# Determining the Optimal Embedding Technique for Mapping Gene Expression Samples into a Distributed Ontology Space

Filip Jevtic[1], Christopher A Mancuso[1], Arjun Krishnan[1,2]

[1]Department of Computational Mathematics, Science and Engineering, [2]Department of Biochemistry and Molecular Biology

## Overview

- There are more than 1 million publicly available human gene expression samples. These samples together constitute an extremely valuable resource that any researcher can use to gain new biological insights about genes and cellular mechanisms. However, most of these samples lack systematic annotation of the context in which they were generated including, tissue and cell type of origin as well as disease state. Developing computational methods that can automatically label a sample is a grand challenge in biomedical research.

- A new labeling method that develops upon previous methods has emerged, termed *OnClass*[1], which uses a neural network to first generate a mapping between gene expression space to a label space that is created through a low dimensional embedding of an ontology. This mapping simultaneously captures the relationships of different terms in the ontology while also allowing predictions for terms that were never seen during training the model.

- While the *OnClass* method has shown state-of-the-art results for single cell prediction, there was little investigation into how the method works on other classification tasks or how to find the best hyperparameters.

- We investigated this question by developing software that can easily input new data, labels, hyperparameters and embeddings. The train, validation and test sets are produced in a rigorous way. The software is a work in progress but this poster highlights some of the features that will be released.
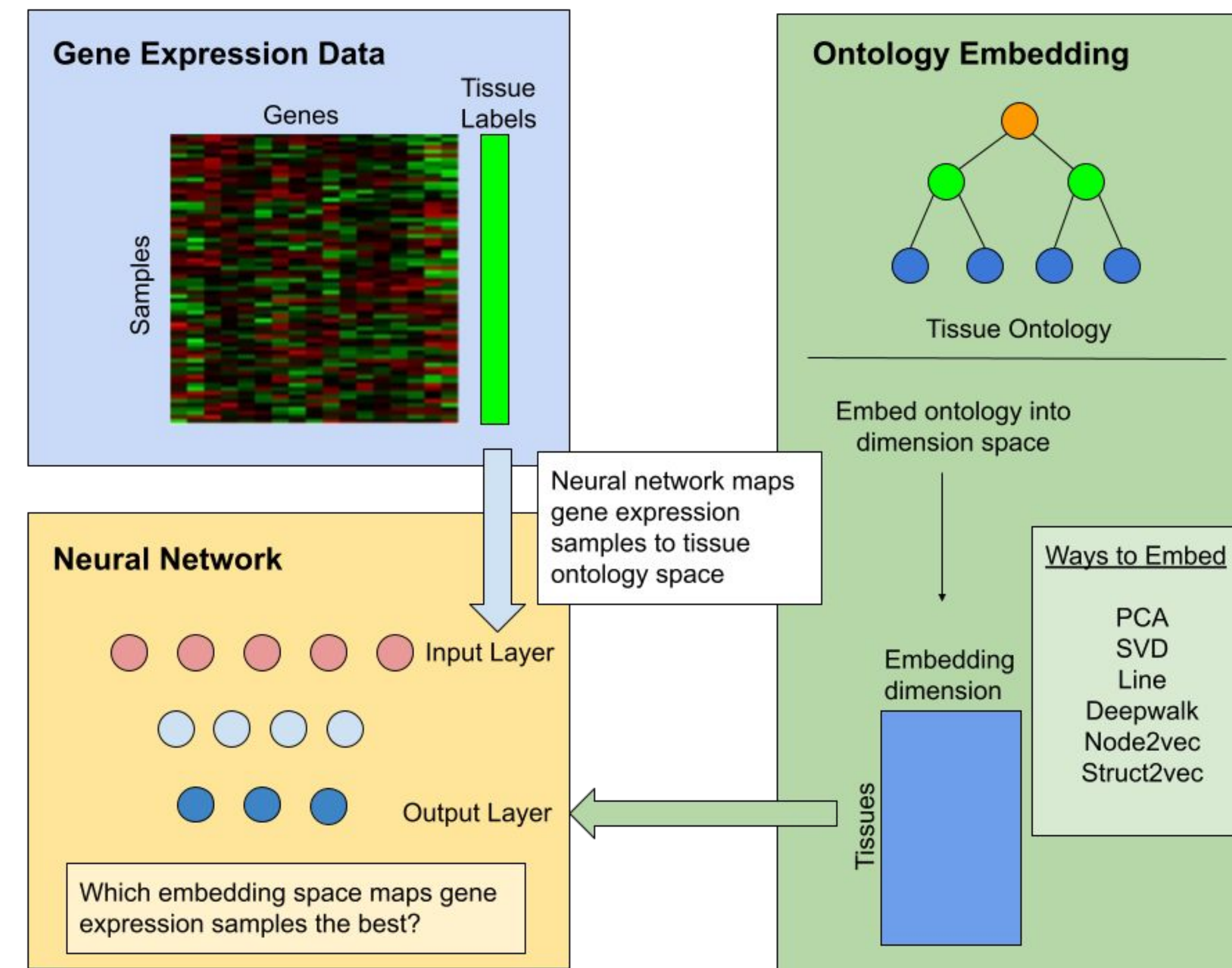
## Previous Methods

