

1. Introduction

1.1 Motivation

- Protein-Protein Interaction (PPI) networks for model organisms are getting **ever-larger**.
- The size of PPI networks of species which are **not model organisms** is **much smaller**.

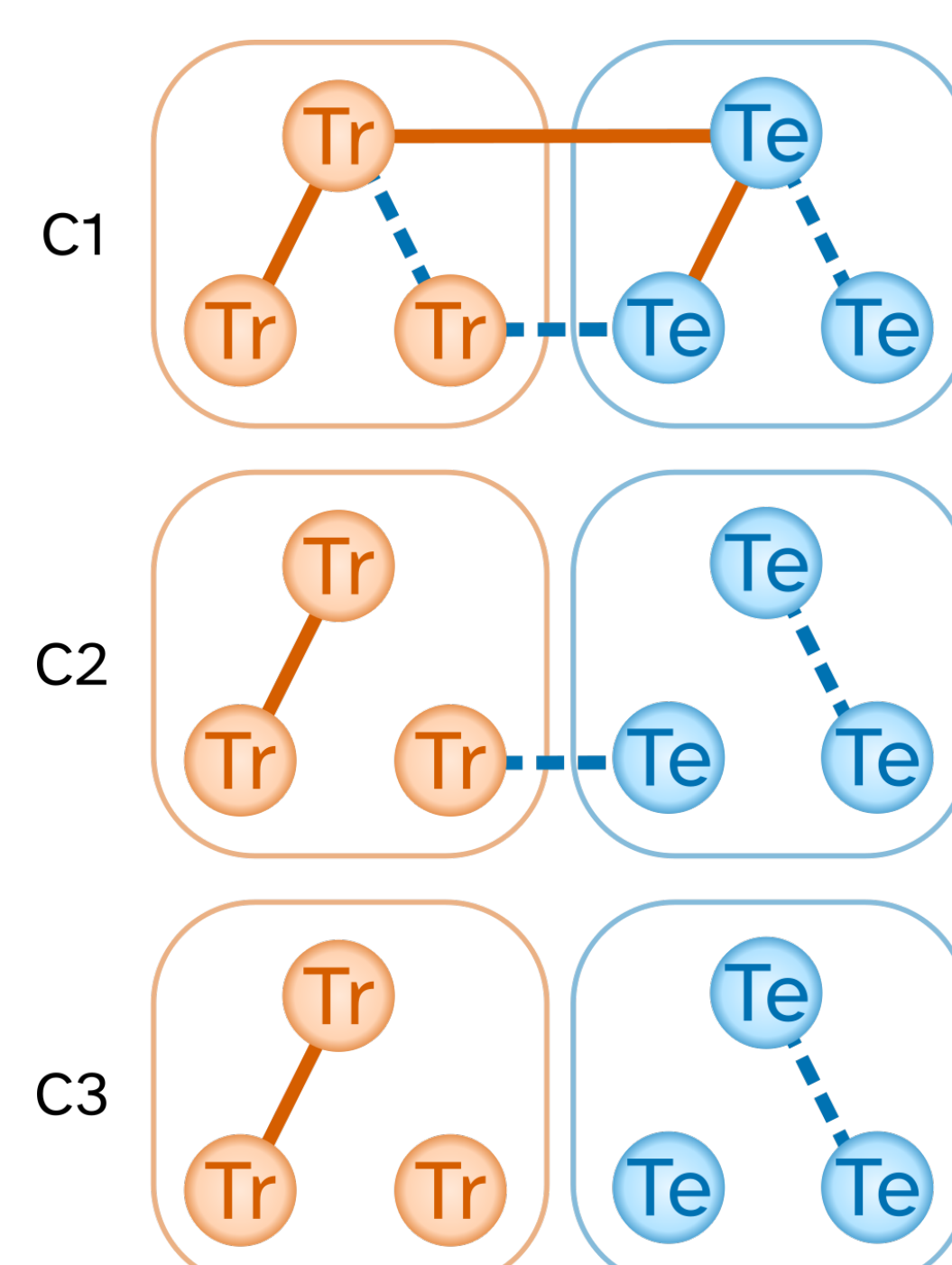
1.2 Out-of-Distribution Predictions

- The organisms with the **most incomplete graphs** often have **too little data** to train good PPI prediction models.
- It's desirable to therefore **train on model organisms** and **testing on species with little data**.
- Many ML models **do not make accurate out-of-distribution (OOD)** predictions [1].

