1    **RAPPPID: Towards Generalisable Protein Interaction Prediction with AWD-LSTM Twin Networks**

2

3

4

5    Joseph Szymborski[1,2] and Amin Emad[1,2,3,*]

[1] Department of Electrical and Computer Engineering, McGill University, Montréal, QC, Canada

[2] Mila, Quebec AI Institute, Montréal, QC, Canada

[3] The Rosalind and Morris Goodman Cancer Institute, Montréal, QC, Canada

6    * Corresponding Author:

7    Amin Emad

8    755 McConnell Engineering Building

9    3480 University Street

10    Montréal, QC, Canada, H3A 0E9

11    Email: amin.emad@mcgill.ca

1

12 **ABSTRACT**

13 **Motivation:** Computational methods for the prediction of protein-protein interactions, while

14 important tools for researchers, are plagued by challenges in generalising to unseen proteins.

15 Datasets used for modelling protein-protein predictions are particularly predisposed to

16 information leakage and sampling biases.

17 **Results:** In this study, we introduce RAPPPID, a method for the Regularised Automatic

18 Prediction of Protein-Protein Interactions using Deep Learning. RAPPPID is a twin AWD-

19 LSTM network which employs multiple regularisation methods during training time to learn

20 generalised weights. Testing on stringent interaction datasets composed of proteins not seen

21 during training, RAPPPID outperforms state-of-the-art methods. Further experiments show that

22 RAPPPID's performance holds regardless of the particular proteins in the testing set and its

23 performance is higher for biologically supported edges. This study serves to demonstrate that

24 appropriate regularisation is an important component of overcoming the challenges of creating

25 models for protein-protein interaction prediction that generalise to unseen proteins. Additionally,

26 as part of this study, we provide datasets corresponding to several data splits of various

27 strictness, in order to facilitate assessment of PPI reconstruction methods by others in the future.

28 **Availability and Implementation:** Code and datasets are freely available at

29 https://github.com/jszym/rapppid.

30 **Contact:** amin.emad@mcgill.ca

31 **Supplementary Information:** Online-only supplementary data is available at the journal's

32 website.

33

34 **INTRODUCTION**

35  Interactions of proteins with other proteins and their surroundings are fundamental to the internal

36  machinery of a cell. These interactions are of particular interest, as it is essential for a bevy of

37  diverse cellular functions: from organising cell structure to generating metabolic energy (Huttlin

38  *et al.*, 2017). These interactions are typically validated with a high degree of confidence by the

39  many biological assays commonly employed today, each with their own specific advantages and

40  challenges (Snider *et al.*, 2015). Assays for validating protein interactions range from the

41  venerable yeast two hybrid (Y2H) (Vidal and Fields, 2014) which researchers have relied on for

42  the past decades, to more recent Biotin-related techniques such as BioID-MS (Roux *et al.*, 2012).

43  A characteristic of all these assays, however, is that they are costly in terms of time, labour, and

44  materials. A further complication of protein-protein interaction studies are the multiple sources

45  of biases that plague small and large datasets alike. The choice of proteins to include in a study

46  poses a particular threat of so-called "bait bias" in smaller datasets, while large datasets suffer

47  from biases in the discovered protein interactions ("prey bias") as well as the exacerbation of

48  laboratory biases when included in aggregated datasets (Gillis *et al.*, 2014). Entire classes of

49  proteins, such as membrane proteins, which are sometimes difficult to experimentally validate

50  are often under-represented in these interaction studies as well.

51

52  Computational approaches to predict protein-protein interactions (PPIs) are therefore useful to

53  help towards reducing the number of costly experiments researchers are required to perform.

54  Researchers have deployed many diverse approaches to solve the task of protein sequence-based

55  interaction prediction. Most sequence-based methods rely on the understanding that co-evolution

56  and co-expression of proteins are both tied to protein interaction and sequence similarity (Cong

3

57    *et al.*, 2019; Jansen, 2003). Some methods rely on substitution matrices for sequence alignment

58    such as BLOSUM or PAM in combination with machine learning methods to predict interactions

59    (Henikoff and Henikoff, 1992; Ding *et al.*, 2016). Other methods utilise Support Vector

60    Machines (SVMs) with kernels specifically designed for use with protein sequences (Ben-Hur

61    and Noble, 2005). Other statistical methods including naïve bayes (NB) and *k*-nearest neighbors

62    (kNN) have been used to predict protein interactions from protein sequences (Browne *et al.*,

63    2007). Some of the most successful PPI prediction methods belong to the family of methods that

64    rely on novel substring search algorithms (Li and Ilie, 2017; Dick *et al.*, 2020), operating

65    similarly to sequence search tools like BLAST (Altschul *et al.*, 1990). Deep learning models

66    have also been designed for predicting protein interactions (Chen *et al.*, 2019). These deep

67    approaches commonly either learn wide networks or share weights in a twin design; the latter

68    being shown to be both more efficient and effective (Richoux *et al.*, 2019).

69

70    Such methods, however, face many challenges due to the nature of the data on which they train.

71    Arguably the most pervasive of which is the ability of models to generalise and predict the

72    interactions of proteins previously unseen by the prediction method. To ensure such

73    generalisability, careful cross-validation techniques must be used to avoid data leakage. While

74    the necessity of appropriate cross-validation techniques is not unique to this area of research, the

75    application of these networks (e.g., PPIs, transcriptional regulatory networks) to obtain

76    biological insights makes it particularly important to address this challenge in the task of

77    biological network reconstruction (Park and Marcotte, 2012; Tabe-Bordbar *et al.*, 2018).

78

4

79    For PPI reconstruction, developing generalizable models prove particularly difficult. The nature

80    of PPI networks makes it easy to create datasets with testing/training splits which leak

81    information, resulting in inflated performance metrics that cannot properly assess the

82    generalisability of these methods. In particular, simply splitting interaction datasets into training

83    and testing sets using random selection of edges results in the construction of testing datasets that

84    are almost entirely comprised of interactions between proteins found in the training set. Indeed,

85    Park and Marcotte in 2012 found that all the PPI prediction models they surveyed were tested on

86    such naïvely constructed datasets (Park and Marcotte, 2012). The surveyed models, which had

87    optimised their performance on these naïve datasets, that suffered from a large degree of

88    information leakage, were also found to incur precipitous falls in their prediction metrics when

89    tested on datasets where no proteins in the testing set occurred in the training dataset.

90

91    When faced with obstacles in the construction of generalisable models, strategic and targeted

92    applications of regularisation at training time can significantly improve results. This is

93    particularly relevant in the context of deep learning methods which pay for their expressiveness

94    by learning an outsized number of parameters. Among the most common regularisation

95    techniques used in the deep learning context is "dropout" (Srivastava *et al.*, 2014). Applying

96    dropout to a layer consists of randomly zeroing the activations of the previous network with

97    some probability $p$. However, choice of regularisation techniques must be selected with care and

98    in accordance with the architecture. For example, applying dropout directly to the hidden state of

99    Recurrent Neural Networks (RNNs) (Lipton *et al.*, 2015), an architecture which lends itself

100   naturally to sequential inputs such as amino acid chains, impairs its ability to retain its memory

101   of previous inputs (Zaremba *et al.*, 2015).

5

102   Applying regularisation to RNNs requires additional considerations. Recent work by Merity *et*

103   *al*. has demonstrated that randomly zeroing the weights of RNNs ("dropconnect") (Wan *et al.*,

104   2013) rather than their hidden state activation ("dropout") effectively reduces testing error

105   (Merity *et al.*, 2017). Merity *et al.* describe applying dropout to the embedding layer as well as

106   using an averaged optimiser (NT-ASGD) as part of a series of regularisation techniques dubbed

107   Averaged Weight-Dropped Long Short-Term memory (AWD-LSTM). The regularisation

108   techniques used by AWD-LSTM models are specifically selected for their suitability in the

109   context of training RNNs.

