# Systematic tissue annotations of genomic samples by modeling unstructured metadata

Nathaniel T. Hawkins[1], Marc Maldaver[1], Arjun Krishnan[1,2]

[1]Dept. of Computational Mathematics, Science and Engineering. Michigan State Univeristy. East Lansing, MI, United States.
[2]Dept. of Biochemistry and Molecular Biology. Michigan State University. East Lansing, MI, United States.
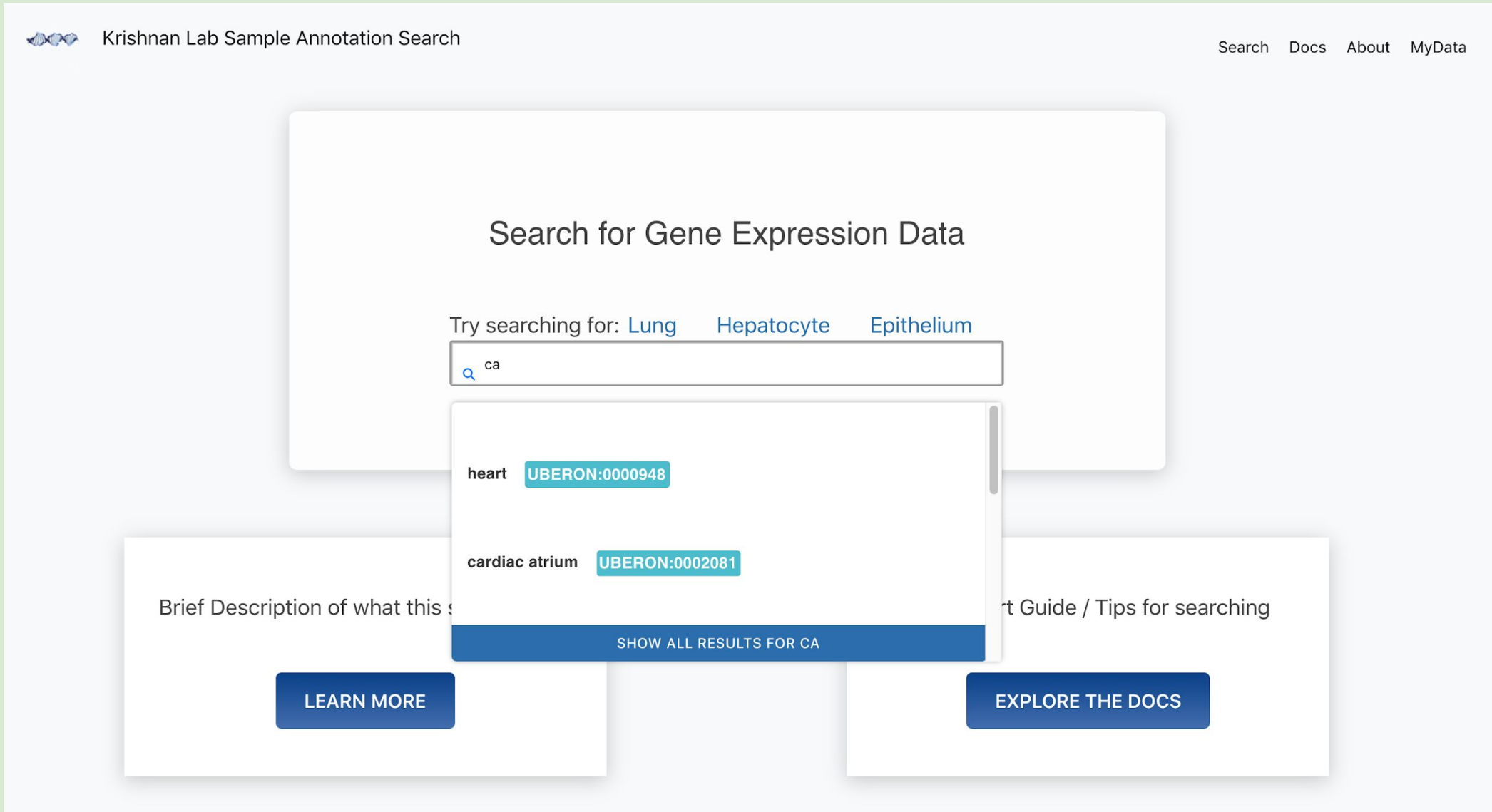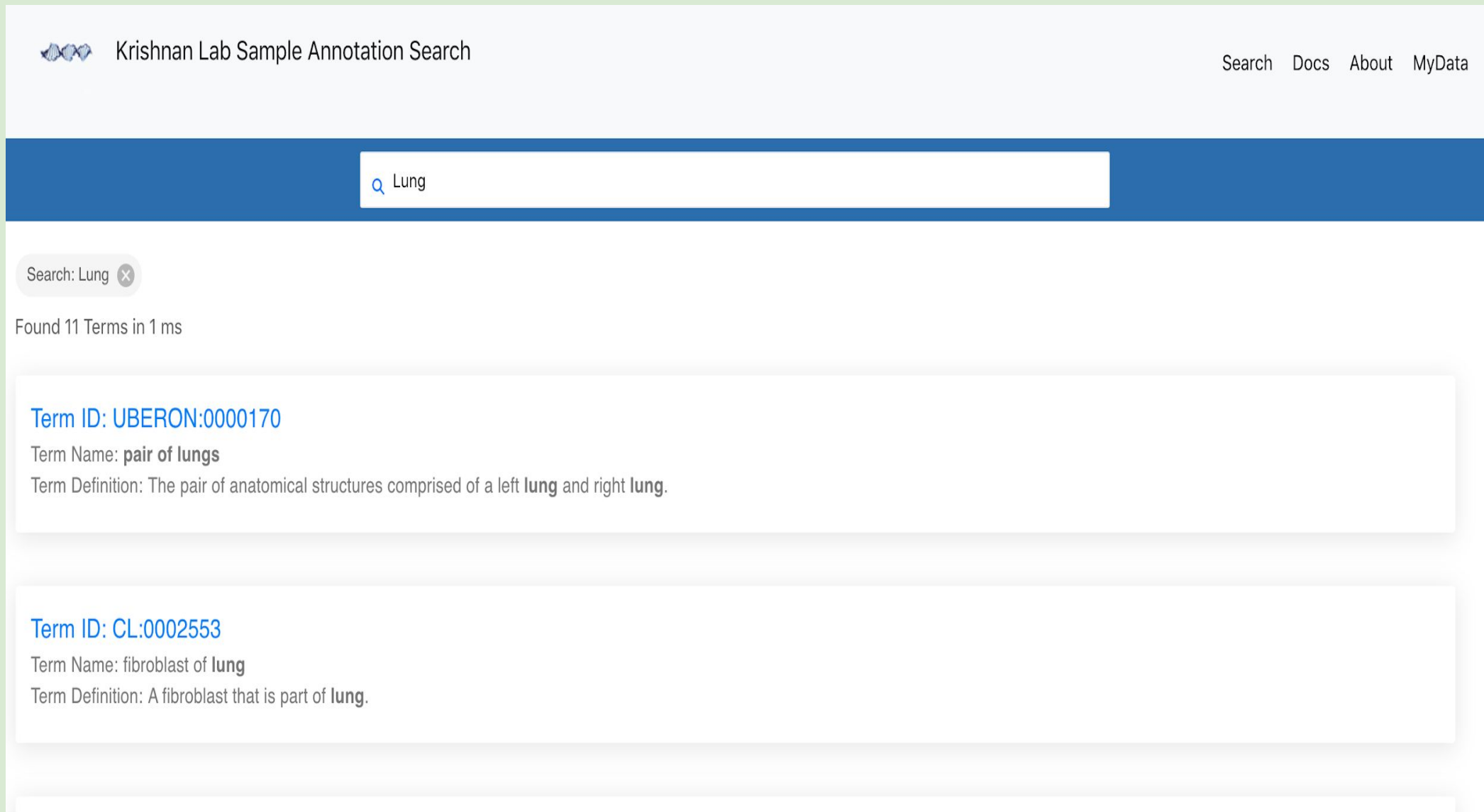
## Overview