2. Methodology

2.2 Special Considerations for Validation & Testing Dataset

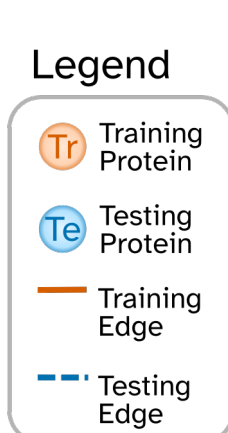
Park & Marcotte identified an information leakage problem with PPI prediction validation techniques [1].



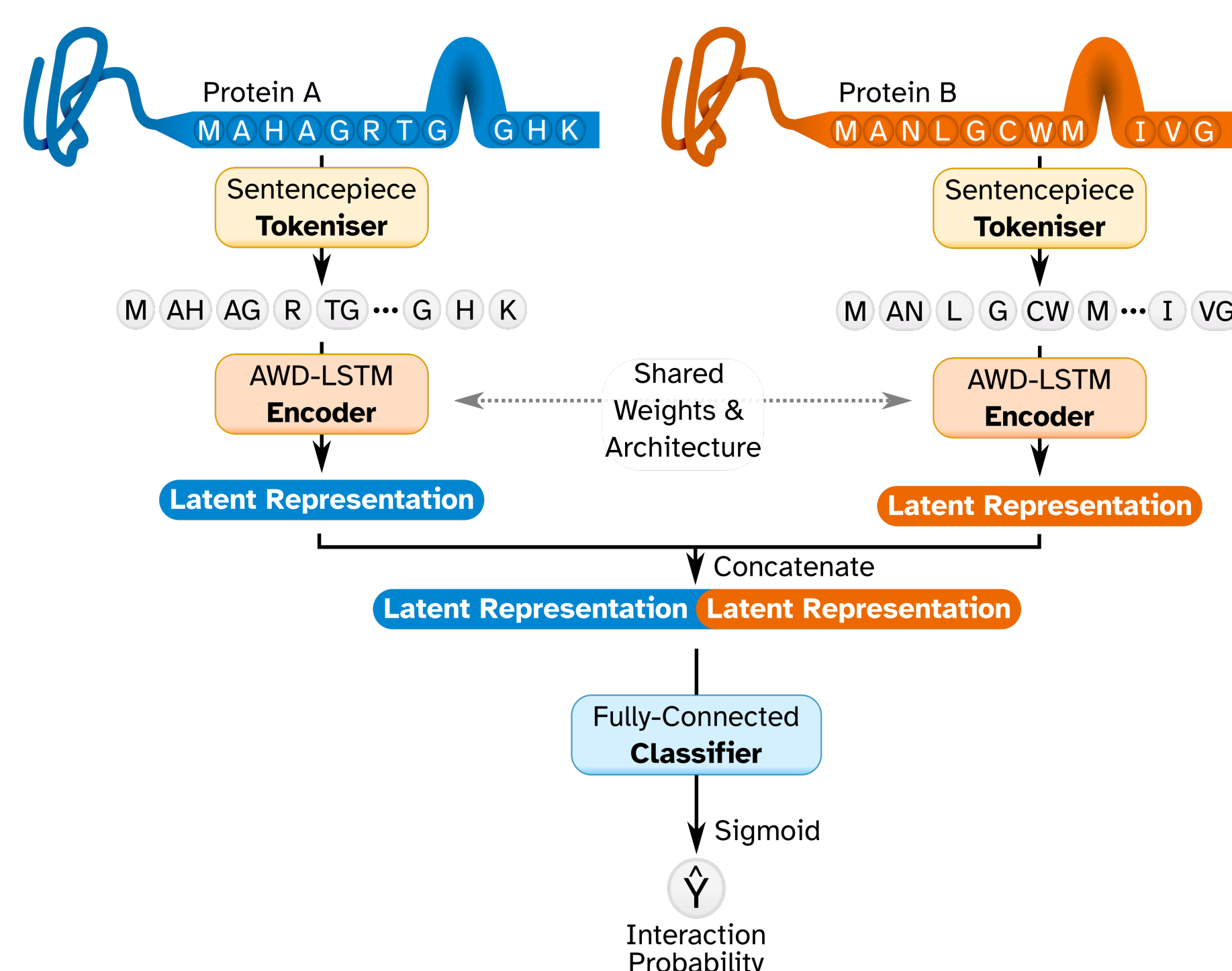
C3 assures no proteins in the testing set.