110

111   To meet the generalisation challenges posed by PPI prediction tasks, we developed a method

112   called the Regularised Automatic Prediction of Protein-Protein Interactions using Deep

113   Learning, or RAPPPID. RAPPPID addresses the challenges in creating generalised models for

114   PPI prediction by adopting (with modification) the AWD-LSTM, a regularised recurrent neural

115   network training routine (Merity *et al.*, 2017). In this study, we showed that RAPPPID

116   outperforms state-of-the-art PPI prediction methods on strict validation datasets constructed in

117   accordance with guidelines set out by Park & Marcotte (Park and Marcotte, 2012). Additionally,

118   we performed various analyses and applied RAPPID to different use-cases to show-case the

119   effect of different components of its architecture and training on its performance and its

120   applicability to different real-world scenarios.

121

122   It is worth mentioning that in addition to RAPPPID's code, we have made various pre-processed

123   datasets freely available in the hope that they facilitate the evaluation of protein interaction

6

124 prediction methods on datasets that both mitigate information leakage and are appropriately large

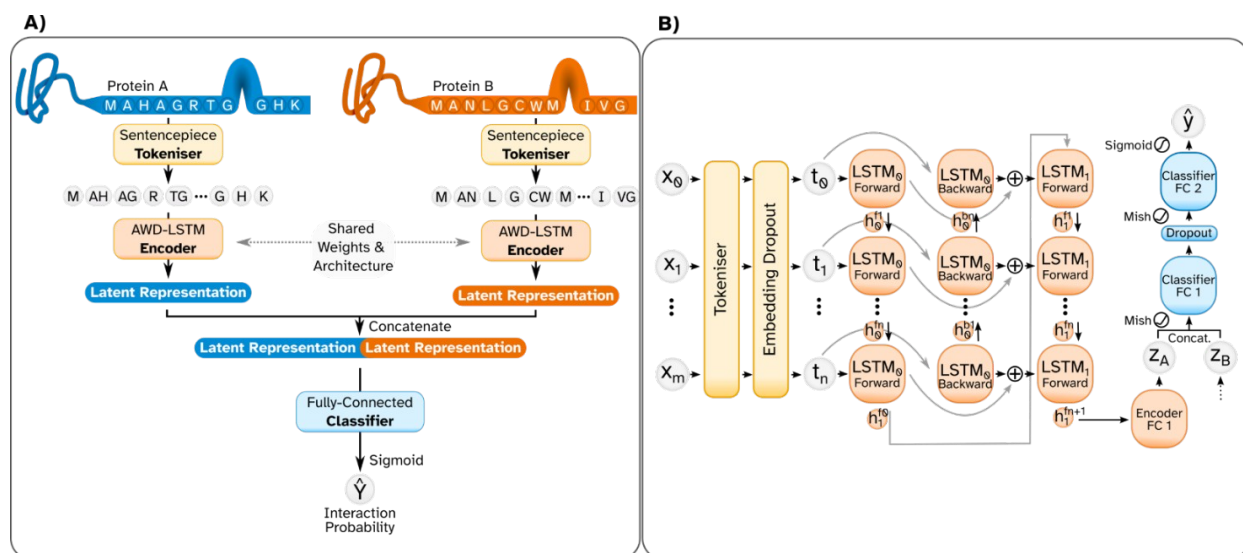125 for deep learning applications.

126

127

128 **METHODS**

129 **PPI prediction using AWD-LSTM Twin Networks**

130 RAPPPID is first trained by considering pairs of amino acid sequences of proteins along with a

131 label indicating whether they do or do not interact. The amino acid sequences are first tokenised

132 using the Sentencepiece algorithm (Kudo and Richardson, 2018), which allows for better

133 recognition of common groupings of amino acid residues that make up the secondary structure

134 and motifs of proteins. Fixed-length latent vector representations of the token sequences of both

135 proteins are then computed by twin neural networks (Bromley *et al.*, 1993), forming the encoder

136 of RAPPPID's pipeline. These twin networks have shared architectures and weights and are

137 trained jointly. An overview of the pipeline is provided in Figure 1.

138

139 Each twin network consists of a two-layer bidirectional AWD-LSTM network (Merity *et al.*,

140 2017), which takes a tokenised amino acid sequence as its input and generates a fixed-length

141 latent vector representation as its output. AWD-LSTMs are architecturally identical to the

142 LSTMs (Hochreiter and Schmidhuber, 1997), however at inference time, several regularisation

143 techniques are employed while training to promote learning generalised weights. Among these

144 regularisation techniques are an averaged optimizer and dropout applied to embeddings and

145 LSTM weights (Athiwaratkun *et al.*, 2019; Wan *et al.*, 2013). Using an AWD-LSTM encoder

146  enables RAPPPID to leverage the strong inductive biases of LSTMs, while ensuring that the

147  learned weights are generalised.



148

149  **Figure 1: Overview of the RAPPPID pipeline and architecture.** (A) The pipeline of the
150  RAPPPID begins with two protein sequences. Each sequence is first tokenised by the
151  Sentencepiece tokeniser (yellow). Each sequence of tokens is then inputted separately into an
152  AWD-LSTM encoder layer which results in a latent representation for each protein. The latent
153  representations of each protein are inputted into a fully-connected classifier layer. The classifier
154  layer outputs the predicted probability of the two proteins interacting. (B) Taking a closer look at
155  the architecture of the RAPPPID, individual residues $x_0$ to $x_m$ for a protein comprised of $m$
156  residues are tokenized into $n$ tokens $t_0, t_1, \ldots, t_n$. The embedding dropout layer randomly assigns
157  random tokens from the total vocabulary to zero. The encoder layer is comprised of a multi-layer
158  bidirectional LSTM whose last hidden state is fed to a fully connected layer before outputting a
159  latent representation $z_A$. $z_A$ is then concatenated with the latent representation of a second protein
160  ($z_b$) before being inputted into the two-layer fully-connected classifier. The output of the
161  classifier is activated by the sigmoid function to produce a probability of interaction.

162

163  The final hidden state of the AWD-LSTM is passed to a single fully-connected layer whose

164  output, once activated by the Mish function (Misra, 2020), is the latent representation of the

165  amino acid sequences. The latent representations of both proteins are then concatenated and

166  provided as inputs to the classifier network which generates an interaction probability for each

167  protein pair. The classifier network is a two-layer fully-connected network that outputs a single

8

168 logit whose sigmoid activation serves as the probability of the two proteins interacting. The

169 activated logit is then used to calculate the mean binary cross-entropy loss. Relegating the

170 pairwise comparison of proteins to the shallower classifier network allows RAPPPID to infer

171 protein interactions in an efficient manner. Figure 1 provides an overview of RAPPPID's

172 pipeline.

173

174

175 **Sequence Segmentation and Tokenisation**

176 As mentioned earlier, RAPPPID utilises the Sentencepiece algorithm (Kudo and Richardson,

177 2018) to tokenize amino acid sequences. While words form the basis of many natural languages

178 and may break up sentences and phrases into discrete units, no such higher-order segmentation is

179 as immediately apparent in amino acids. Motifs and protein domains possess many analogous

180 qualities to words in natural languages; they appear repeatedly in amino acid sequences and their

181 combination and relative position in these sequences play important roles in the protein structure

182 and function (Anfinsen, 1973).

