

RAPPPID: Improving Protein Interaction Prediction on Unseen Proteins

SCAN ME



Joseph Szymborski & Amin Emad
BIRS 2022: Deep Learning for Genetics,
Genomics and Metagenomics



Banff International Research Station
for Mathematical Innovation and Discovery

Introduction

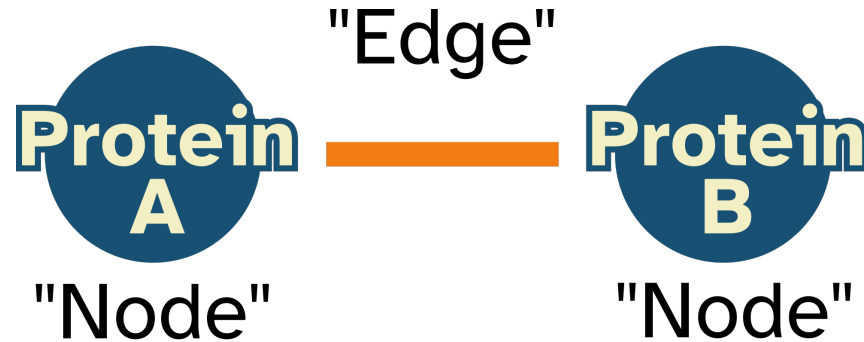
- Joseph Szymborski
 - McGill University,
Department of Electrical & Computer Engineering
 - Mila, Quebec AI Institute
 - PhD Student in Amin Emad's COMBINE Lab



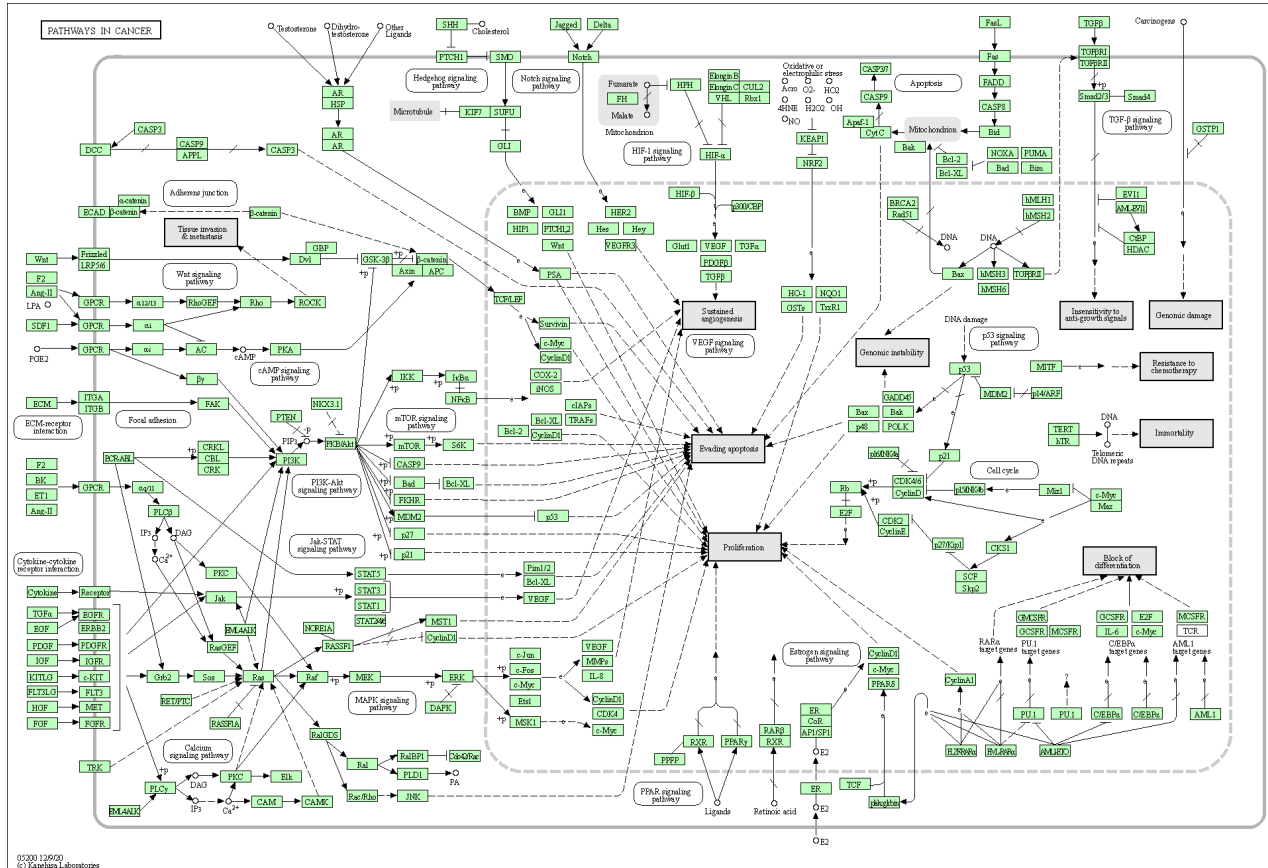
Background: Protein-Protein Interactions

- I've spent the last few years thinking about Protein-Protein Interactions (PPIs).
- Bio' processes as an undirected graph of PPIs.

* An incomplete model,
but it's gotten us pretty far.



Background: Protein-Protein Interactions

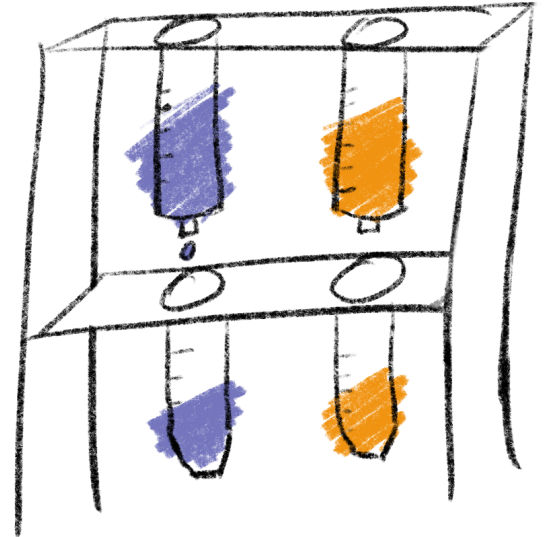


See:
 Kanehisa M. et al.
 10.1093/nar/gkr988



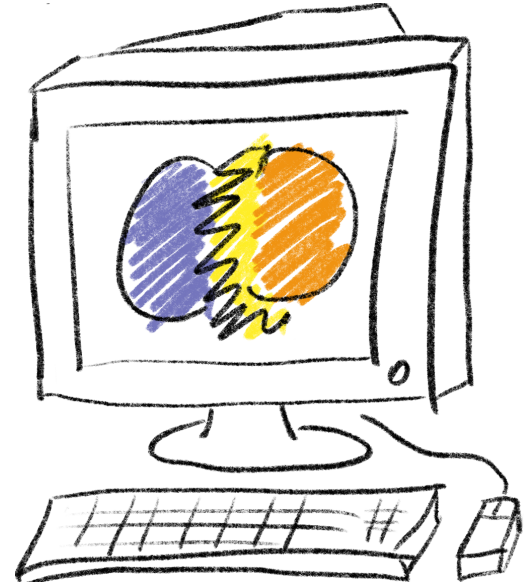
Background: Protein-Protein Interactions

- Protein interactions are typically identified through **“wet lab” experiments**.
- These experiments typically:
 - Take days/weeks.
 - Expensive reagents.
 - Often produce a lot of plastic waste.
 - Are quite definitive.



Background: Protein-Protein Interactions

- Predicting protein interactions using **computational models** try to address some of the **trade-offs of lab experiments**.
 - Take seconds/minutes.
 - Low-to-no cost.
 - Consume electricity and produces e-waste.
 - Not yet definitive.



Background: Protein-Protein Interactions

Given two proteins, do they interact?