C2 assure no more than one protein in training pairs are in the testing or validation set.

C1 training pairs may contain one or two proteins found in the testing or validation set.



2.3 Overview of the RAPPID Architecture



RAPPID is a regularised twin neural network that adopts a modified AWD-LSTM [2].

RAPPID considers **pairs of amino acid (AA) sequences** with an interaction label [3].

AA sequences are first tokenised with the **Sentencepiece algorithm** [4].

Fixed-length latent vector representations are computed for each sequence using **bi-directional AWD-LSTMs**.

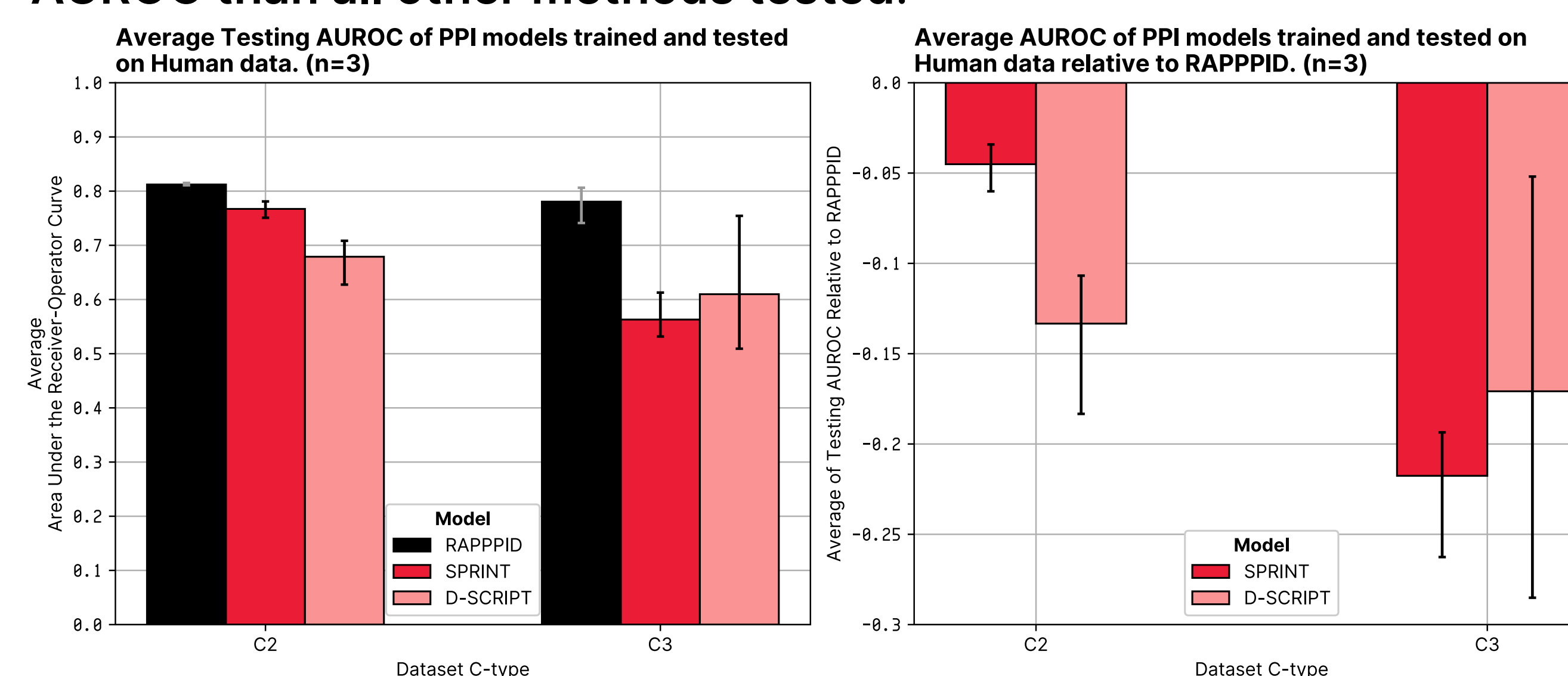
Latent vectors are **concatenated and are inputted** into a two-layer fully-connected **classification head**.