The reuse of publicly available omics data is limited by the lack of systematic annotations for tissues or cell types of origin. Previous work has used the underlying omics data to predict annotations for a given sample, but these tools are not widely usable due to the need to retrain for different data types. We propose a novel approach that leverages free-text, natural-language metadata to identify the tissue or cell-type of origin of omics samples. Using sample metadata files from ArrayExpress, we create a numerical representation for each sample using flair, a python NLP library. These representations are used train OVR logistic regression models, one for each tissue or cell type we have labels for, to predict a sample's annotations given its metadata. We compare our models' predictions to annotations predicted by models trained on gene expression data and those from two existing tools: TAGGER, an NER tool, and MetaSRA, a method that uses prior knowledge and ontologies.
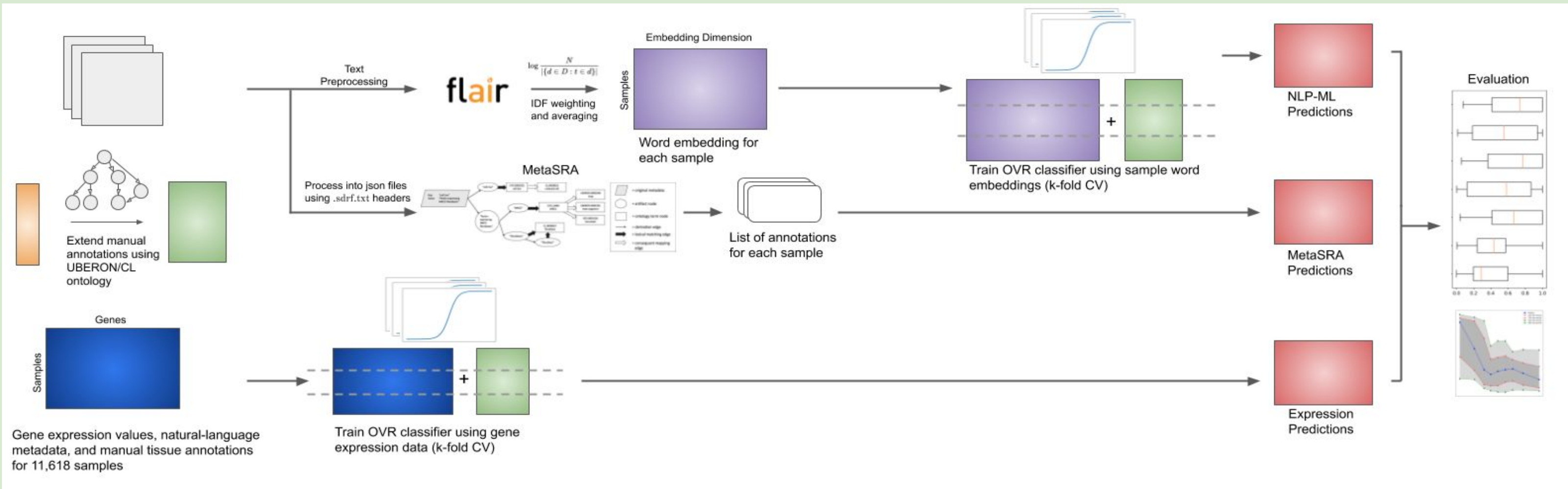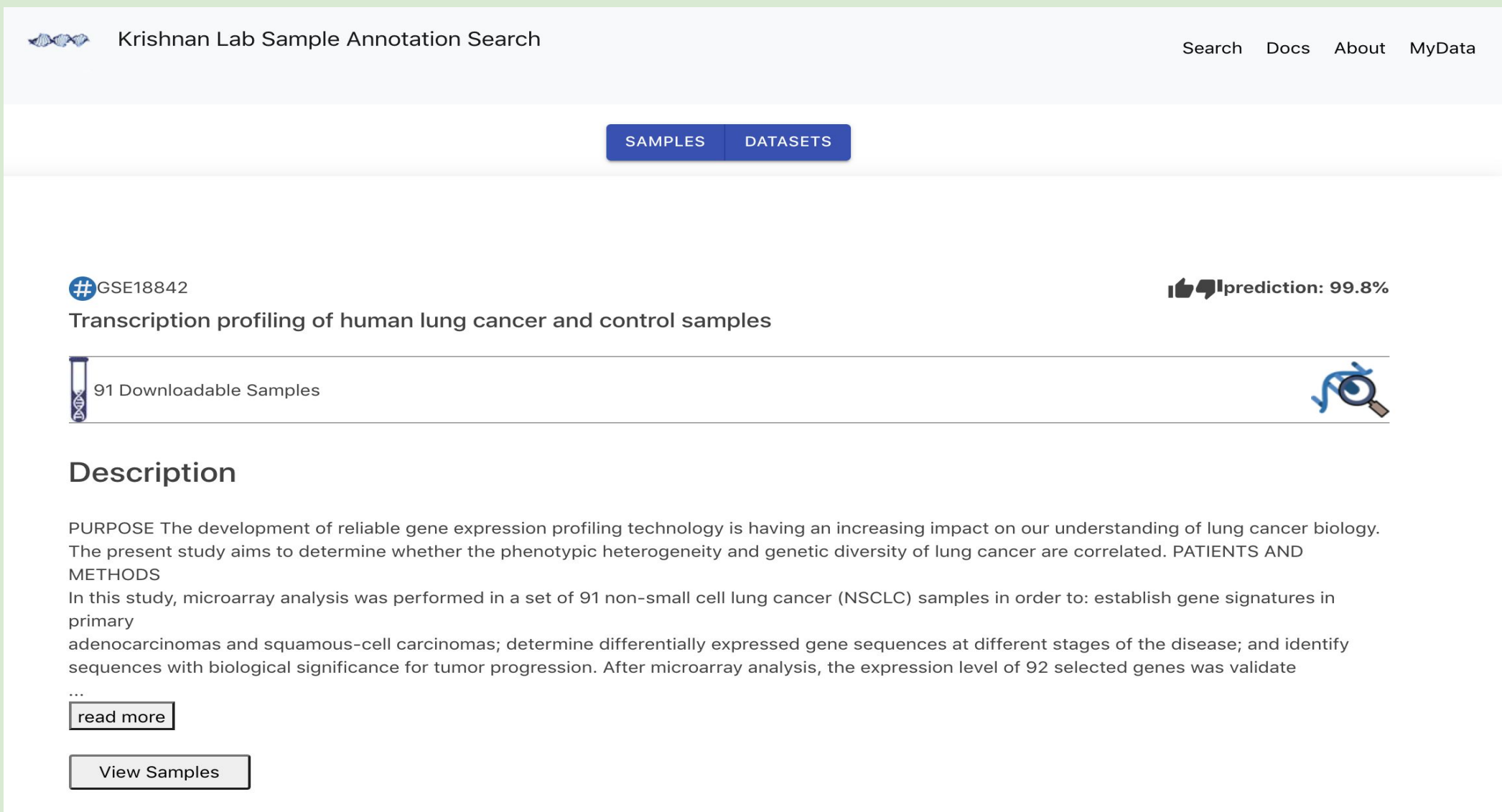
## WEB SERVER DEVELOPMENT