Background: Protein-Protein Interactions

Human succinyl CoA-transferase
E. coli acetate Co-A transferase α
E. coli acetate Co-A transferase β



Homology

Marcotte *et al.*, 1999

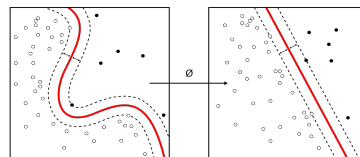
Background: Protein-Protein Interactions

Human succinyl CoA-transferase
E. coli acetate Co-A transferase α
E. coli acetate Co-A transferase β



Homology

Marcotte *et al.*, 1999



Support Vector Machines

Ben-Hur & Noble, 2005

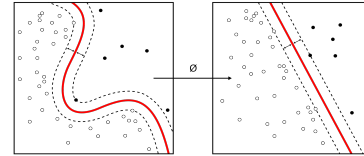
Background: Protein-Protein Interactions

Human succinyl CoA-transferase
E. coli acetate Co-A transferase α
E. coli acetate Co-A transferase β



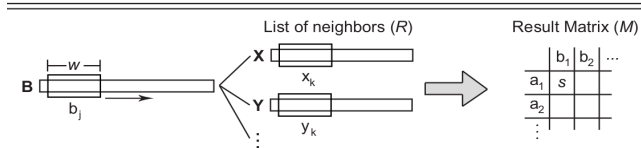
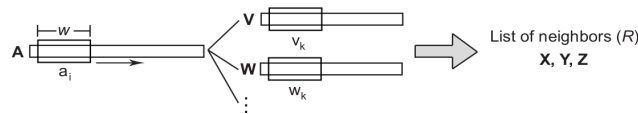
Homology

Marcotte *et al.*, 1999



Support Vector Machines

Ben-Hur & Noble, 2005



Sequence Similarity

Pitre *et al.*, 2006

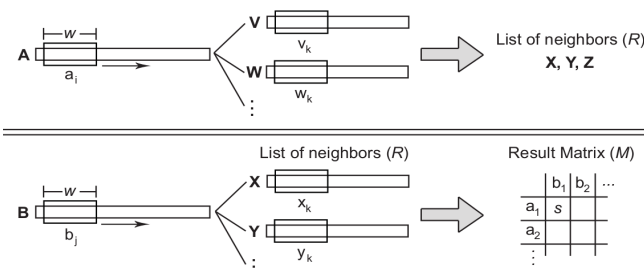
Background: Protein-Protein Interactions

Human succinyl CoA-transferase
E. coli acetate Co-A transferase α
E. coli acetate Co-A transferase β



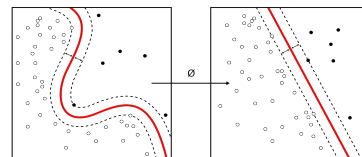
Homology

Marcotte *et al.*, 1999



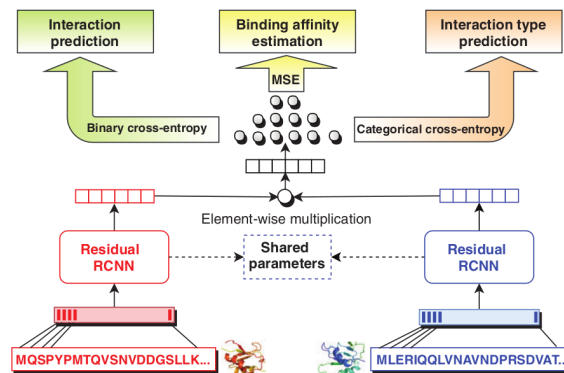
Sequence Similarity

Pitre *et al.*, 2006



Support Vector Machines

Ben-Hur & Noble, 2005

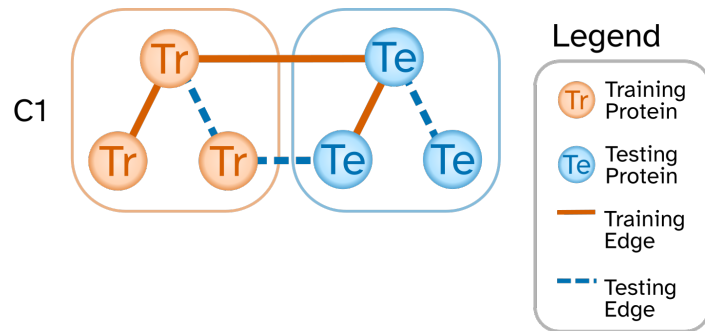


Deep Learning

Chen *et al.*, 2019

Background: Protein-Protein Interactions

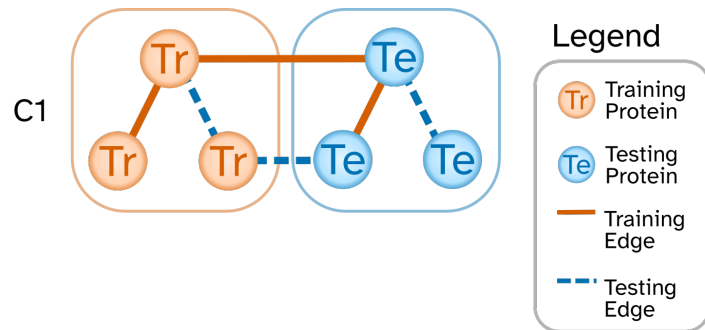
AUROC on *H. sapiens*



Background: Protein-Protein Interactions

AUROC on *H. sapiens*

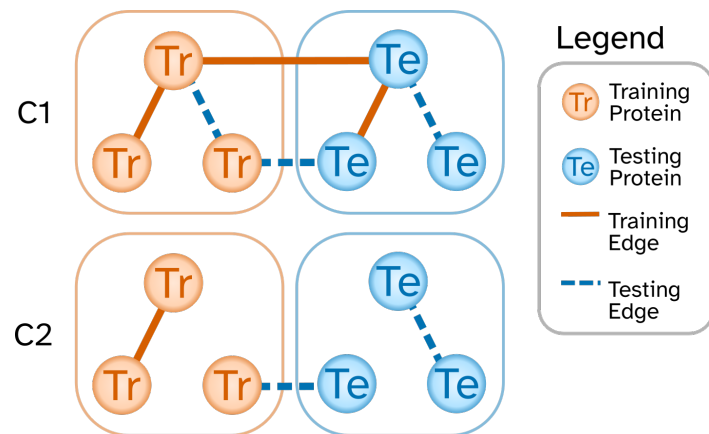
| |
|-----------|
| C1 |
| 0.81±0.01 |
| 0.85±0.01 |
| 0.64±0.01 |
| 0.77±0.01 |
| 0.81±0.01 |
| 0.77±0.01 |
| 0.56±0.01 |



Background: Protein-Protein Interactions

AUROC on *H. sapiens*

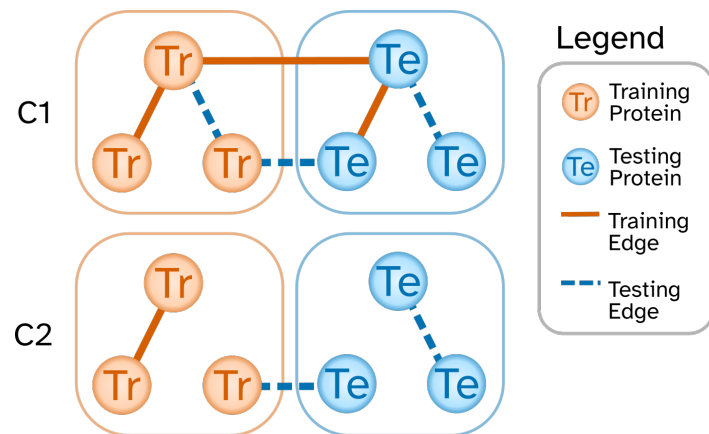
| |
|-----------|
| C1 |
| 0.81±0.01 |
| 0.85±0.01 |
| 0.64±0.01 |
| 0.77±0.01 |
| 0.81±0.01 |
| 0.77±0.01 |
| 0.56±0.01 |



Background: Protein-Protein Interactions

AUROC on *H. sapiens*

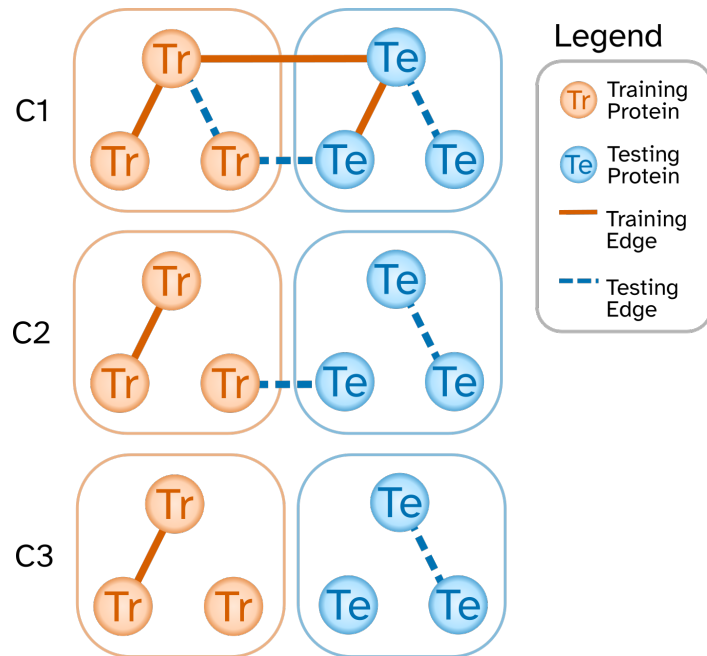
| C1 | C2 |
|-----------|-----------|
| 0.81±0.01 | 0.61±0.01 |
| 0.85±0.01 | 0.60±0.01 |
| 0.64±0.01 | 0.55±0.01 |
| 0.77±0.01 | 0.57±0.02 |
| 0.81±0.01 | 0.59±0.01 |
| 0.77±0.01 | 0.64±0.01 |
| 0.56±0.01 | 0.53±0.01 |