Output of the classifier is the **interaction probability**.

3. Results

3.1 RAPPID generalises better than leading PPI prediction methods.

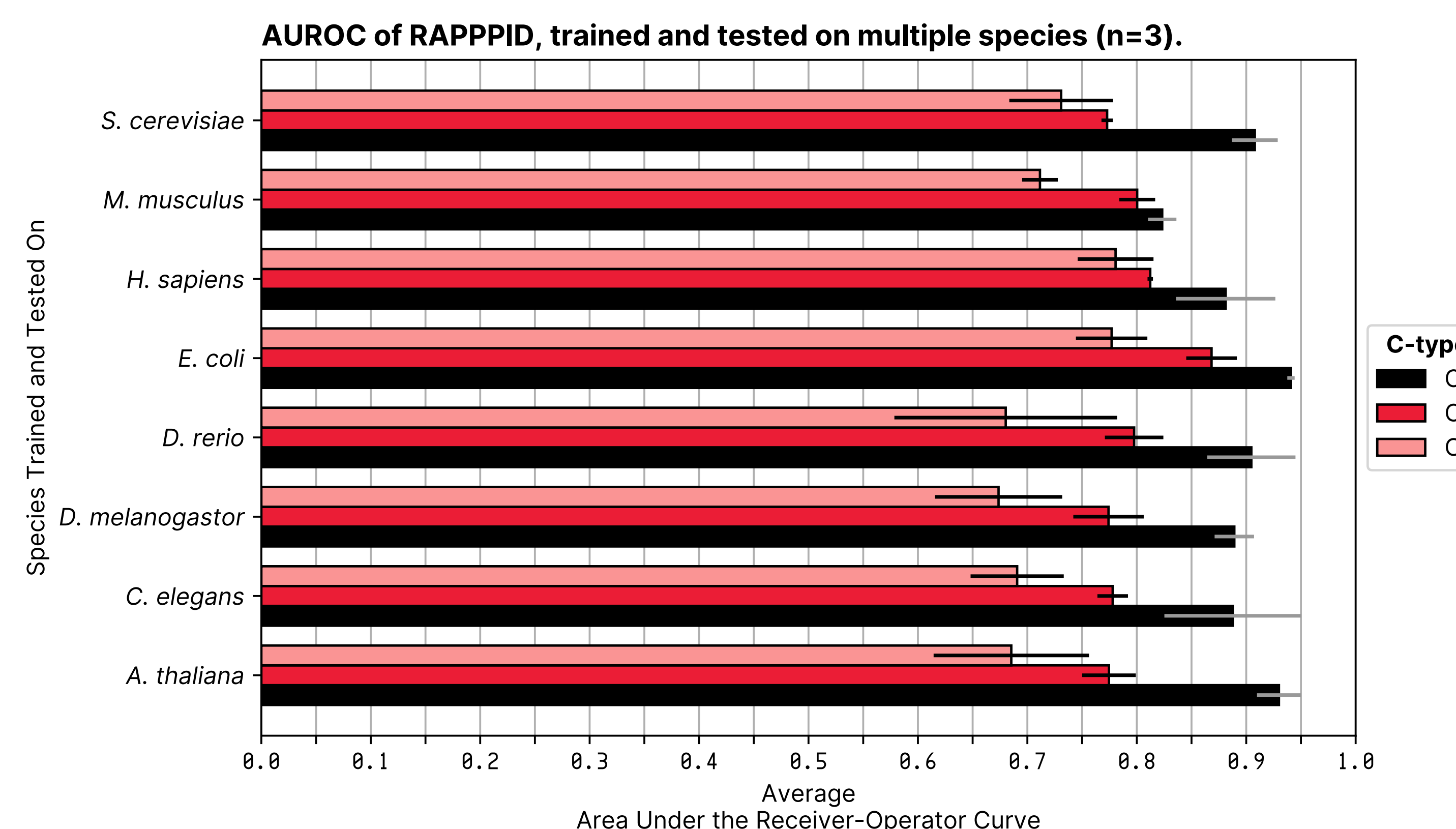
Across C2 and C3 *Homo sapiens* datasets, RAPPID achieves a higher AUROC than all other methods tested.