### Binary Classification



### Standard Multilabel Classification



### Predictions in Ontology Embedding Space



## Method



**Figure 1.** Workflow for systematic evaluation of ontology embedding techniques.

### Parameters of the Model

**Learning Rate** : How much the model changes weights

**Hidden Layer Units**: The size of the hidden layer

**Data Scaling**: How or if the data in standardized

**Training Label**: How the training loss function is implemented

**Embedding Method**: How the ontology is embedded

**Embedding Size**: The size of the embedding

**Embedding Parameters**: a few examples are diffusion restart value, walk lengths, walk style …

**Evaluation Scheme:** How the predictions are evaluated ex. log2prior score, AUROC …

## Code Base Instructions



## Example Visualizations

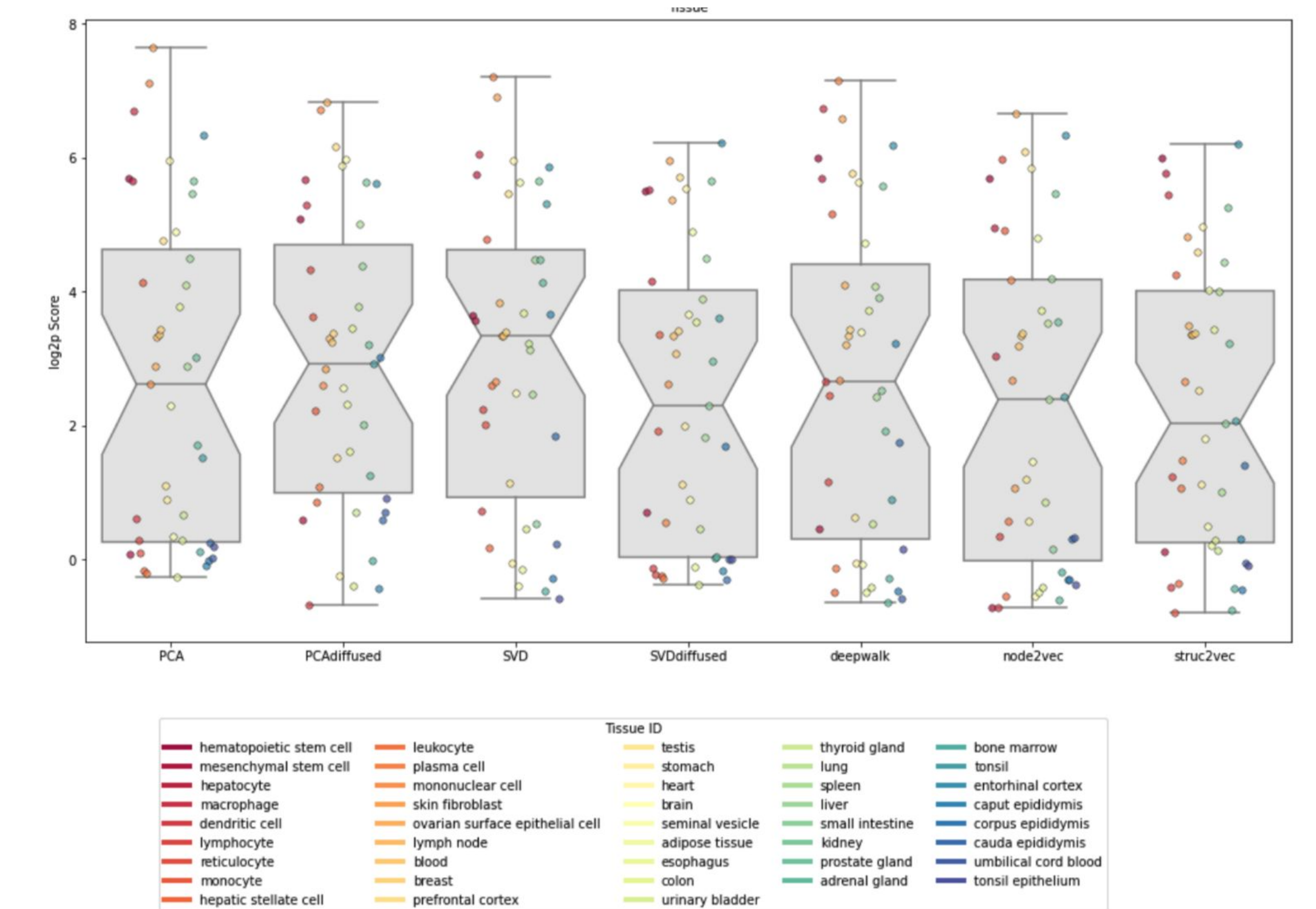**Figure 2.** Best validation model displayed by tissue.