Example user query. Search text matched to available models using Elasticsearch



Or user can search a term and select the most appropriate tissue from list of available models



Query results. Series listed by predicted probability from NLP-ML models along with associated metadata. Individual samples in series can be selected



## PERFORMANCE COMPARISON

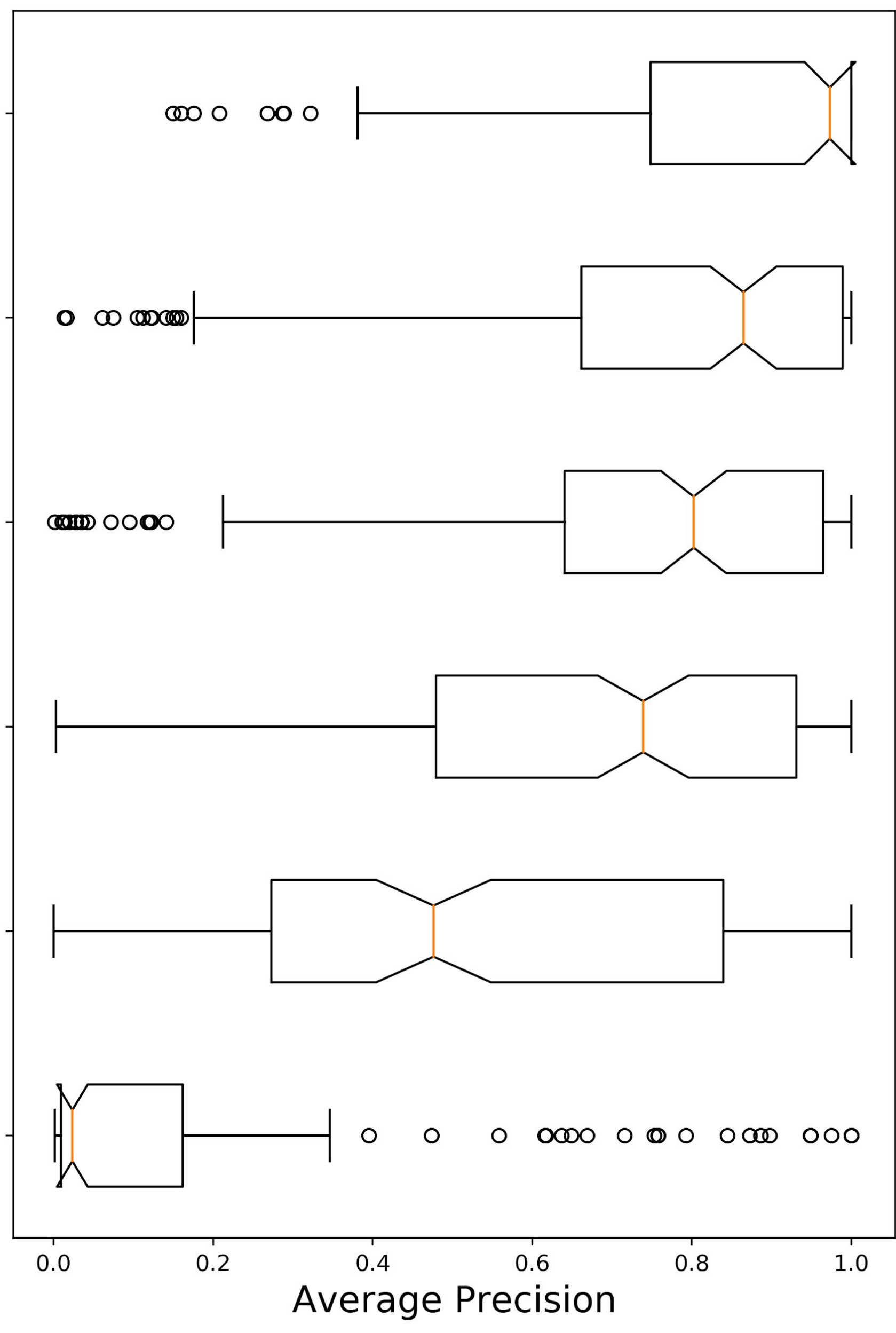Combination of MetaSRA, NLP-ML, and gene expression models weighted by F1 score

Combination of NLP-ML and gene expression models weighted by F1 score

Gene expression models: logistic regression models trained on normalized gene expression data

NLP-ML models: logistic regression models trained on word embedding representations of sample descriptions

MetaSRA: existing tool that constructs knowledge graph from text using ontologies and prior information

