.0021	0.0031	0.0030	
				8	в	tiance di28	4,728	× 250	A 664	8258	825	, ait	ance deA	all	ince edge	
:	SHBG)		¢	phap Duff M.	PC5-8810.92V	Dephap ae	GTEX AL	offet At	GTEXPO	Omics	S. P. DVITTON	PC5-65-0-848	15 J <sup>b a</sup>	PC5-04-0.9548	118.2"	
		Embedding														

21



#### **Evaluating with Linear Regression**

- Almost all embeddings perform worse
- Exception: DepMap, but still only slight increase
- Omics robust against regression model selection
- **RNA Autoencoder** also **robust** in both setups, but not as good
- Hypothesis: Assumptions of linear regression model might not hold





#### **Further Research**

- 1. Comparable Dimensions of Embeddings
- 2. Embeddings Generated from VDTFM
- 3. Combining Embeddings from Different Data Sources

## ПΠ

#### References

[1] F. Brechtmann, T. Bechtler, S. Londhe, C. Mertes, and J. Gagneur, "Evaluation of input data modality choices on functional gene embeddings," NAR Genomics and Bioinformatics, vol. 5, no. 4, p. lqad095, 2023.

[2] A. Tsherniak, F. Vazquez, P. G. Montgomery, B. A. Weir, G. Kryukov, G. S. Cowley, S. Gill, W. F. Harrington, S. Pantel, J. M. Krill-Burger, et al., "Defining a cancer dependency map," Cell, vol. 170, no. 3, pp. 564–576, 2017.

[3] G. Consortium, "The gtex consortium atlas of genetic regulatory effects across human tissues," Science, vol. 369, no. 6509, pp. 1318–1330, 2020.

[4] M. Uhlén, L. Fagerberg, B. M. Hallström, C. Lindskog, P. Oksvold, A. Mardinoglu, Å. Sivertsson, C. Kampf, E. Sjöstedt, A. Asplund, et al.,

"Tissue-based map of the human proteome," Science, vol. 347, no. 6220, p. 1260419, 2015.

[5] T. T. S. Consortium\*, R. C. Jones, J. Karkanias, M. A. Krasnow, A. O. Pisco, S. R. Quake, J. Salzman, N. Yosef, B. Bulthaup, P. Brown, et al., "The tabula sapiens: A multiple-organ, single-cell transcriptomic atlas of humans," Science,

vol. 376, no. 6594, p. eabl4896, 2022.

[6] M. Wainberg, R. A. Kamber, A. Balsubramani, R. M. Meyers, N. Sinnott-Armstrong, D. Hornburg, L. Jiang, J. Chan, R. Jian, M. Gu, et al., "A genomewide atlas of co-essential modules assigns function to uncharacterized genes," Nature genetics, vol. 53, no. 5, pp. 638–649, 2021.

[7] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2017.

[8] L. Biewald, "Experiment tracking with weights and biases," 2020. Software available from wandb.com.

[9] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2019.

[10] Pan-UKB team, "Pan-uk biobank," 2020.

[11] C. Bycroft, C. Freeman, D. Petkova, G. Band, L. T. Elliott, K. Sharp, A. Motyer, D. Vukcevic, O. Delaneau, J. O'Connell, et al., "The uk biobank resource with deep phenotyping and genomic data," Nature, vol. 562, no. 7726, pp. 203–209, 2018.

[12] C. A. de Leeuw, J. M. Mooij, T. Heskes, and D. Posthuma, "Magma: generalized gene-set analysis of gwas data," PLoS computational biology, vol. 11, no. 4, p. e1004219, 2015.

[13] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, ACM, Aug. 2016.



Thank You for the Attention...



• Traits selected in GWAS Signal Prediction

Trait	Max. ind.	Trait	Max. ind.	Trait	Max. ind.
	set		set		set
Alanine Amino-	True	Albumin	True	Apolipoprotein A	True
transferate					
Apoliprotein B	True	C-reactive Protein	True	Calcium 30680	True
Calcium 100024	True	Cholesterol	True	Creatinine	True
Direct Bilirubin	False	Glucose	Flase	HDL Cholesterol	True
HDL	False	IGF-1	True	LDL Direct Adjusted	Flase
				by Medication	
LDL Direct	True	LDLC	False	Lipoprotein	False
MCH False		Mean Corpuscular	False	Phosphate	True
		Haemoglobin			
RBC	ABC False		False	SHBG	True
		throcyte Count			
Testosterone	Flase	Total Bilirubin	True	Total Protein	True
Triglycerides	True	Urate	False	Vitamin D	True

Full Models vs. Covariate Model (α=0.05)

## Supplement



ТШ



ТШ



Delta R <sup>2</sup>								
Alanine_aminotransferase	0.0140	0.0149						
Albumi	0.0163	0.0222	- 0.025					
Apolipoprotein_/	0.0105	0.0114	0.023					
Apolipoprotein_E	0.0097	0.0123						
C-reactive_protein	0.0053	0.0094						
Calcium_100024	0.0000	0.0000						
Calcium_3068	0.0108	0.0107						
Cholesterc	0.0064	0.0120	- 0.020					
Creatinin	0.0112	0.0211						
Direct_bilirubir	-0.0017	0.0025						
Glucos	0.0099	0.0085						
HDI	0.0091	0.0086						
HDL_cholestero	0.0088	0.0085	- 0.015					
IGF-	0.0076	0.0115						
志 LDLC	0.0140	0.0183						
⊨ LDL_direc	0.0065	0.0101						
LDL_direct_adjusted_by_medication	0.0139	0.0206						
Lipoprotein_/	-0.0008	-0.0015	- 0.010					
MCH	0.0153	0.0148						
Mean_corpuscular_haemoglobin	0.0200	0.0211						
Phosphate	0.0031	0.0064						
RBO	0.0160	0.0226						
Red_blood_cell_erythrocyte_coun	0.0179	0.0220	- 0.005					
SHBO	0.0184	0.0269	0.000					
Testosteron	0.0067	0.0129						
Total_bilirubi	-0.0008	0.0008						
Total_protein	0.0214	0.0247						
Triglyceride	0.0055	0.0051						
Urati	0.0078	0.0106	- 0.000					
Vitamin_E	0.0055	0.0098						
	State of the State	Seres - Elle						
	Mo	del						