183

184 Much like words, however, motifs and protein domains present an "out-of-vocabulary" problem,

185 where unseen examples are difficult to handle. Attempts to solve this problem in natural

186 language processing tasks has resulted in "subword" segmentation algorithms, particularly in

187 difficult-to-segment languages such as Japanese which do not separate words by spaces (Schuster

188 and Nakajima, 2012). Here, we employ the Sentencepiece algorithm (Kudo and Richardson,

189 2018) to sample tokens from "subword" vocabularies generated by the Unigram algorithm

190 (Kudo, 2018). The unigram and sentencepiece algorithms construct vocabularies of arbitrary size

9

191   by modelling the probabilities of subwords and provides a principled manner for sampling from

192   this distribution to reconstruct sequences. The multi-residue tokens that comprise the vocabulary

193   subdivide low-entropy areas and reduce the overall length of the sequences encoded.

194

195

196   **Generalising Protein Sequence Encoding with AWD-LSTM**

197   The task of protein-protein interaction prediction on unseen proteins is a difficult problem prone

198   to overfitting, as demonstrated by the poor testing performance of various methods on unseen

199   proteins (Park and Marcotte, 2012). For this reason, a training and optimisation methodology that

200   allows efficient regularisation is desirable. AWD-LSTM was recently devised to enable efficient

201   training of generalisable recurrent neural networks (RNNs) (Merity *et al.*, 2017). This approach

202   deploys several regularisation techniques during training to achieve this goal. RAPPPID adopts,

203   with modification, the training methodology of AWD-LSTM.

204

205   RAPPPID utilises Embedding Dropout, DropConnect (Wan *et al.*, 2013), and Weight Decay

206   (Loshchilov and Hutter, 2019) on the LSTM weights, as described by AWD-LSTM, in the

207   encoder during training. Both AWD-LSTM and RAPPPID optimise over average weights,

208   however the optimisers are quite different. AWD-LSTM makes use of the non-monotonically

209   triggered averaged stochastic gradient descent (NT-SGD) optimiser which switches between

210   stochastic gradient descent (SGD) and the averaged variant (ASGD). RAPPPID uses the recent

211   Stochastic Weight Averaging (SWA) strategy in combination with the Ranger21 optimiser

212   (Athiwaratkun *et al.*, 2019; Wright and Demeure, 2021).

213

10

214 SWA has been shown to promote generalisable models in part by overcoming the challenges of

215 finding best solutions within flat loss basins (Izmailov *et al.*, 2019). Ranger21 is an optimiser that

216 applies the "Lookahead mechanism" (Zhang *et al.*, 2019) to the AdamW optimiser (Loshchilov

217 and Hutter, 2019) and includes several optimisation techniques (Wright and Demeure, 2021).

218 These techniques, which include gradient centralisation and adaptive gradient clipping, enable us

219 to further improve our ability to learn generalised weights and smooth training trajectories

220 (Brock *et al.*, 2021; Yong *et al.*, 2020). Finally, since RAPPPID does not rely on the timestep

221 outputs of the LSTM network, the Temporal Regularisation (TAR) described by AWD-LSTM is

222 not applicable.

223

224

225 **Details of RAPPPID's architecture and hyperparameter tuning**

226 The dimensionality of the AWD-LSTM hidden state is made equal to the dimensionality of the

227 embeddings (*i.e.*, 64). The output of the fully-connected layer is of equal dimensionality to that

228 of the AWD-LSTM hidden state and is activated by the Mish activation function (Misra, 2020).

229 The output of the first fully-connected layer is half the size of the embedding dimension,

230 activated by the Mish function, and regularised by a dropout layer (Srivastava *et al.*, 2014).

231 RAPPPID trains on a vocabulary of 250 tokens that is generated by the Sentencepiece algorithm.

232 New vocabularies are generated for each dataset before training.

233

234 The number of LSTM layers, L2 coefficient, and various dropout rates are defined as hyper-

235 parameters that require tuning between different datasets. Hyper-parameters were selected by

236 cross-validation. The range of considered hyperparameters and their selected values are provided

11

237    in the Supplementary Tables S1-S2 (in Supplementary File 1). The chosen hyper-parameter

238    ranges are a compromise between common, reasonable values and maintaining a manageable

239    hyper-parameter space which can be explored practically.

240

241

242    **Protein Interaction Datasets**

243    We obtained protein-protein interactions (PPIs) and protein sequences from version 11 of the

244    Search Tool for the Retrieval of Interacting Genes/Proteins (STRING) database (Szklarczyk *et*

245    *al.*, 2019), from the official STRING website. Edges were downloaded from https://stringdb-

246    static.org/download/protein.links.detailed.v11.0.txt.gz and sequences were downloaded from

247    https://stringdb-static.org/download/protein.sequences.v11.0.fa.gz. In this dataset, the association

248    between any two proteins is assigned a confidence score depending on the source of the

249    information (called a "channel"). In our analyses, only associations with a combined STRING-

250    score above 95% (equivalent to a score above 950, obtained from combining different channels)

251    were considered as positive edges. We also retained the channel-specific scores for further

252    analysis.

253

254    To obtain an estimate of the false-positive rate of the confidence score-filtered STRING dataset,

255    we leveraged the curated and experimentally validated non-interacting protein pairs from the

256    Negatome dataset (Blohm *et al.*, 2014). By comparing the set of proteins that are in both

257    STRING and Negatome, and evaluating the number of negative edges in Negatome that were

258    considered a positive edge in this intersection, we estimated the false-positive rate of our

12

259  STRING dataset to be 4.01%. This false-positive rate is within the expected 5% upper-bound

260  given by our 95% confidence threshold.

261

262  PPI graphs are understood to be scale-free in the general case. This property of PPI graphs can

263  make them challenging datasets upon which to train a generalised model, as some proteins can

264  be over-represented. To characterise the extent to which proteins are represented in the dataset,

265  we calculated the distribution of the relative degree of proteins in the network (Supplementary

266  Figure S1A-C). The vast majority of proteins (upwards of 85%) have edges with fewer than 1%

267  of proteins within their dataset split (*i.e.*, train/validation/test). The protein with the highest

268  degree is CDC5L which has a relative degree of just over 10%, while the second highest has a

269  relative degree of 6.2% (Supplementary Figure S1D-F).

270

271

272  **Negative Examples**

273  The preparation of datasets of pairs of proteins which are known not to interact with one another

274  is a fraught process. Various methods are typically deployed to create such datasets, which are

275  integral for machine learning methods as they typically require negative examples. Ideally, only

276  negative examples are entirely composed of non-interacting protein pairs which are

277  experimentally verified and manually curated, such as the Negatome database (Blohm *et al.*,

278  2014).

279

280  Unfortunately, non-interacting pairs are difficult to experimentally validate. As a result, there are

281  far fewer *H. sapiens* negative protein pairs in such datasets such as Negatome than positive

13

282    protein pairs in databases such as STRING. Precisely, there are 1,191 negative *H. sapiens* pairs

283    in Negatome, and 263,130 positive pairs above a 95% confidence threshold in STRING (only

284    pairs comprised of proteins present in UniprotKB were included in this analysis). The resulting

285    class imbalance is detrimental to learning performant, generalisable models.

286

287    Synthetic negative pairs are often used to compensate for the small number of experimentally

288    verified non-interacting pairs is construction of synthetic negative pairs. A common method for

289    constructing negative pairs is to select pairs of proteins from distinct sub-cellular compartments.

290    This method has unfortunately been shown to result in biased samples according to multiple

291    measures (Ben-Hur and Noble, 2006). Selecting pairs at random from the space of pairs not

292    known to interact has proven to evade such biases, but runs the risk of capturing yet unknown

293    false-negatives. Unless otherwise noted, RAPPPID utilises synthetic random pairs of proteins

294    which are not known to interact.