3.2 RAPPID intra-species performance.

RAPPID models trained and tested on various species maintain a high degree of performance.

The average AUROC of three different random seeds and data splits are reported.



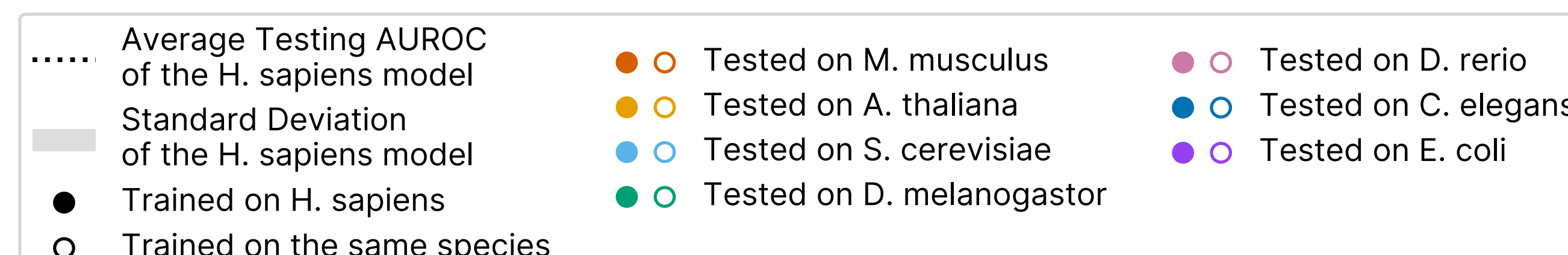
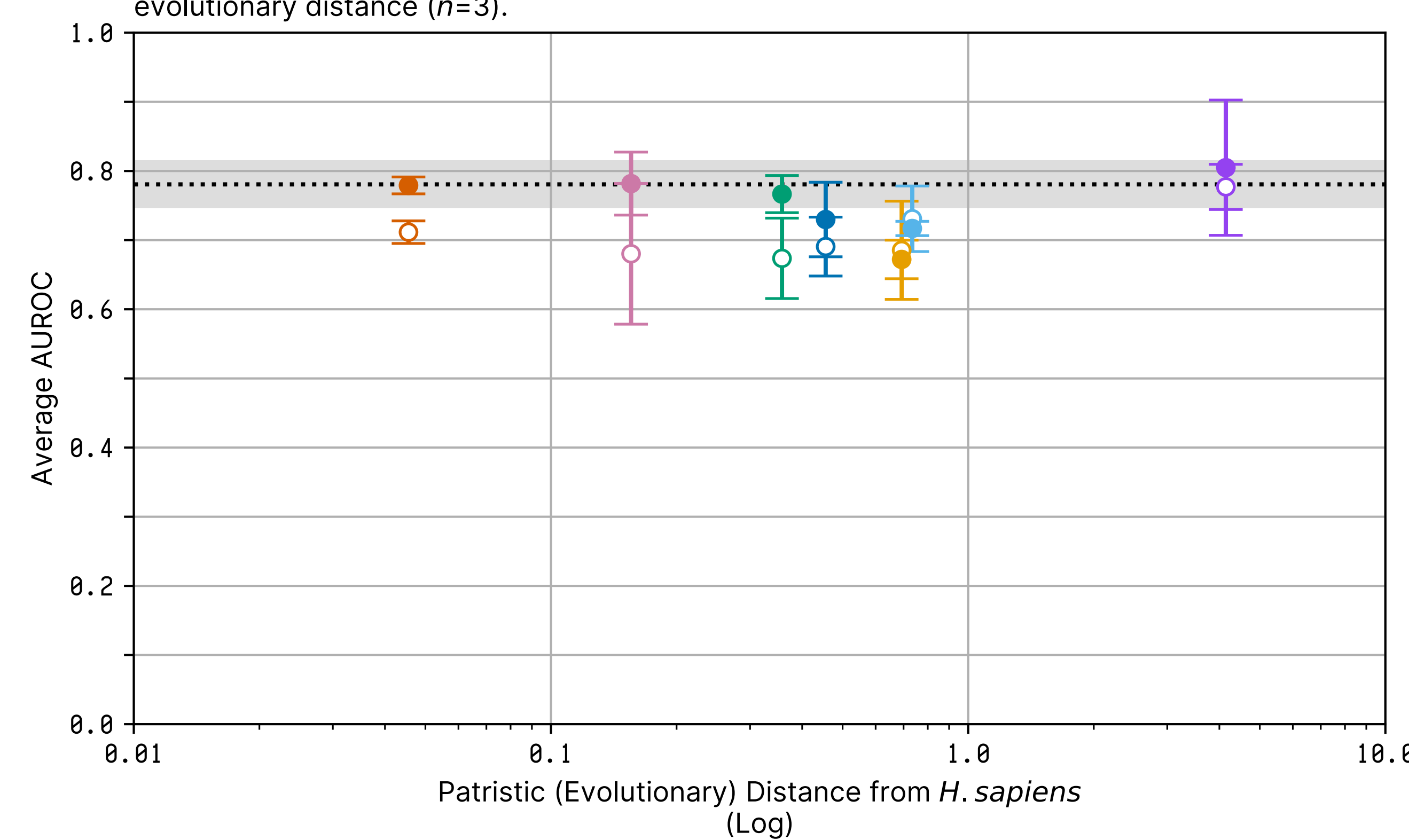
3.3 RAPPID cross-species performance.