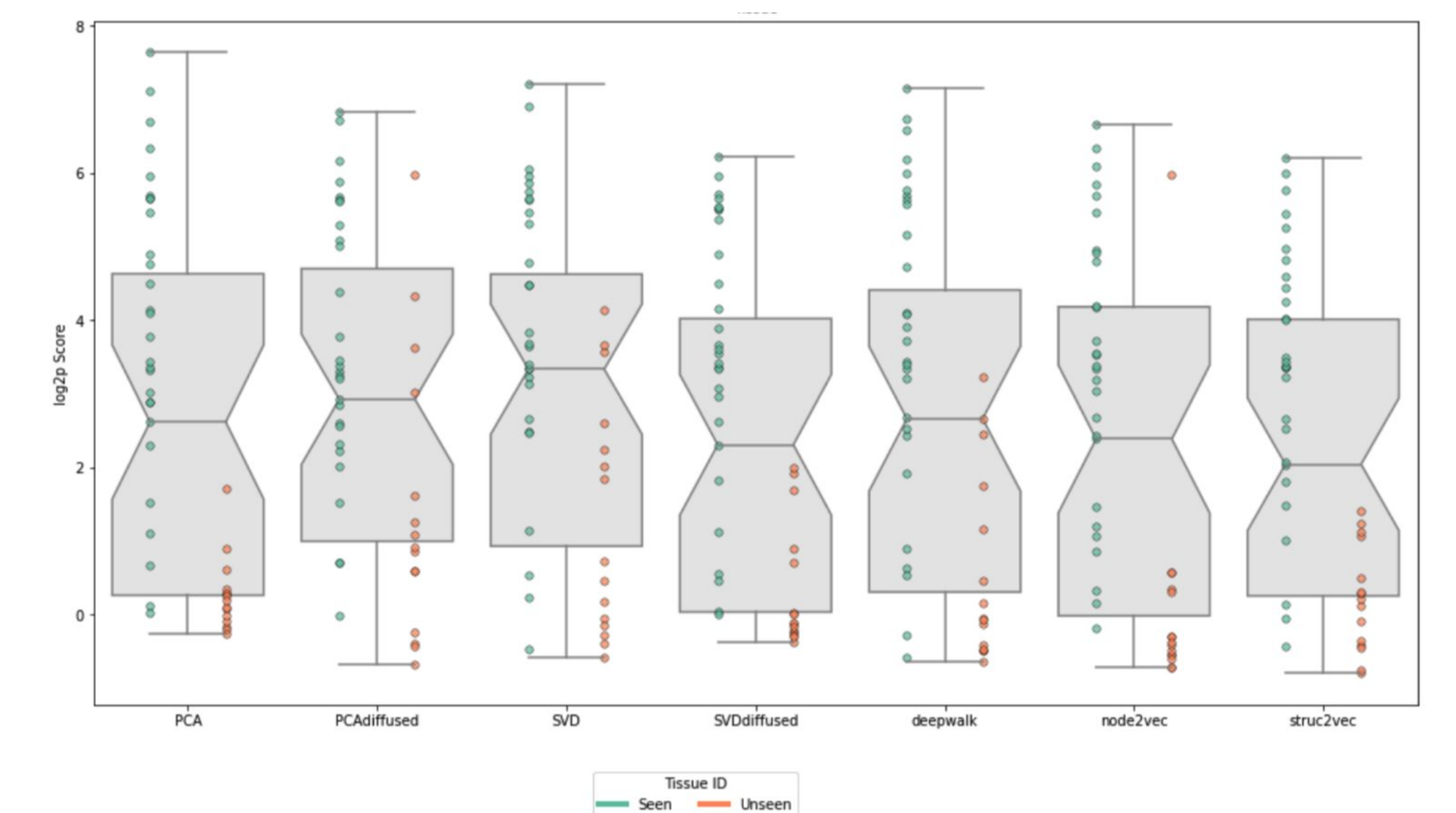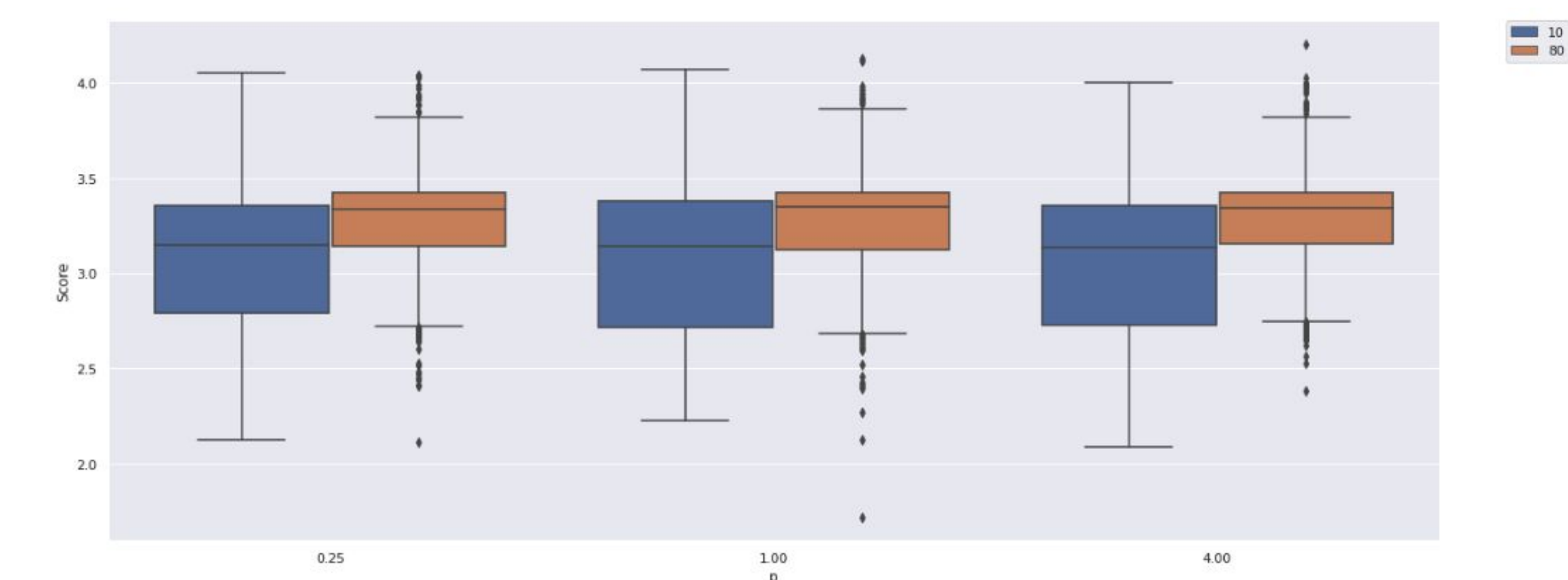**Figure 3.** Best validation model displayed by tissues seen/unseen in training.