295

296

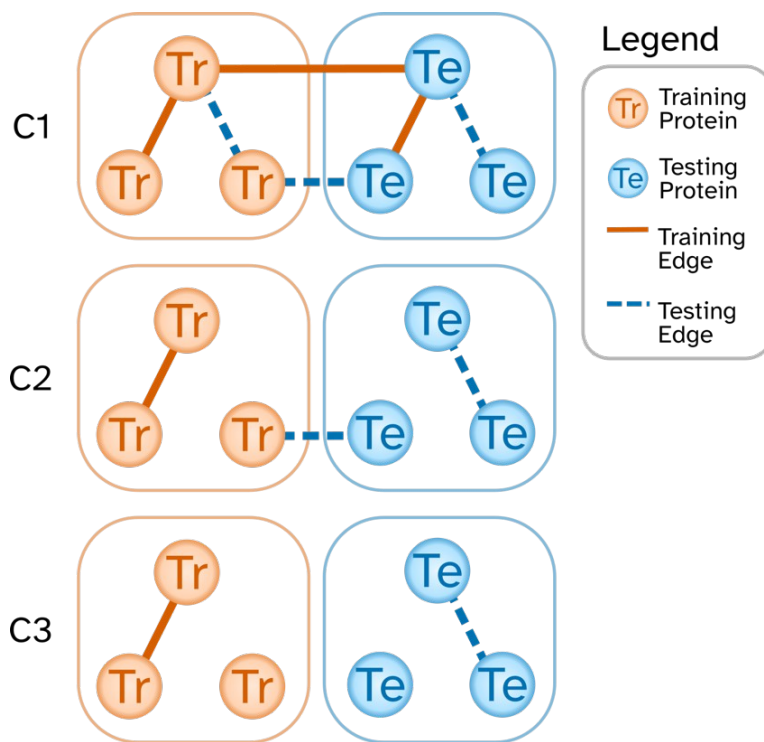297    **Training, Validation and Testing Set Construction**

298    As identified by Park and Marcotte (Park and Marcotte, 2012), methods which consider the

299    interaction of proteins in a pairwise fashion must (and have historically failed to) take additional

300    care to avoid information leakage when constructing training and testing datasets. Following

301    their suggestion, here we use three different classes of testing and training sets to evaluate the

302    performance of RAPPPID and other algorithms (Figure 2).

303

14

304  1) "C1" refers to the evaluation scheme in which edges (pairs of proteins) are randomly selected

305  to form the training or testing sets. Since the selection criterion is based on edges, both proteins

306  in the pair may be present in both the testing and training sets (due to the presence of other edges

307  adjacent to each protein).

308  2) "C2" refers to the evaluation scheme in which proteins are randomly selected to form the

309  training or testing sets. In this scheme, only one protein in a pair may be present in both testing

310  and training sets (but never both). This evaluation scheme mimics the scenario in which a model

311  trained on the interactome is used to predict the interaction of known proteins with a newly

312  discovered protein (that was not used to train the model).

313  3) "C3" refers to the evaluation scheme in which proteins are randomly selected to form the

314  training or testing sets. However, unlike C2, proteins which appear in the training set never

315  appear in the testing set. This is the most strict evaluation scheme.

**Figure 2: Illustration of differences in edges between C1, C2, and C3 datasets.** Differences between C1, C2, and C3 datasets are most visible by first dividing the population of all proteins in the dataset into training (orange, left) and testing (blue, right). In the case of the strict C3 dataset (bottom row), edges known at training time (orange, solid) only occur between training proteins. Similarly, C3 datasets are evaluated on testing edges (blue, dotted) that only occur between testing proteins. The C2 dataset has all the edges present in the C3 dataset, but also includes testing edges between testing proteins and training proteins. Finally, the pervasive C1 datasets allows all possible training and testing edges that are not identical.

While most methods have typically reported on models validated with datasets in the C1 class, they often perform much worse on similar datasets in the more conservative C2 and C3 class. This is likely due to the information leakage between training and testing sets present in the C1 class and, to a lesser extent, in the C2 class. In our evaluations, we report the performance of different methods using all three evaluation schemes above, but we are most interested in the results of C3 due to a lack of information leakage.

16

333    In the C-type datasets we've created, there are approximately 9.3 thousand proteins. The number

334    of edges in the datasets range from just over 263 thousand edges in the C1 dataset to just under

335    174 thousand edges in the C3 dataset. Full dataset statistics are presented in Supplementary

336    Table S3. These C-type datasets are made freely available to the public, with instructions on how

337    to download them at https://github.com/jszym/rapppid/tree/main/data. We've made these

338    datasets available in the hope that it facilitates the evaluation of protein interaction prediction

339    methods on datasets that both mitigate information leakage and are appropriately large for deep

340    learning applications. The adoption of C-type datasets for training and evaluation of PPI

341    prediction methods is important for accurate and representative benchmarking.

342

343

344    **Implementation**

345    RAPPPID was implemented in the Python computer language using the PyTorch and PyTorch

346    Lightning deep learning framework (Paszke *et al.*, 2019; Falcon *et al.*, 2020). Embedding

347    Dropout and DropConnect implementation were obtained from the AWD-LSTM code base

348    (Merity *et al.*, 2017). The source code for RAPPPID can be found by visiting

349    https://github.com/jszym/rapppid. RAPPPID was trained at a rate of approximately 2.7 training

350    steps per second at a batch size of 80 protein pairs on an NVIDIA RTX 2080 GPU and 32 CPU

351    cores clocked at 2.2 GHz. The C3 model takes approximately 2.4 hours to train, while C1 and C2

352    take approximately 7.8 hours.

353

354

355

17

356 **Protein Similarity Experiments**

357 In the analysis of the sequence similarity between testing and training proteins, "Percent

358 Identity" was measured between proteins using NCBI's PSI-BLAST tool running locally as part

359 of version 2.12.0 of NCBI's BLAST+ software suite (Altschul *et al.*, 1997). In addition to the

360 Percent Identity, for two proteins to be considered similar, an E-value cut-off of at most 5 and an

361 alignment length of more than 30% of the query sequence were considered necessary. The 64-bit

362 Linux binaries of the BLAST+ suite were obtained from the link

363 ftp.ncbi.nlm.nih.gov/blast/executables/blast+/2.12.0/.
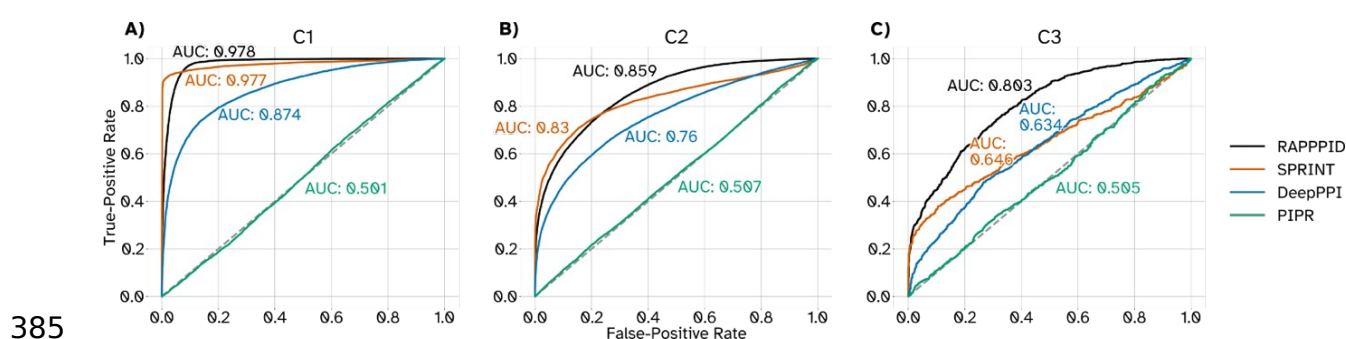
364

365

366 **Results**

367 **Performance evaluation of RAPPPID and other algorithms**

368 To establish the ability of RAPPPID to correctly predict protein-protein interactions within the

369 current landscape of PPI prediction methods, we compared it against three recent methods

370 (Figure 3). The first of these is the Scoring Protein INTeractions (SPRINT) method, which

371 belongs to the family of methods that predict interactions according to measures of sequence

372 similarity (Li and Ilie, 2017). SPRINT was shown to outperform support-vector machine

373 (SVM), random-forest (RF), and sequence similarity-based methods across C1, C2, and C3-like

374 datasets.