Background: Protein-Protein Interactions

AUROC on *H. sapiens*

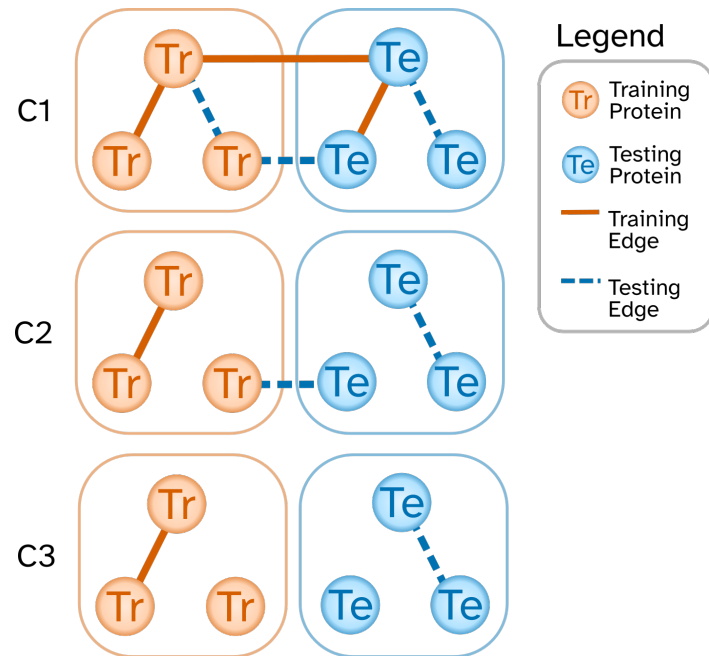
| C1 | C2 |
|-----------------|-----------------|
| 0.81 ± 0.01 | 0.61 ± 0.01 |
| 0.85 ± 0.01 | 0.60 ± 0.01 |
| 0.64 ± 0.01 | 0.55 ± 0.01 |
| 0.77 ± 0.01 | 0.57 ± 0.02 |
| 0.81 ± 0.01 | 0.59 ± 0.01 |
| 0.77 ± 0.01 | 0.64 ± 0.01 |
| 0.56 ± 0.01 | 0.53 ± 0.01 |



Background: Protein-Protein Interactions

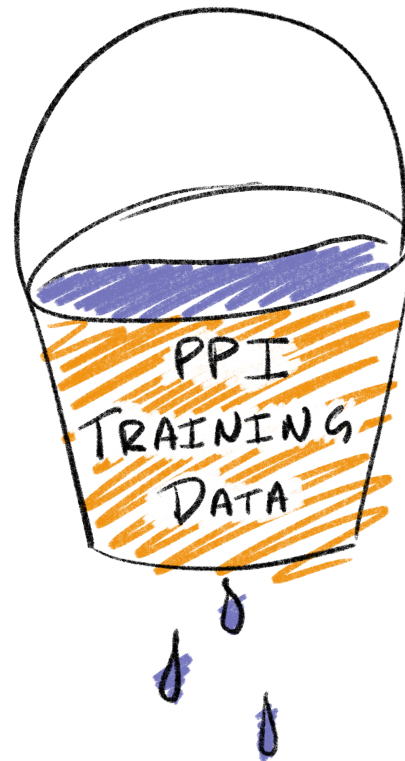
AUROC on *H. sapiens*

| C1 | C2 | C3 |
|-----------------|-----------------|-----------------|
| 0.81 ± 0.01 | 0.61 ± 0.01 | 0.58 ± 0.03 |
| 0.85 ± 0.01 | 0.60 ± 0.01 | 0.58 ± 0.02 |
| 0.64 ± 0.01 | 0.55 ± 0.01 | 0.50 ± 0.00 |
| 0.77 ± 0.01 | 0.57 ± 0.02 | 0.53 ± 0.02 |
| 0.81 ± 0.01 | 0.59 ± 0.01 | 0.54 ± 0.02 |
| 0.77 ± 0.01 | 0.64 ± 0.01 | 0.59 ± 0.02 |
| 0.56 ± 0.01 | 0.53 ± 0.01 | 0.54 ± 0.02 |



The Problem?

- It's hard to plug data leaks in PPI datasets.
- Many models depend on these leaks for their performance.
- How do we plug the leak?



Introducing RAPPID

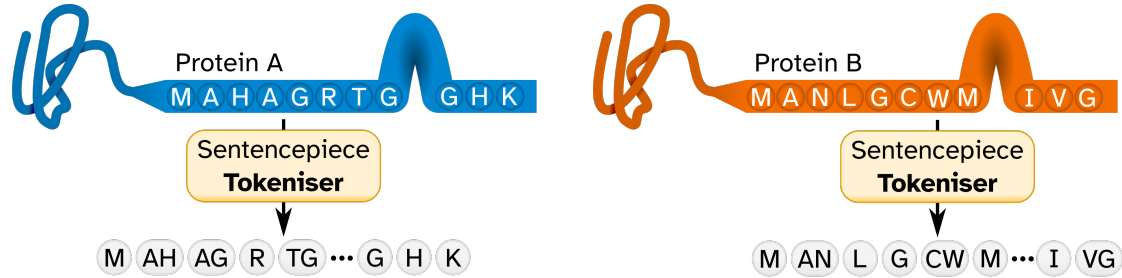
Regularised **A**utomatic **P**rediction of **P**rotein-**P**rotein **I**nteractions using **D**eep Learning

Szyborski, J. & Emad, A. RAPPID: Towards Generalisable Protein Interaction Prediction with AWD-LSTM Twin Networks. bioRxiv 2021.08.13.456309 (2021) doi:10.1101/2021.08.13.456309.