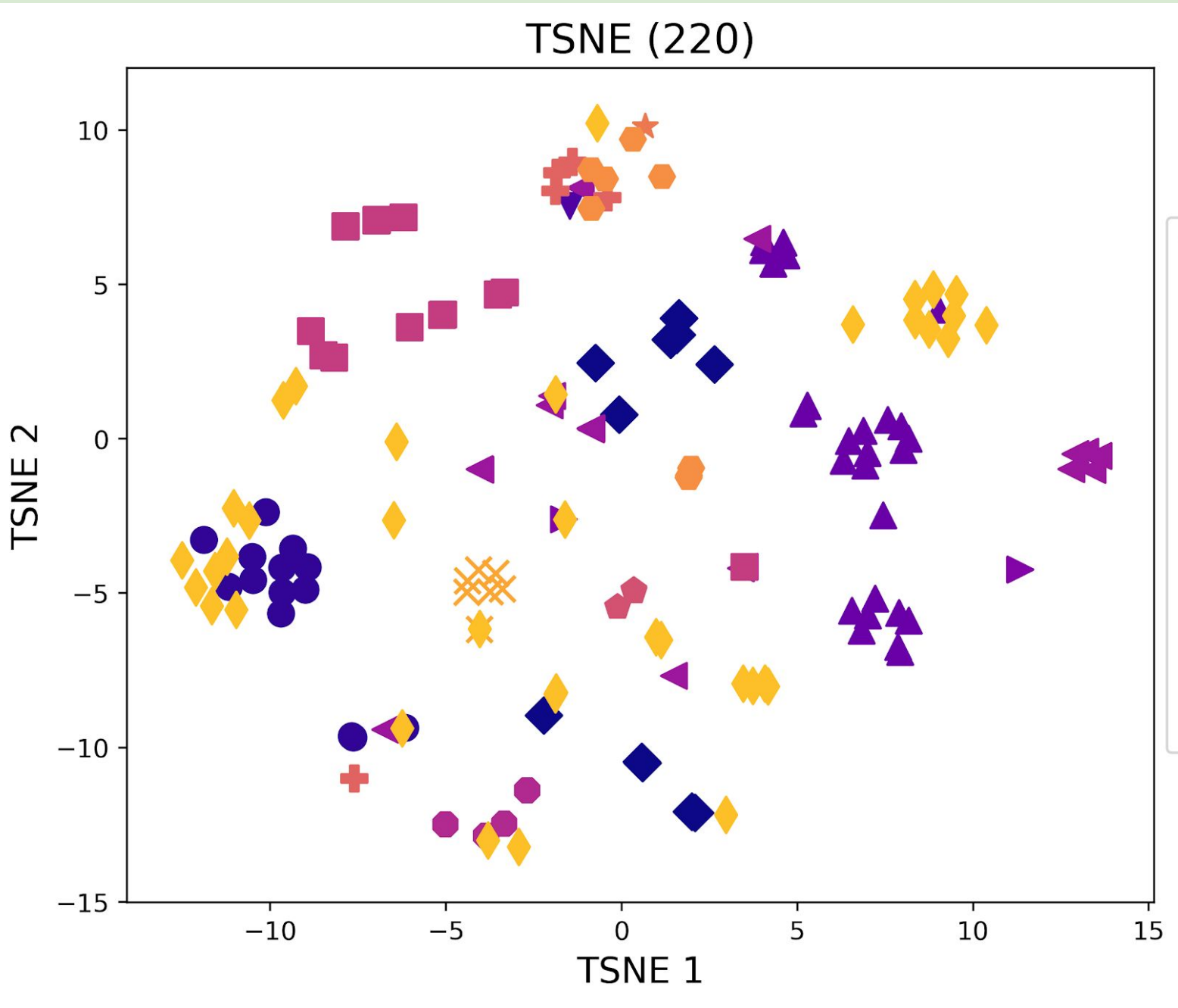TAGGER: named entity recognition tool that directly matches text to a dictionary of ontology terms



Each point in the boxplot is a model to provide annotations for a single tissue/cell-type

## INTERPRETABILITY

Top predicted disease ontology terms for brain model using word embedding representations of ontology terms' description

| DOID | Probability |
| --- | --- |
| neuronal ceroid lipofuscinosis | 0.99999999 |
| complex cortical dysplasia with other brain malformations | 0.99999812 |
| neurodegeneration with brain iron accumulation | 0.99998791 |
| hypomyelinating leukodystrophy | 0.99992988 |
| Parkinson's disease | 0.99992542 |
| Joubert syndrome | 0.99900774 |
| Ritscher-Schinzel syndrome | 0.99848593 |
| holoprosencephaly | 0.99726263 |
| autosomal dominant nocturnal frontal lobe epilepsy | 0.99222049 |
| advanced sleep phase syndrome | 0.97340695 |

tSNE plot of normalized model coefficients colored by high-level anatomical system

1. Bernstein et al., 2017, *Bioinformatics*, 33, 18, 2914-2923.
2. Lee et al., 2013, *Bioinformatics*, 29, 23, 3036–3044.
3. Jensen, 2016, *bioRxiv*, doi.org/10.1101/067132.
4. Akbik et al., 2019, *NAACL Proceedings,* 724–728.

hawki2235@msu.edu