375

376 The two other methods, PIPR and DeepPPI, are deep learning methods that similar to RAPPPID

377 utilize twin networks (Chen *et al.*, 2019; Richoux *et al.*, 2019). PIPR uses a residual recurrent

378 convolutional neural network (RCNN) for its encoder with the goal of more effectively

18

379   summarising both local and global features. We compared RAPPPID against the best performing

380   iteration of DeepPPI, whose encoder comprises of a convolutional neural network feature

381   extractor followed by an LSTM network. All three methods required retraining, as the available

382   weights of these models belong to diverse datasets, all of which are either not C2 or C3 type, or

383   (in the case of SPRINT) are insufficiently large for training a generalisable deep learning model.

384



385

386   **Figure 3. Receiver-Operator curves across methods and datasets.** The receiver-operator
387   curves (ROCs) for all four methods tested across C1 (A), C2 (B), and C3 (C) datasets.

388

389   Across C1, C2, and C3 testing datasets, RAPPPID achieved higher area under the receiver-

390   operator curve (AUROC) than all other methods tested (Table 1). The margin between RAPPPID

391   and the second highest performing method (SPRINT in all cases) was highest when performed

392   on the stricter C3 dataset, resulting in approximately a 24.3% improvement. The improvement

393   obtained by RAPPPID compared to SPRINT was lower on the C2 dataset (approximately 3.4%),

394   and finally nearly equivalent on the least strict C1 dataset.

395

396   With regards to the area under the precision-recall curve (AUPR), this trend across dataset types

397   persisted. RAPPPID's AUPR was higher than all other methods for the C3 dataset, with a margin

19

398    to the second highest method of 0.094 (equivalent to an approximately 14.6% improvement).

399    RAPPPID's AUPR score was matched by SPRINT in experiments conducted on the C2 dataset,

400    but outperformed DeepPPI and PIPR. SPRINT achieved the highest AUPR of all the methods in

401    the C1 dataset, outperforming RAPPPID with a margin of 0.009 (equivalent to an approximately

402    0.9% improvement).

403    **Table 1: Comparison of PPI prediction performance on C1, C2, and C3 datasets.** The
404    testing AUROC and AUPR of four different PPI prediction methods is reported across the three
405    different dataset types described by Park & Marcotte (Park and Marcotte, 2012).

| Dataset | Method | Testing AUROC | Testing AUPR |
|---|---|---|---|
| C1 | RAPPPID | **0.978** | 0.974 |
| | SPRINT | 0.977 | **0.983** |
| | DeepPPI | 0.874 | 0.881 |
| | PIPR | 0.501 | 0.405 |
| C2 | RAPPPID | **0.859** | **0.868** |
| | SPRINT | 0.830 | **0.868** |
| | DeepPPI | 0.760 | 0.787 |
| | PIPR | 0.507 | 0.508 |
| C3 | RAPPPID | **0.803** | **0.810** |
| | SPRINT | 0.646 | 0.716 |
| | DeepPPI | 0.574 | 0.590 |
| | PIPR | 0.505 | 0.509 |

406

407    While we were able to replicate results of the PIPR model on the *S. cerevisiae* dataset published

408    as part of the original PIPR publication (Chen *et al.*, 2019), PIPR suffered from convergence

409    issues during training on our *H. sapiens* STRING datasets. We suspect this is due to a variety of

410    factors, with the large differences in dataset characteristics being the most likely cause. The

411    number of proteins and interactions in the *S. cerevisiae* dataset is far smaller than our STRING

20

412  datasets. Furthermore, the *S. cerevisiae* dataset selected pairs of proteins which occupy different

413  subcellular compartments in order to construct negative samples; an approach found to lead to

414  biased estimations of prediction accuracy (Ben-Hur and Noble, 2006).

415

416  Taken together, these results suggest that RAPPPID outperforms alternative methods in the

417  majority of evaluations on C1, C2, and C3 schemes and is particularly effective in the stricter and

418  more difficult C3 evaluation.

419

420

421  **Channel-specific performance of RAPPPID**

422  The STRING database, integrates and annotates protein association data from a wide range of

423  sources. The "database", "text-mining", "experiments", and "coexpression" channels make-up

424  the majority of the edges in our datasets (e.g., 98.4% of all the edges in the C3 dataset).

425

426  The "database" channel is comprised of several curated databases of interactions such as KEGG

427  and Reactome (Kanehisa, 2000; Jassal *et al.*, 2019). Edges in the "text-mining" channel are the

428  result of a statistical analysis of proteins whose names and/or identifiers co-occur in publications.

429  The "experiments" channel is populated by interactions evidenced by high-throughput

430  experiments curated by members of the International Molecular Exchange (IMEx) consortium

431  (Orchard *et al.*, 2012). This includes datasets such as IntAct, DIP, the BioGRID, and many

432  others (Orchard *et al.*, 2014; Salwinski *et al.*, 2004; Oughtred *et al.*, 2020). Finally, the

433  "coexpression" channel arises from proteomic and transcriptomic assays which quantify gene-
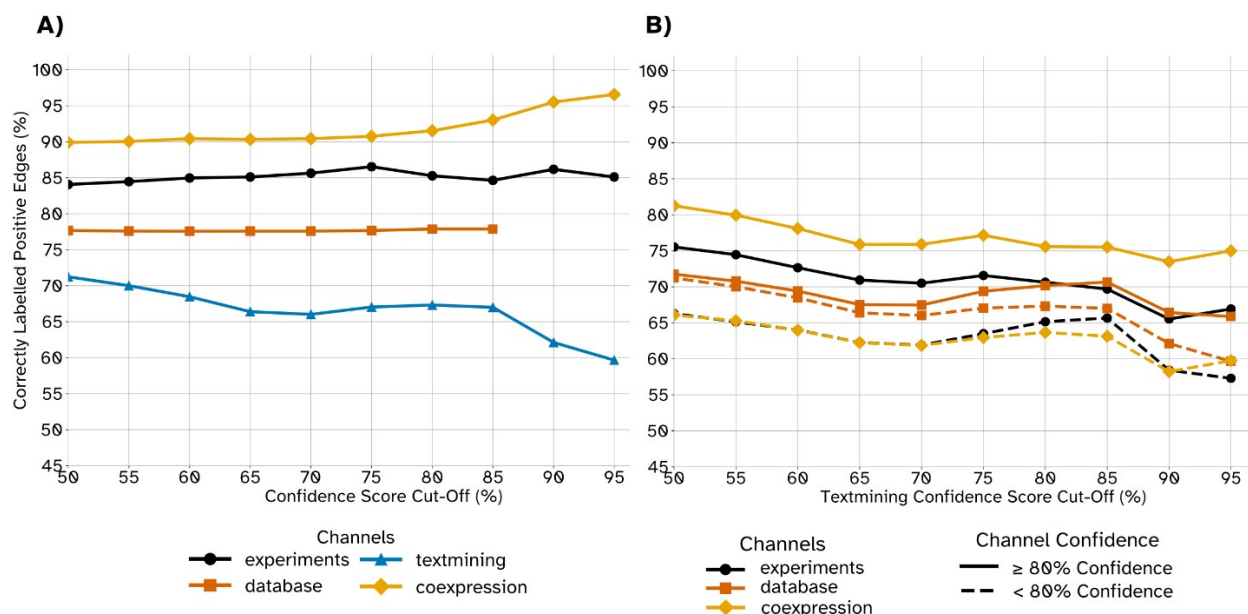
434  by-gene correlations. STRING additionally assigns calibrated confidence scores to each of the

21

435  edges which summarise the evidence supporting an edge. Scores are assigned to each edge by

436  channel, and finally harmonised into a final "combined confidence score" which represents the

437  evidence across all channels present in STRING."