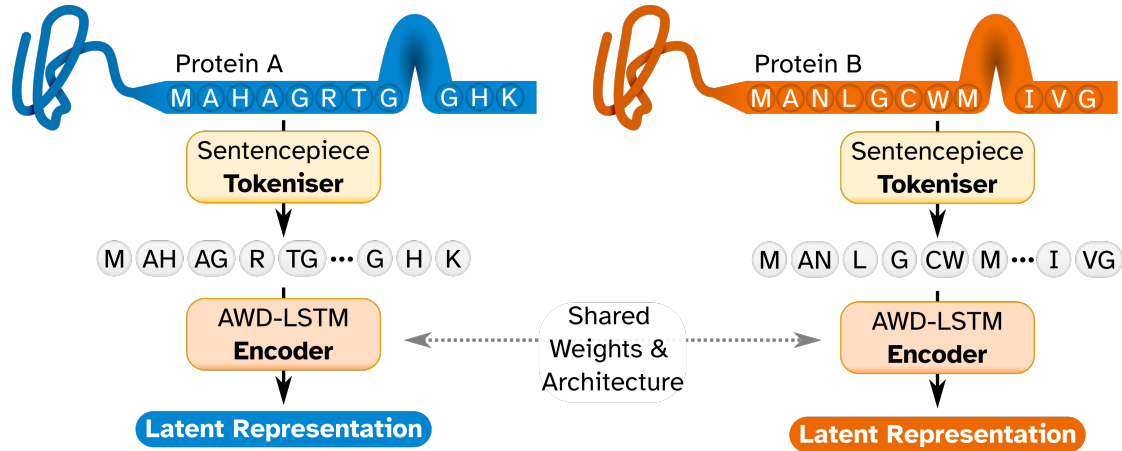
RAPPPID Architecture



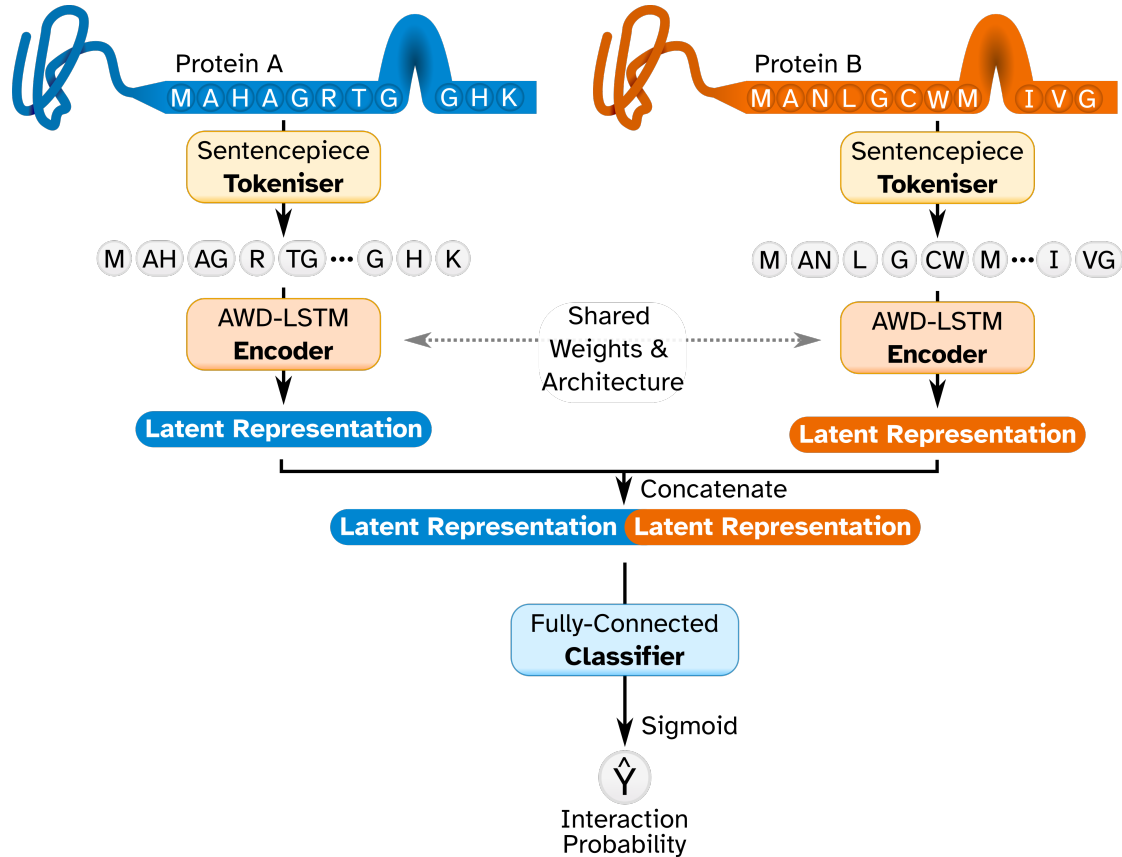
RAPPPID Architecture



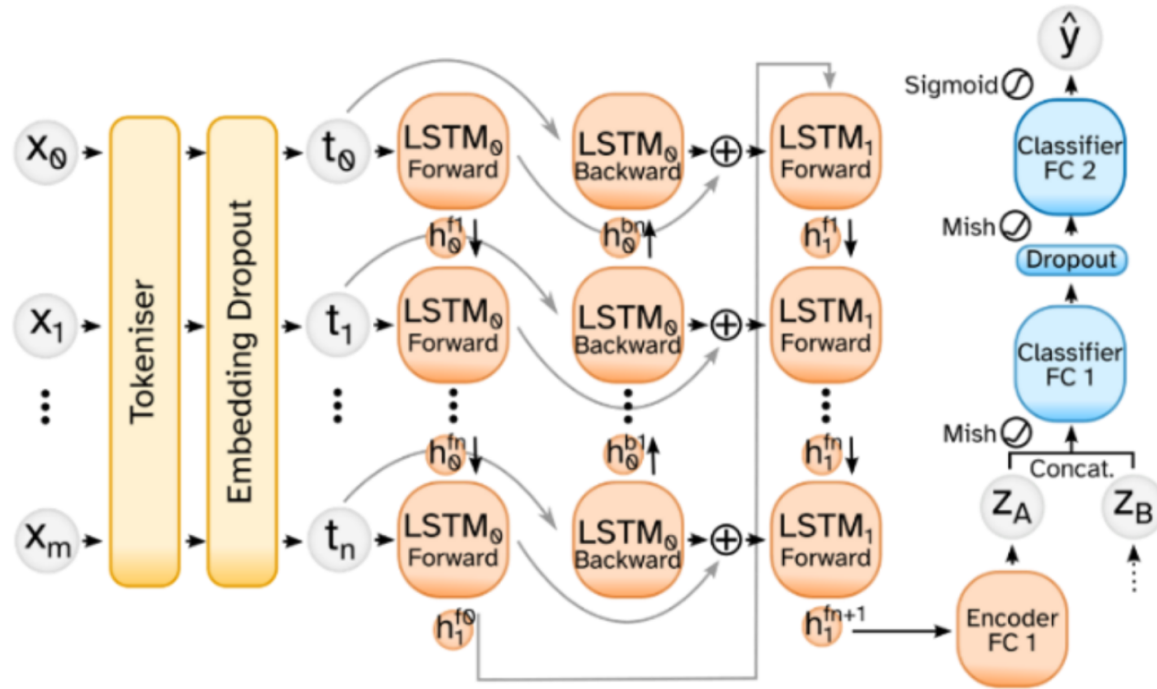
RAPPPID Architecture



RAPPPID Architecture



RAPPPID Architecture



What makes RAPPPID different?

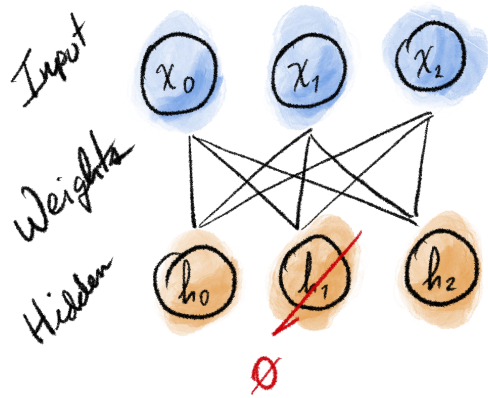
- In short, lots of regularisation
 - AWD-LSTM
 - Embedding dropout
 - Ranger21 Optimiser
 - Stochastic Weight Averaging (SWA)

What makes RAPPPID different?

- In short, lots of regularisation
 - AWD-LSTM
 - Embedding dropout
 - Ranger21 Optimiser
 - Stochastic Weight Averaging (SWA)
- Also
 - Sentencepiece tokenisation

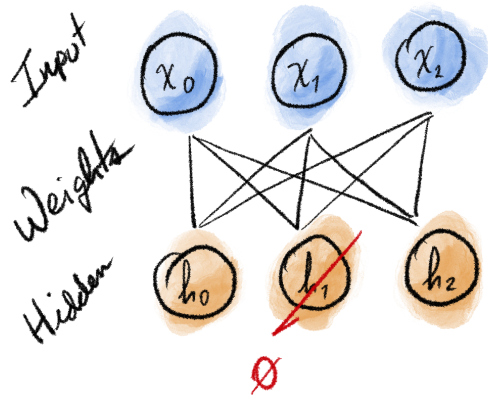
Regularising Recurrent Networks

Dropout

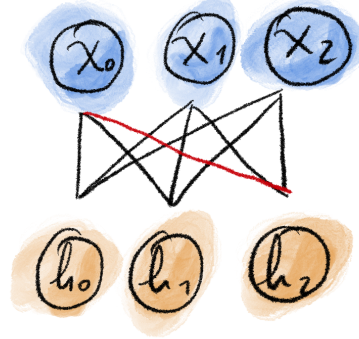


Regularising Recurrent Networks

Dropout

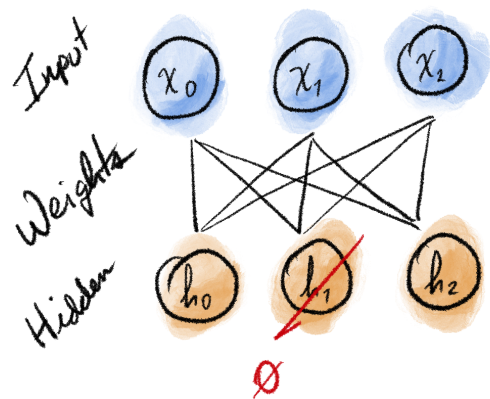


Dropconnect

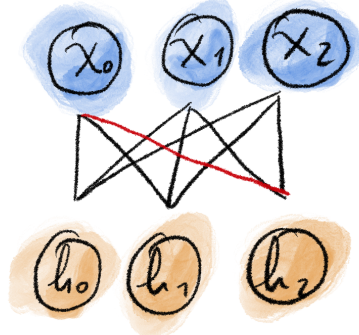


Regularising Recurrent Networks

Dropout



Dropconnect



Embedding Dropout

| | | | | | | | |
|----------|--------------|--------------|--------------|-----|--------------|----------|--------------|
| A | 1 | 0 | 0 | ... | 0 | 0 | 0 |
| C | 0 | 1 | 0 | ... | 0 | 0 | 0 |
| D | 0 | 0 | 1 | ... | 0 | 0 | 0 |
| . | | | | . | | | |
| . | | | | . | | | |
| . | | | | . | | | |
| V | 0 | 0 | 0 | ... | 1 | 0 | 0 |
| W | 0 | 0 | 0 | ... | 0 | 1 | 0 |
| Y | 0 | 0 | 0 | ... | 0 | 0 | 1 |