The performance of a C3 RAPPID model trained on *H. sapiens* PPI data was measured on testing sets composed of PPIs of other species.

This was compared to RAPPID trained & tested on the same species.

RAPPID trained on Human PPIs maintains its performance when tested on PPIs of other species.

Average AUROC of a Human RAPPID model tested on other species as a function of their evolutionary distance (n=3).



3. Results

3.4 RAPPID transfer-learning performance.

To evaluate how RAPPID out-of-distribution performance changes with fine-tuning, RAPPID was trained on *H. sapiens*, and fine-tuned on *Escherichia coli* PPI data. Finally, the model was tested on *E. coli* PPIs.

The means of three different random seeds and data splits are reported below.

Training Species	Testing Species	Fine-Tuning Species	AUROC (Mean)	APR (Mean)
E. coli	E. coli	None	0.818	0.840
H. sapiens	E. coli	None	0.839	0.877
H. sapiens	E. coli	E. coli	0.867	0.890

4. Future Directions

4.1 Online full-proteome prediction server.

We are currently working on developing an online server for PPI prediction using RAPPID.

We will leverage RAPPID's efficiency to enable the prediction of an inputted amino acid sequence against entire known proteomes of a specified organism.

4.2 Tools for therapeutic peptide discovery.

RAPPID's out-of-distribution performance on unseen proteins across species lends itself quite naturally to the context of therapeutic peptide discovery.

We're developing tools to aid in the discovery of therapeutic peptides using RAPPID's online interface.

5. Acknowledgments



6. References

- Park Y, Marcotte EM. Flaws in evaluation schemes for pair-input computational predictions. *Nat Methods*. 2012 Dec;9(12):1134–6.
- Merity S, Keskar NS, Socher R. Regularizing and Optimizing LSTM Language Models. arXiv:170802182 [cs] [Internet].
- Szymborski J, Emad A. RAPPID: towards generalizable protein interaction prediction with AWD-LSTM twin networks. *Bioinformatics*. 2022 Aug 15;38(16):3958–67.
- Kudo T, Richardson J. SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In: *EMNLP 2018*

See more information online
https://jszym.com/meetings/2023_cshl