438

439  To better characterise the results of our protein-protein prediction tests, we sought to identify the

440  source of the testing edges RAPPPID correctly and incorrectly identified. Figure 4A and

441  Supplementary Figure S2 show that RAPPPID can accurately predict the testing set edges that

442  have a high confidence score in biologically supported channels of co-expression, experiments,

443  and database. However, the accuracy for the edges that have a high confidence score in the text-

444  mining channel is inferior to the other channels. Since edges that are only supported by text-

445  mining (but not by the other channels) are arguably the ones most prone to error, we expected

446  RAPPPID to have an inferior performance on such edges (since the edges themselves may not be

447  reliable). To test whether the inferior performance of RAPPPID in the text-mining channel in

448  Figure 4A are indeed due to such edges, for a fixed threshold k ($50 \leq k \leq 95$), we divided the

449  testing edges with a text-mining confidence score at least equal to k into two groups: a group

450  with "experiments" confidence score at least equal to 80% and a group with "experiments"

451  confidence score smaller than 80% (Figure 4B). Evaluating the testing edges in C1, C2, and C3

452  showed that for this channel, the accuracy on the former group is higher than the latter group

453  (sometimes as large as ~22% higher, Supplementary Figure S3). Repeating the analysis for co-

454  expression and database channels also confirmed this trend (Figure 4B, Supplementary Figure

455  S3). Taken together, these results suggest that the inferior performance of RAPPPID on the text-

456  mining channel in Figure 4A is indeed due to the edges that are supported only by text-mining

457  and not by other biologically identified channels.

22

458



459

**Figure 4: Accuracy of positive edges across edge confidence stratified by STRING channels.** (A) The percentage of correctly labelled positive edges are plotted for each major STRING channel. The x-axis denotes the channel edge confidence cut-off score for each curve's respective channel. (B) Here, we see a similar chart but rather than using each channel's respective score as a confidence cut-off, edges are excluded according to their text-mining confidence (x-axis). The solid curves include edges which have a channel confidence $\geq 80\%$ for the channel indicated by the curve's colour. Dashed curves conversely include edges whose channel confidence is ¿80% for the channel indicated by the curve's colour. In both (A) and (B) data shown reflects the C2 model/dataset.
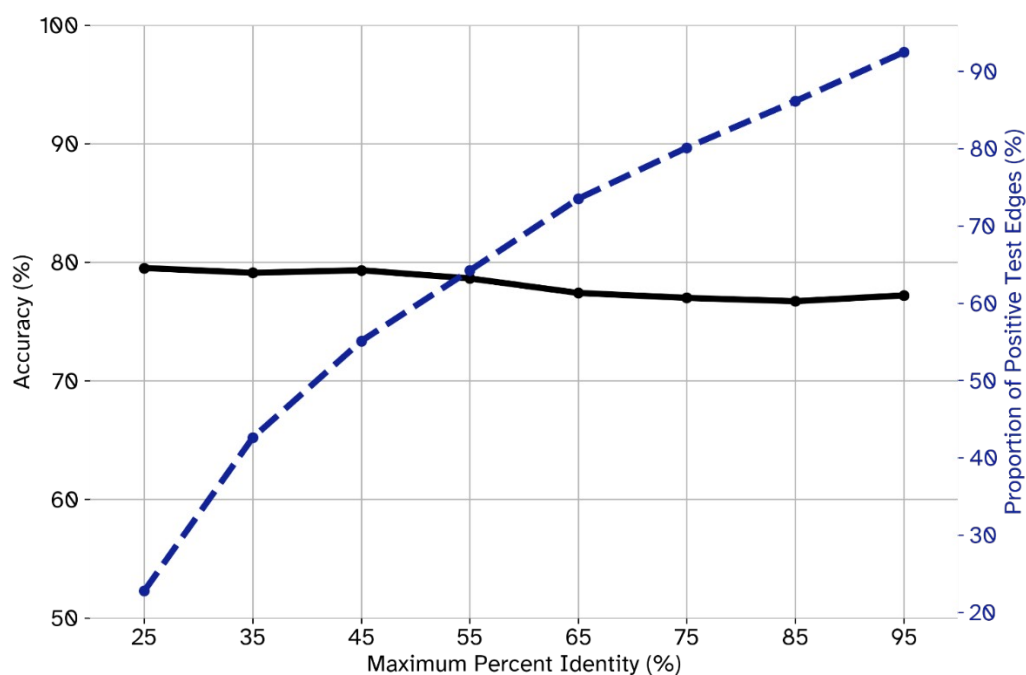
469

470

**Role of Protein Similarity on RAPPPID's Performance**

The procedures for C2 and C3 datasets were devised to reduce the information leakage by

avoiding testing on edges which contain proteins which are known to an algorithm during its

training. This safeguard against information leakage, however, does not account for proteins

which are known by different identifiers but share near identical protein sequences. Further

complicating the matter, sequence similarity is a valid PPI prediction feature and is often used by

23

477    methods as a proxy measure of co-evolution and conserved functional domains (Cong *et al.*,

478    2019; Jansen, 2003). Indeed, this strategy is leveraged by the SPRINT algorithm against which

479    RAPPPID was compared.



480
481    **Figure 5: Accuracy of positive edges as a function of similarity between testing and training**
482    **proteins in C2.** The similarity between testing and training proteins was measured using their
483    percent identity as computed by NCBI's PSI-BLAST software. The highest percent identity
484    between any training protein and a testing protein in a testing edge was considered to be that
485    testing edge's "maximum percent identity". The percentage of accurately labelled positive edges
486    (black curve, left y-axis) is reported for edges with maximum percent identities lower than the
487    threshold reported on the x-axis. The proportion of testing edges for each threshold values is
488    reported by the dashed blue curve and the right y-axis.
489

490    In spite of the challenges above, we sought to determine whether the superior performance of

491    RAPPPID (particularly in the strict datasets of C2 and C3) is due to sequence similarity between

492    testing and training proteins or not. For this purpose, we used PSI-BLAST algorithm (Altschul *et*

493    *al.*, 1997) to evaluate sequence similarities between each pair of testing/training proteins. Figure

24

494  5 shows the accuracy of RAPPPID on the C2 dataset when different degrees of restriction on

495  sequence similarity are imposed. More specifically, a threshold $t$ (x-axis in Figure 5) determined

496  the maximum allowable Percent Identity score between a testing protein and any of the training

497  proteins that were candidates to be similar to it (see Methods for details). Any testing protein that

498  did not satisfy this condition for the threshold $t$ was excluded from the calculation of accuracy.

499  As one moves towards larger values of $t$, the sequence similarity constraint loosens and $t=100\%$

500  is equivalent to the complete C2 dataset. Our analysis on C2 (Figure 5) and C3 (Supplementary

501  Figure S4) revealed that RAPPPID's accuracy is largely independent of the sequence similarities

502  between testing and training proteins and the performance of RAPPPID does not deteriorate

503  when removing testing proteins that have a highly similar training protein.