Weight Decay

- Just a fancy name for L2 weight regularisation.

$$L = l + \lambda \|w\|_2$$

Regularised Loss

Loss

Weight Decay Parameter

L2 Norm of Model Weights

Averaged Stochastic Gradient Descent (ASGD)

- ASGD simply keeps a running average of the weights.
 - often through each epoch.
- SGD is then applied on those averaged weights instead.

$$\bar{w}_{t:T} := \frac{1}{T-t} \sum_{t'=t}^{T-1} w_{t'} .$$

Stochastic Weight Averaging (SWA)

- Very similar to ASGD but keeps a pair of weights:
 - One that the optimiser minimises (w).
 - Another that is a running average of the previous weight (w_{SWA}).

$$w_{SWA} \leftarrow \frac{w_{SWA} \cdot n_{models} + w}{n_{models} + 1}$$

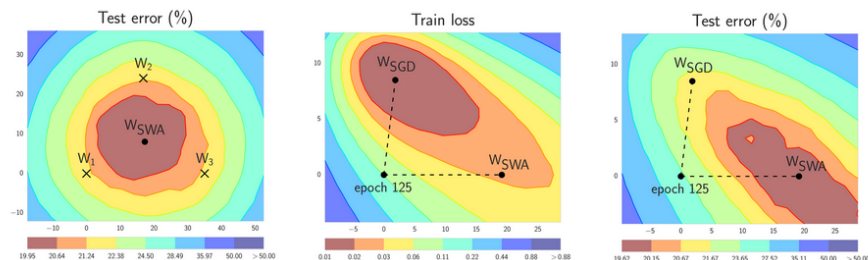
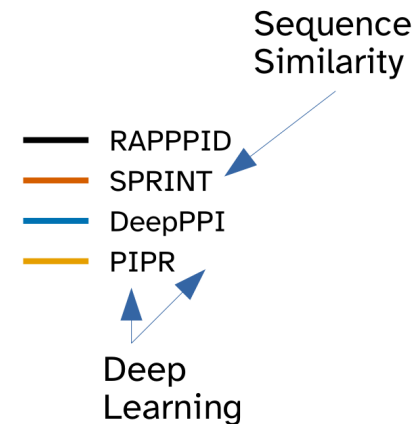
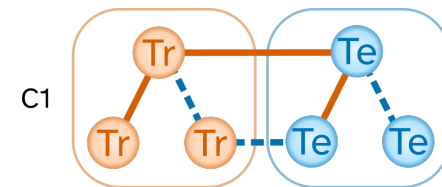
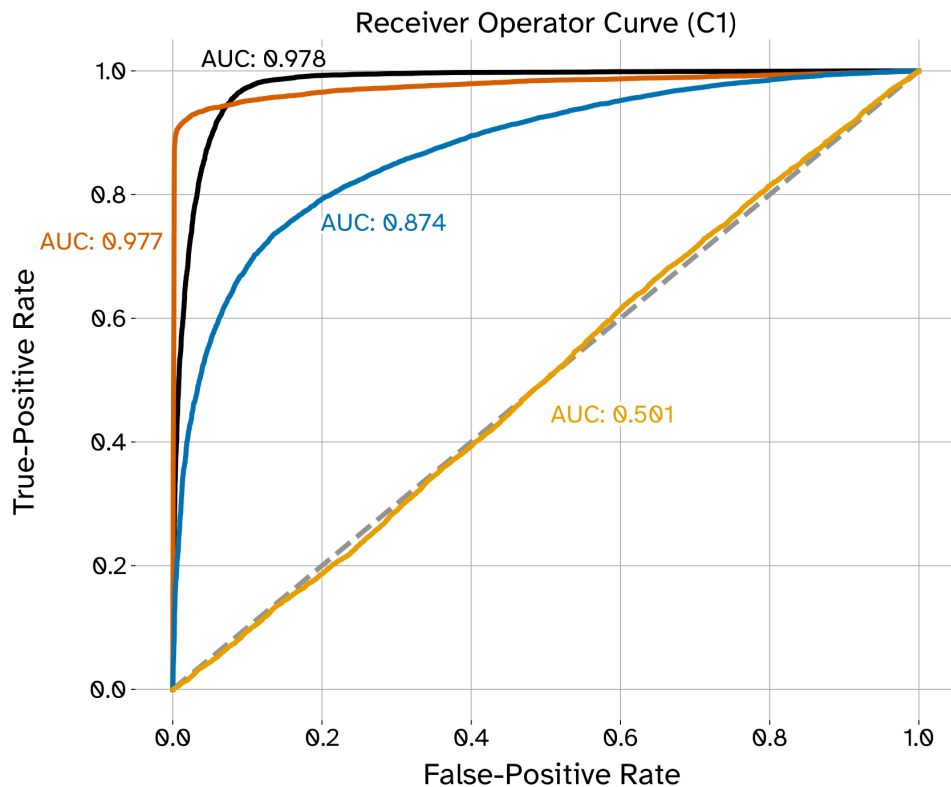
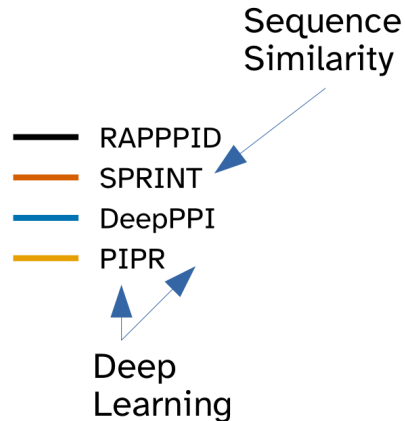
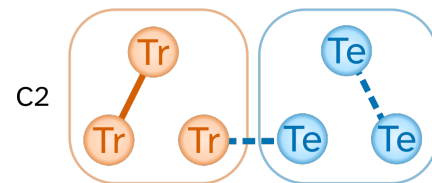
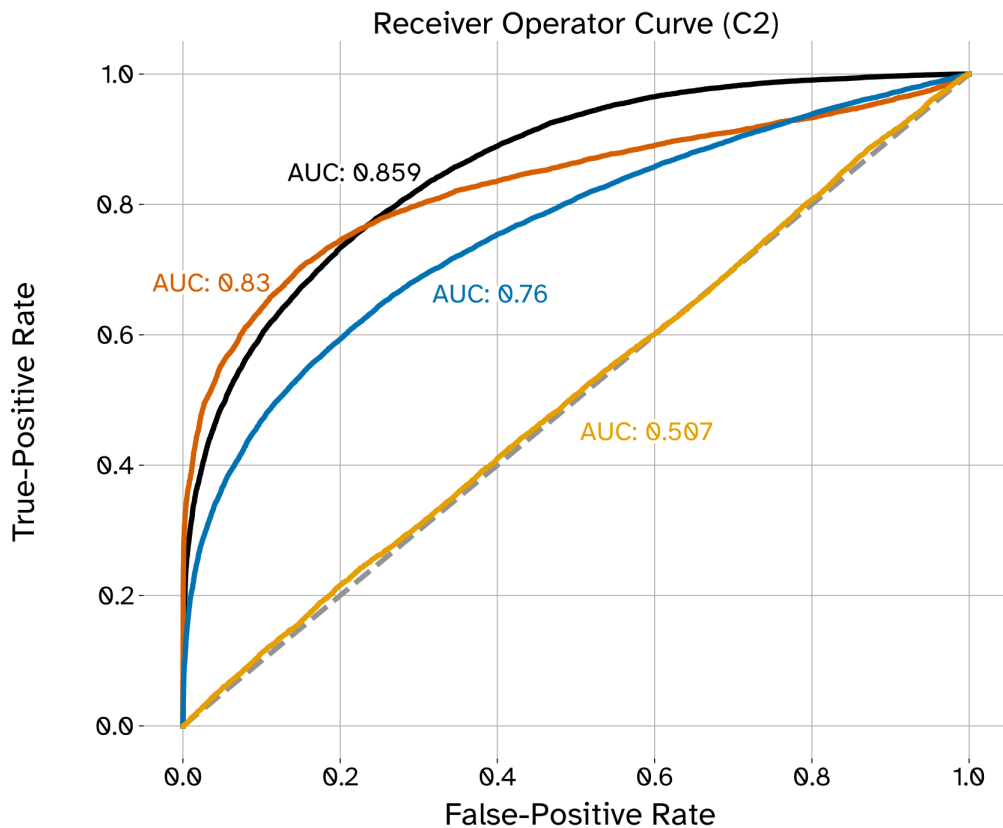


Figure 1. Illustrations of SWA and SGD with a Preactivation ResNet-164 on CIFAR-100 [1]. **Left:** test error surface for three FGE samples and the corresponding SWA solution (averaging in weight space). **Middle and Right:** test error and train loss surfaces showing the weights proposed by SGD (at convergence) and SWA, starting from the same initialization of SGD after 125 training epochs. Please see [1] for details on how these figures were constructed.