**Figure 4.** Performance of node2vec Models Based on P Value, Walk Length and Walk Style
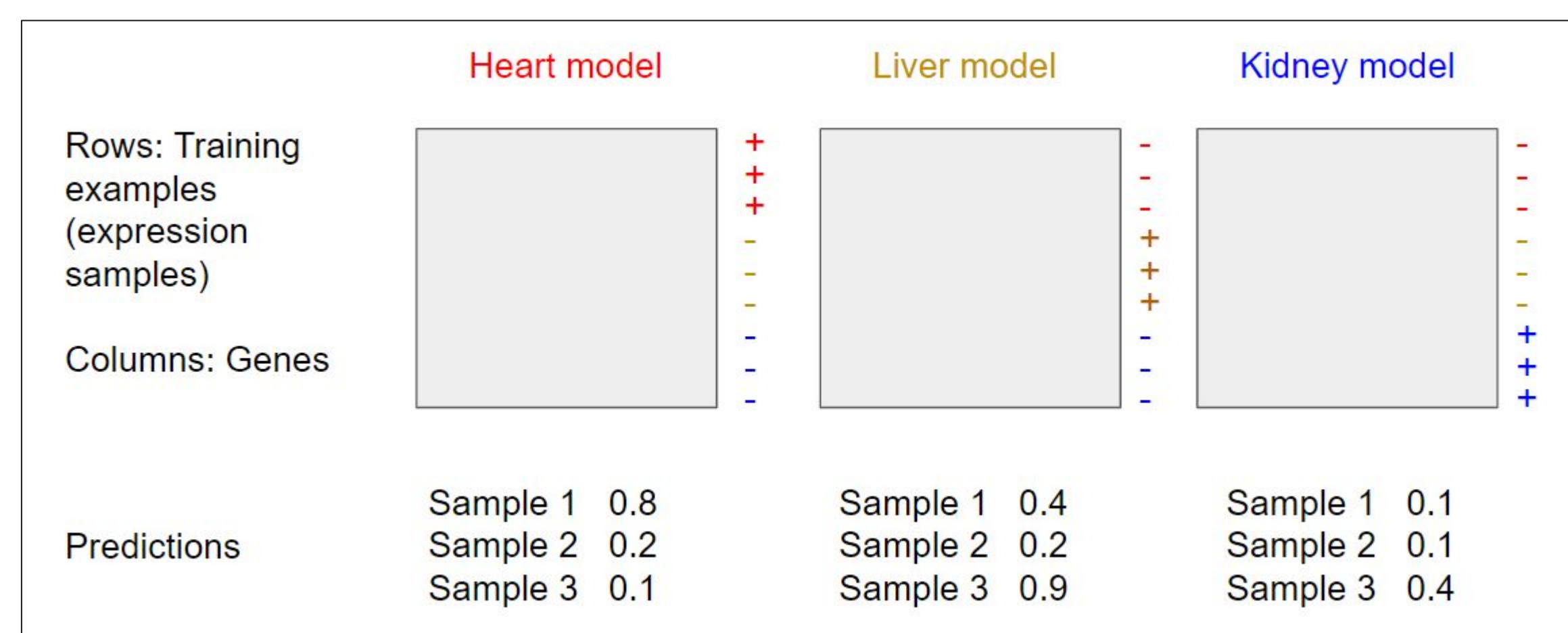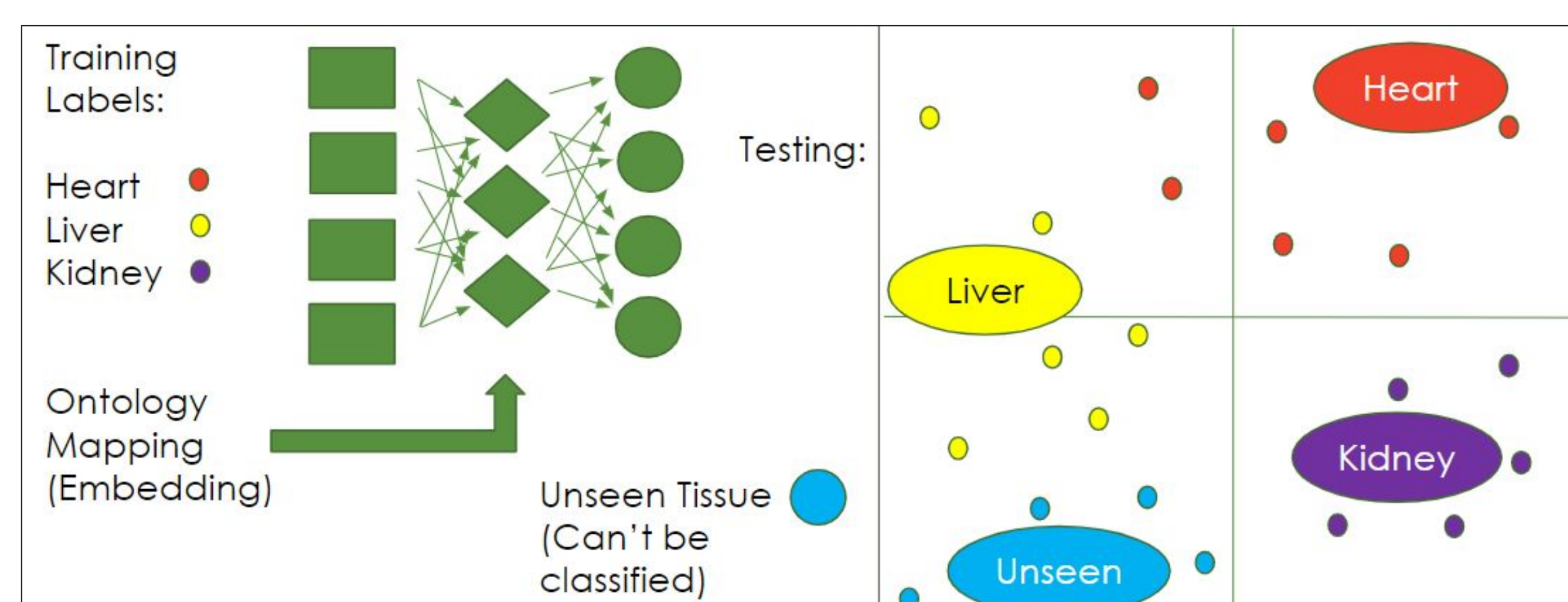


## Reference and Contact and Funding

[1]Wang et.al., (2021) Leveraging the Cell Ontology to classify unseen cell types. Nat. Comm. **12** 5556

E-mail: jevticfi@msu.edu