504

505

506  **Effect of different components on RAPPPID's performance**

507  To understand RAPPPID's performance when trained and tested on negative examples which are

508  experimentally validated and manually curated, we created a new C3 dataset where the negative

509  examples are provided by the Negatome database of non-interacting pairs. To combat the large

510  class imbalance caused by the relatively few pairs of human proteins in Negatome (see Methods

511  for details), additional random pair negatives were supplemented until both positive and negative

512  classes had an equal number of samples. After 13 epochs, RAPPPID achieved a testing AUROC

513  of 0.802 and a testing AUPR of 0.799 on C3 dataset, which is comparable to the performance of

514  RAPPPID using randomly selected negative edges.

515

25

516    Since RAPPPID utilizes random initialisation, mini-batch sampling, and token sampling, there is

517    some stochasticity present in its performance. First, we sought to test the effect of these

518    stochastic components and to ensure that the specific choice of training and test sets were not

519    responsible for the superior performance of RAPPPID. For this purpose, we ran RAPPPID twice

520    on three additional C3 datasets whose training, validation, and testing proteins were chosen at

521    random (a total of six times). In each of these runs, different seeds were used to assess the effect

522    of stochastic components of RAPPPID (Table 2). Overall, the average testing AUROC across

523    models trained on these additional three datasets was 0.792 (±0.007). Results from these

524    repeatability experiments illustrate that RAPPPID's strong performance on C3 datasets is not

525    tied to a specific set of testing or validation proteins, nor certain weight initialisation states.

526

527    Next, we conducted an ablation study (Table 2) using the same three C3 datasets and two runs of

528    each variant per dataset (a total of six models per variant). Each variant substitutes one

529    component of RAPPPID's regularisation, architecture or optimisation, allowing us to quantify

530    the difference in performance for which these components are responsible. In this table,

531    "RAPPPID-SWA" variant removes SWA, "RAPPPID+Adam" replaces Ranger21 with Adam

532    optimizer (learning rate selected by hyperparameter tuning), "RAPPPID-AWD" removes the

533    AWD regularisation, and "RAPPPID-SentencePiece" substitutes the SentencePiece tokens with

534    amino-acid-residue-level tokens. As can be seen in this table, the choice of optimizer has

535    negligible effect on the performance. However, SWA, AWD, and SentencePiece tokens all

536    contribute to the superior performance of RAPPPID, with the tokeniser having the largest

537    contribution.  One advantage that SentencePiece affords RAPPPID is the regularising effect of

538    random sampling tokens. Additionally, the SentencePiece tokeniser reduces the total length of

26

539    sequences inputted into the tokeniser. The protein sequences on which RAPPPID is trained,

540    while limited to 1,500 residues, are sufficiently large enough for gradients to vanish quiet

541    substantially; a phenomenon long known to plague RNNs when trained on long sequences

542    (Hochreiter, 1998).

543

544    **Table 3: Results from an ablation study conducted on RAPPPID.** Each model is
545    trained/tested twice on three randomly generated C3 datasets. The performance metrics
546    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% |

547

548

549    In acknowledgment of the recent prevalence and performance of attention-based methods in deep

550    learning (Vaswani *et al.*, 2017), the sixth and seventh variants replace the AWD-LSTM encoder

551    with transformers of varying sizes. "RAPPPID+TransfLG" has six layers and eight heads and

552    feed-forward networks with 2048 dimensions, while "RAPPPID+TransfSM" has fewer

553    parameters (two layers, two heads, and 80 dimensions). Both transformer variants have a dropout

554    rate of 20%. The number of heads and layers in RAPPPID+TransfLG were chosen to closely

555    match those used in the BERT language model (Devlin *et al.*, 2019) whose architecture has been

556    previously used for predicting protein function (Elnaggar *et al.*, 2021). The number of heads and

557    layers for RAPPPID+TransfSM were chosen to be representative of reasonable lower bounds for

27

558 those hyper-parameters. The results of these analyses (Table 2) reveal that both of these models

559 have a worse performance compared to the original RAPPPID model. The superior performance

560 of the small transformer compared to the large one suggests that the large number of parameters

561 required by transformers are leading to severe overfitting in this task and are the reason behind

562 the performance deterioration

563

564

565 **Transfer Learning on Protein-Ligand Data from X-Ray Crystallography Experiments**

566 RAPPPID's strong performance on C3 datasets demonstrate that it is capable of generalising to

567 predictions between proteins absent from the training set. Additionally, the ability to generalise

568 these models to make predictions on test datasets that are not representative of the training data

569 (e.g., different types of experiments or protocols are used to generate them) is also often

570 desirable, yet is often much more challenging. Transfer learning is one approach that can be used

571 to overcome the challenges imposed by differences between training and test datasets. Transfer

572 learning is the process by which a network is first trained on a large, high-quality dataset. Many

573 of these weights are then "frozen" such that they are no longer backpropagated through, while

574 the remaining "unfrozen" weights are trained on a dataset that is often smaller and/or belong to a

575 different input distribution (Yosinski *et al.*, 2014). Here, we use this approach to tackle an

576 interaction prediction task on a dataset that differs fundamentally from STRING: protein-ligand

577 interactions and sequences as determined from X-ray crystallography data.

578

579 STRING and similar PPI datasets curate interactions from a wide range of modalities and

580 provide sequences from reference proteomes. We've constructed a protein-ligand interaction task

28

581    by leveraging BioLiP, a semi-manually curated list of X-ray crystallography experiments

582    recorded in the Protein Data Bank (PDB) which reflect interactions between proteins and ligands

583    (Yang *et al.*, 2012; Berman, 2000). Using these datasets, we were able to extract sequences from

584    the PDB records and, after filtering out interactions with low-quality sequences, construct a

585    novel protein-ligand dataset. This dataset is available for download at

586    https://github.com/jszym/rapppid. The sequences extracted from PDB are measured by

587    fundamentally different modalities from those present in STRING/UniprotKB, and are often

588    truncated and incomplete as a result. Furthermore, the nature of the interactions recorded by

589    BioLiP are fundamentally different from STRING. Whereas STRING catalogues interactions

590    between large classes of proteins, the nature of X-ray crystallography biases captured

591    interactions to those with slower molecular dynamics and which do not primarily exist in

592    aliphatic environments (Carpenter *et al.*, 2008). These fundamental differences between BioLiP

593    and STRING make them ideal datasets to illustrate the ability to use transfer learning with

594    RAPPPID.

595

596    We first pre-trained RAPPPID on data from STRING, and then fine-tuned it on the BioLiP

597    dataset after the LSTM encoder weights were frozen, leaving only the fully-connected classifier

598    to be trained. To ensure no data leakage between the STRING and BioLiP datasets exist, protein

599    sequences in STRING determined to have more than 90% identity with those in BioLiP were

600    moved from the pre-training dataset to the fine-tuning dataset. C3-type training, validation, and

601    testing splits were constructed for both the pretraining and fine-tuning datasets for evaluation

602    purposes. Using this approach, RAPPPID pre-trained on the STRING dataset and fine-tuned on a

603    portion of the BioLip dataset achieved an AUROC of 0.909. To test our hypothesis that transfer

29

604    learning is indeed necessary to achieve good performance on this dataset, we also directly

605    applied the pre-trained RAPPPID model (without any fine-tuning on BioLip) to the BioLip test

606    split. As expected, this resulted in almost random results (AUROC = 0.548). These results

607    highlight the difficulties of generalising a model to fundamentally different datasets and

608    emphasise the utility of transfer learning to achieve good performance using RAPPPID in such

609    scenarios.