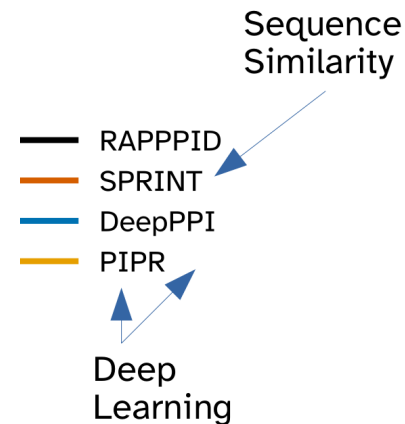
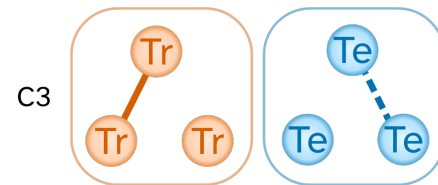
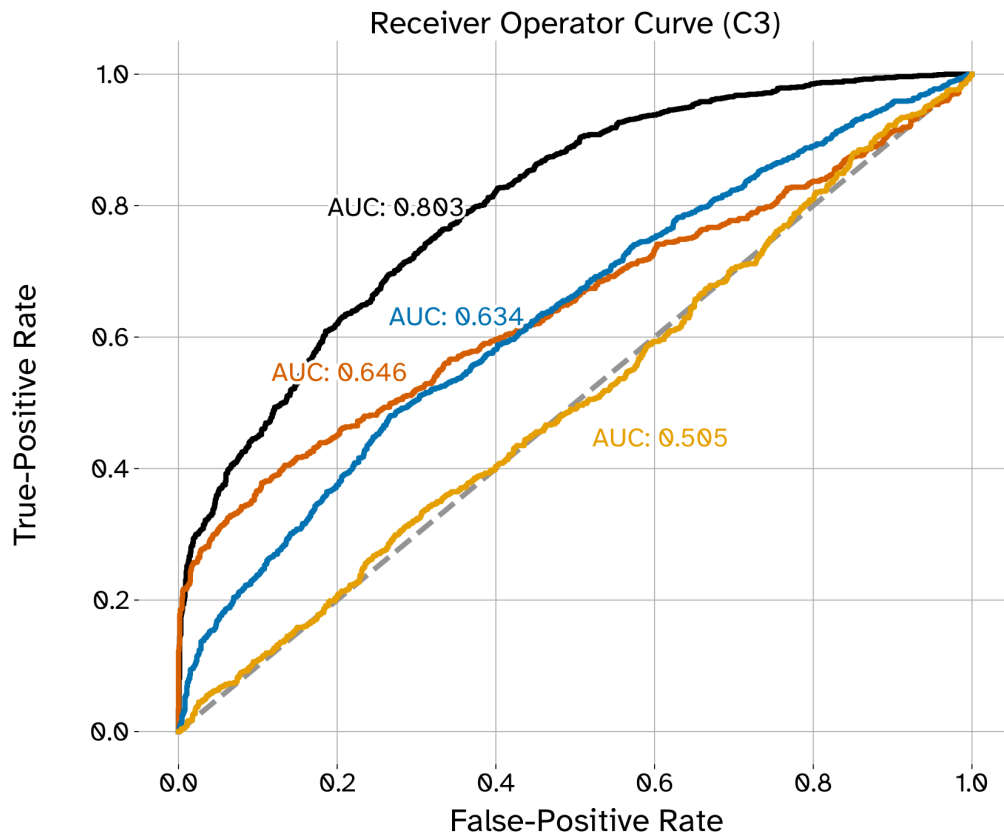
How does RAPPID perform?



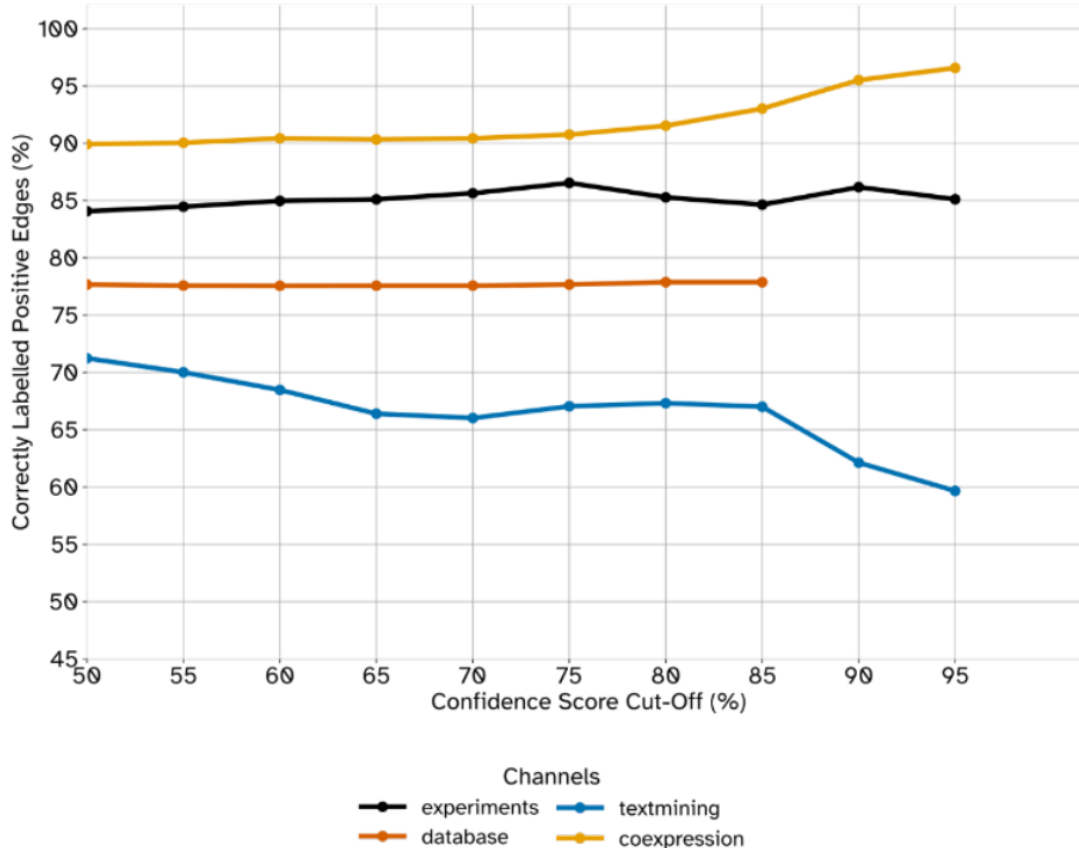
How does RAPPID perform?



How does RAPPID perform?



RAPPID performance vs. data providence



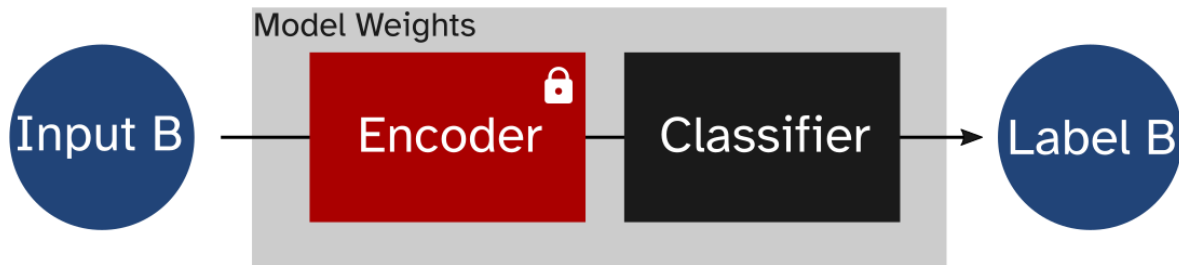
RAPPPID performance vs. data providence

Results from an ablation study conducted on RAPPPID. Each model is trained/tested twice on three randomly generated C3 datasets. The performance metrics correspond to held-out test sets.

| | RAPPPID (original) | RAPPPID -SWA | RAPPPID +Adam | RAPPPID-AWD | RAPPPID-SentencePiece | RAPPPID + TransfLG | RAPPPID + TransfSM |
|------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| Test AUROC | 0.792 (± 0.007) | 0.782 (± 0.007) | 0.791 (± 0.025) | 0.762 (± 0.020) | 0.749 (± 0.009) | 0.670 (± 0.030) | 0.747 (± 0.026) |
| AUROC Diff | N/A | -1.20% | -0.100% | -3.70% | -5.37% | -15.3% | -5.68% |
| Test APR | 0.794 (± 0.009) | 0.783 (± 0.007) | 0.792 (± 0.032) | 0.757 (± 0.022) | 0.748 (± 0.011) | 0.686 (± 0.040) | 0.758 (± 0.025) |
| APR Diff | N/A | -1.37% | -0.273% | -4.62% | -5.85% | -13.6% | -4.61% |

Transfer Learning on X-Ray Crystallography Data

- BioLIP dataset: semi-curated dataset of Protein/Ligand interactions based on the PDB
- We pretrain on STRINGDB, then fine-tune on BioLIP

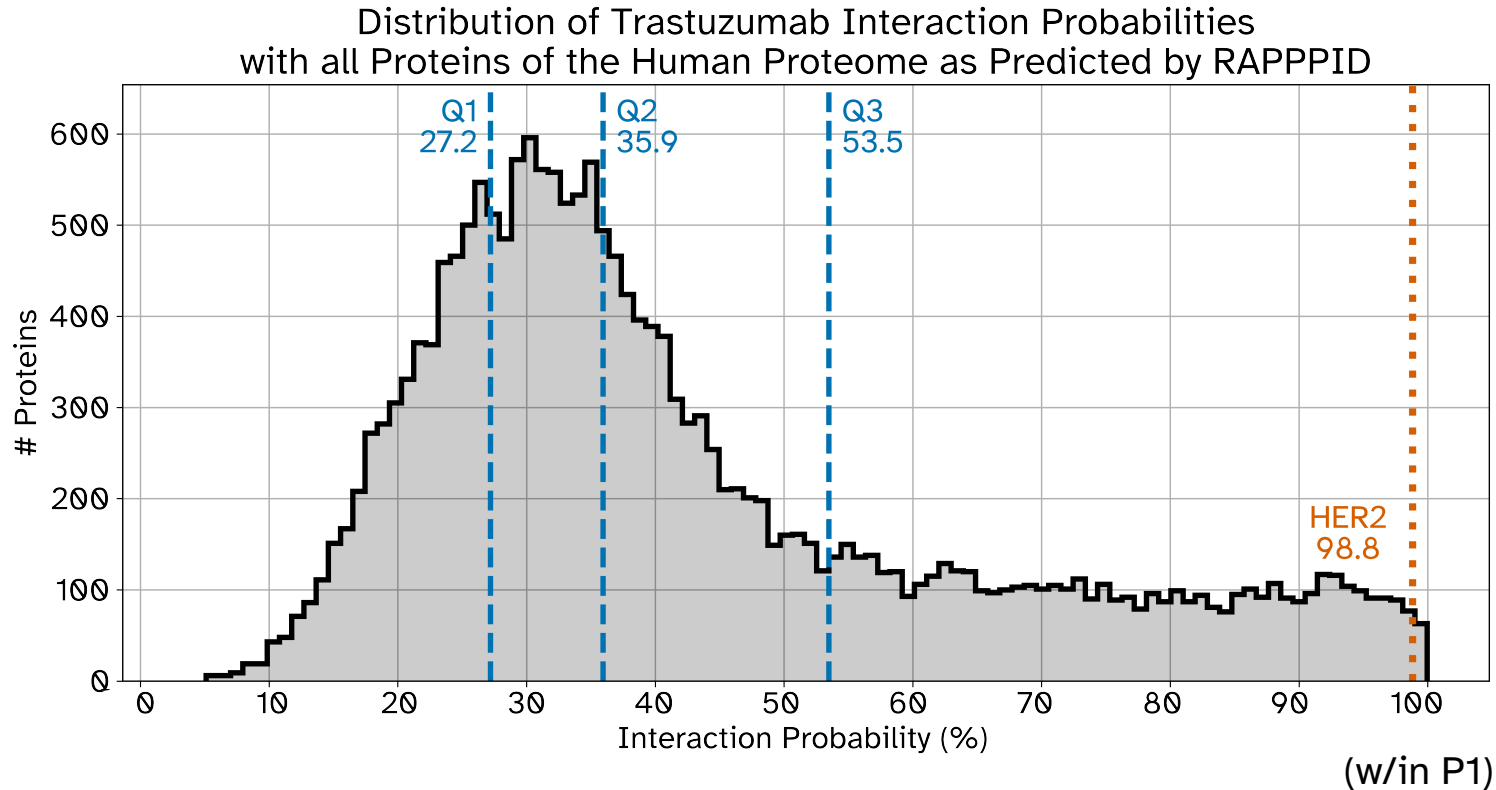


- Training on STRING DB, fine-tuning on BioLIP, and testing on BioLIP:
 - AUROC of **0.909**

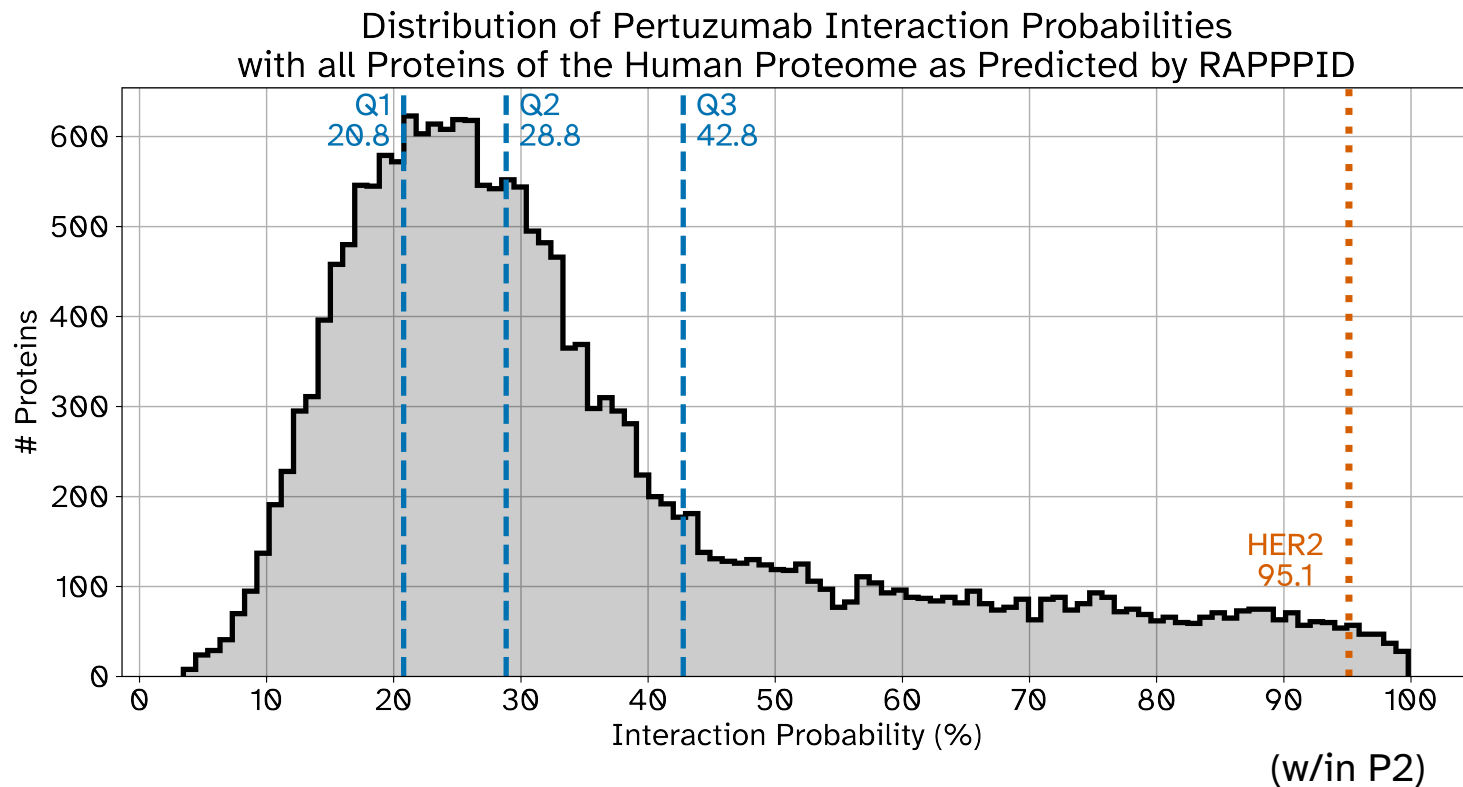
RAPPPID predicts interaction of HER2 with Trastuzumab and Pertuzumab

- How might one use RAPPPID to validate hypothesized interactions between:
 - Target proteins
 - Candidate therapeutic proteins and peptides
- Two examples: Trastuzumab and Pertuzumab.
 - Recombinant humanised monoclonal antibodies
 - Used for HER2-positive metastatic breast cancer

RAPPPID predicts interaction of HER2 with Trastuzumab and Pertuzumab



RAPPPID predicts interaction of HER2 with Trastuzumab and Pertuzumab



Acknowledgements

- Thanks to the members of the COMBINE lab for their feedback and support
 - P.D.F: Antoine Soulé
 - Ph.D.: David E. Hostallero, Ali Saberi, Yitian Zhang
 - M.Sc.: Mohamed Reda El Khili, Jessica (Yihui) Li, Chen Su, Abulrahman Takiddeen

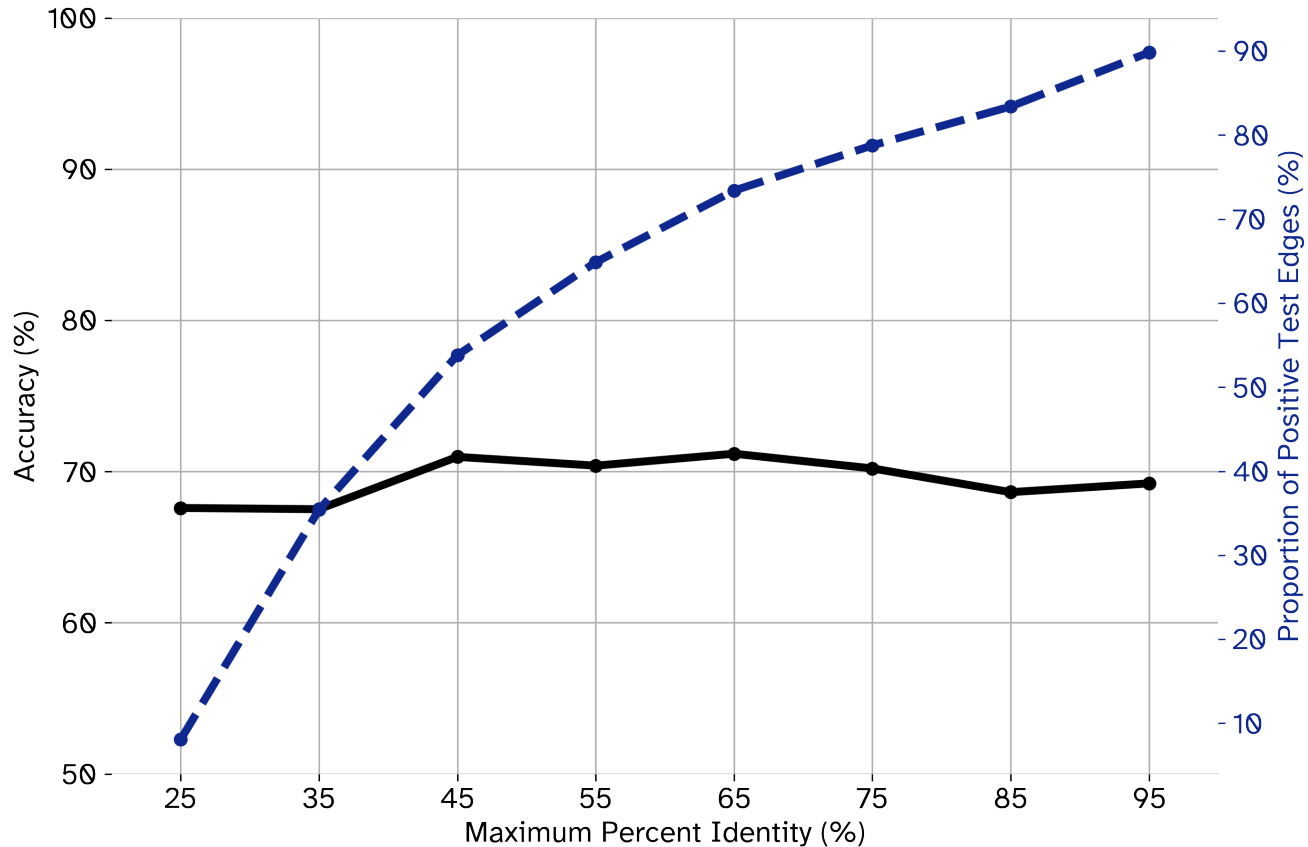
Thanks to our supporters:



Thank you

Questions?

Is RAPPPID just identifying similar sequences?



Existing PPI datasets are not great for Deep Learning.

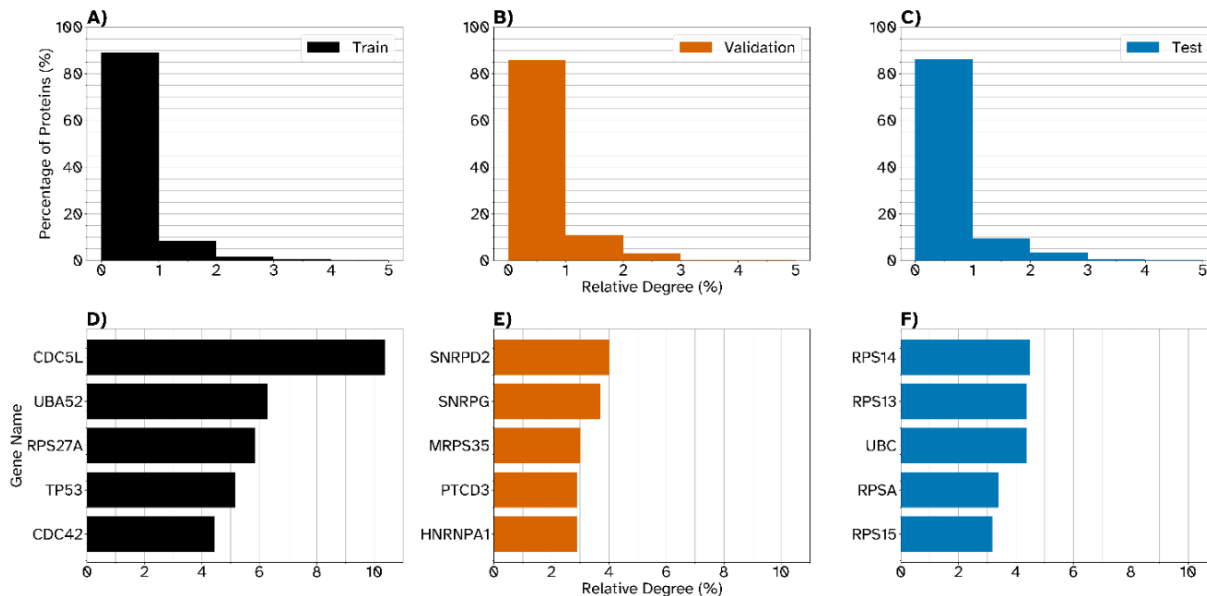
- We wanted to use additional datasets, like HIPPIE and iRefWeb
- Only STRING has enough high-confidence edges for deep learning purposes
 - 98.5% fewer edges in HIPPIE than in STRING (human, 95% confidence)
 - 87.9% fewer edges with an 85% confidence.
 - 75% fewer edges in iRefWeb than in STRING (human, 95% confidence)
- This is made worse by the fact that PPI datasets overfit terribly to begin with

False-Positive Rate

- We evaluated the false-positive rate of confidence score-filtered STRING dataset
 - We used curated and experimentally validated non-interacting protein pairs from **Negatome**
- We compared the set of proteins that are:
 - Both in STRING and Negatome
 - Evaluating the number of negative edges in Negatome that were considered a positive edge in this intersection
- Estimated the false-positive rate of our STRING dataset to be **4.01%**
- Falls within the expected **5%** upper-bound given by our 95% confidence threshold

Protein Over-Representation

- PPI graphs are understood to be scale-free in the general case
- That means that some hub proteins might be over-represented
- But that isn't the case.



Curated negative examples

- We investigated using the curated database “Negatome” for the negative samples
- There are too few (1,191 negative *H. sapiens* pairs; 263,130 positive pairs)