610

611

612    **RAPPPID predicts interaction of HER2 with Trastuzumab and Pertuzumab**

613    Peptides and proteins have emerged as an important class of therapeutics, enabling researchers to

614    target previously "undruggable" targets (Tsomaia, 2015). Hundreds of peptide and protein

615    therapeutics have been approved by the U.S. Food and Drug Administration (Usmani *et al.*,

616    2017) with applications in treating illnesses ranging from cancer to heart disease (Sikder *et al.*,

617    2019; Chen *et al.*, 2012). Here, we sought to illustrate how one might use RAPPPID to validate

618    hypothesised interactions between target proteins and candidate therapeutic proteins and peptides

619    through two examples: Trastuzumab and Pertuzumab.

620

621    Trastuzumab and Pertuzumab are two recombinant humanised monoclonal antibodies which are

622    used in combination in the treatment of metastatic breast cancers which belong to the HER2-

623    positive subtype (Boekhout *et al.*, 2011; Malenfant *et al.*, 2014). Both Trastuzumab and

624    Pertuzumab target distinct domains of the human epidermal growth factor receptor 2 (HER2).

625    We applied RAPPPID, trained on our *H. sapiens* STRING C2 dataset (which is more appropriate

626    for this application), to the sequences of the Trastuzumab and Pertuzumab antibody chains.

30

627    RAPPPID predicted that HER2 interacts with Trastuzumab and Pertuzumab with 86.20% and

628    95.11% probability, respectively. This suggests that RAPPPID may be utilized as part of an

629    interaction-based low-cost filtering or early validation step in the development of therapeutic

630    proteins and peptides.

631

632

633    **Discussion and Conclusion**

634    This study introduced RAPPPID, a deep learning method that addresses the challenges of

635    creating generalisable PPI prediction models posed by inherent characteristics of PPI datasets.

636    By adopting a modified AWD-LSTM training routine, RAPPPID was able to surpass state-of-

637    the-art models under testing conditions that carefully controlled for information leakage and

638    other sources of prediction accuracy inflation. Further experiments were conducted to confirm

639    the results were independent of the specific proteins present in the training and testing splits.

640    RAPPPIDs ability to PPIs in the STRING database was shown to increase with strong biological

641    evidence for the interaction. This relationship between PPI evidence and RAPPPID predictive

642    ability illustrates that RAPPPID accurately reflects our confidence in interactions, and testing

643    performance is not disproportionately inflated by spurious, low-confidence interactions.

644    Moreover, assessment of the sequence similarity between testing and training proteins revealed

645    that the superior performance of RAPPPID is not due to the presence of highly similar protein

646    pairs in testing and training, and the accuracy of RAPPPID was largely stable with a small

647    improvement when highly similar testing proteins were excluded.

648

31

649 Developing appropriate and meaningful benchmark datasets for PPI prediction remains a

650 challenging problem for a number of reasons. Firstly, deep learning tasks rely on large, high-

651 quality datasets to obtain meaningful generalised models. Such datasets are few and far between.

652 Projects like HIPPIE and iRefWeb join STRING in being among the best examples of PPI

653 datasets which integrate multiple sources to assure both quality and quantity of PPI edges

654 (Alanis-Lobato *et al.*, 2017; Turner *et al.*, 2010). Despite this, STRING is the only dataset among

655 these three which has an appropriate number of high-confidence edges for the purpose of

656 learning deep learning model. Specifically, there are 98.5% fewer edges in HIPPIE than in

657 STRING between human proteins at a 95% confidence threshold. Even when confidence

658 thresholds are lowered to 85%, HIPPIE holds 87.9% fewer human edges than STRING.

659 Similarly STRING at a 95% confidence threshold has 75% more human edges when considering

660 iRefWeb edges with 3 supporting references or more.

661

662 Secondly, this reliance on large, representative datasets is further exacerbated by the unique

663 overfitting challenges posed by the characteristics of PPI data, as explored in this work and in

664 (Park and Marcotte, 2012). While RAPPPID mitigates the overfitting tendencies of PPI data

665 through regularisation, it also relies on large, representative training data to do so. In the

666 construction of our C3 datasets, it is necessary to discard edges which are not between proteins

667 of the same validation split. As a result, we observe a further decrease of up to 33.9% in the

668 number of the already-precious-few edges afforded to us by STRING.

669

670 As we've shown in this work, the construction of benchmark datasets is critical for effectively

671 comparing and evaluating PPI prediction tasks. It is thanks to large, quality PPI dataset projects

32

672   like STRING that we might construct meaningful and appropriate datasets against which to

673   evaluate PPI prediction methods like RAPPPID. Additional efforts and methodologies to collect

674   and integrate PPI edges of the scale of STRING at high confidence levels are greatly desired, as

675   they can help identify and mitigate biases that inevitably arise when constructing datasets.

676

677   The task of PPI prediction is related to the problem of protein docking inference, whereby

678   computational models predict the atomic interactions between two proteins. While many

679   methods have historically been based on the fast Fourier transform for energy evaluation (Desta

680   *et al.*, 2020), new and effective deep learning methods have also become available. One such

681   method is AlphaFold-multimer (Evans *et al.*, 2021), which builds upon a Transformer model for

682   the prediction of protein structure (Jumper *et al.*, 2021) to infer the interface of homo and hetero

683   protein dimers at an atomic level. Integrating docking predictions can be an interesting future

684   direction to help improve RAPPPID's PPI prediction generalisation.

685

686   RAPPPID's ability to predict interactions warrants further study into relevant tasks that might

687   benefit from a similar approach. The RAPPPID architecture might be modified for the tasks of

688   binding site and protein function prediction. These tasks are related to PPI prediction and as a

689   result are exposed to similar challenges to which RAPPPID is well suited. However, in all these

690   cases, it is crucial to consider strict rules for cross-validation and data splitting to ensure data

691   leakage is avoided.

692

693

694

33

700

701

702 **Authors' contributions:**

703 AE and JS conceived the study, designed the project and the algorithm, and wrote the

704 manuscript. JS implemented the pipeline and performed the statistical analyses of the results. All

705 authors read and approved the final manuscript.

706

707

708 **References**
709
710 Alanis-Lobato,G. *et al.* (2017) HIPPIE v2.0: enhancing meaningfulness and reliability of
711     protein–protein interaction networks. *Nucleic Acids Res*, **45**, D408–D414.
712 Altschul,S.F. *et al.* (1990) Basic local alignment search tool. *Journal of Molecular Biology*, **215**,
713     403–410.
714 Altschul,S.F. *et al.* (1997) Gapped BLAST and PSI-BLAST: a new generation of protein
715     database search programs. *Nucleic Acids Research*, **25**, 3389–3402.
716 Anfinsen,C.B. (1973) Principles that Govern the Folding of Protein Chains. *Science*, **181**, 223–
717     230.
718 Athiwaratkun,B. *et al.* (2019) There Are Many Consistent Explanations of Unlabeled Data: Why
719     You Should Average. In, *ICLR*.
720 Ben-Hur,A. and Noble,W.S. (2006) Choosing negative examples for the prediction of protein-
721     protein interactions. *BMC Bioinformatics*, **7**, S2–S2.
722 Ben-Hur,A. and Noble,W.S. (2005) Kernel methods for predicting protein–protein interactions.
723     *Bioinformatics*, **21**, i38–i46.
724 Berman,H.M. (2000) The Protein Data Bank. *Nucleic Acids Research*, **28**, 235–242.

725  Blohm,P. *et al.* (2014) Negatome 2.0: a database of non-interacting proteins derived by literature
726      mining, manual annotation and protein structure analysis. *Nucl. Acids Res.*, **42**, D396–
727      D400.
728  Boekhout,A.H. *et al.* (2011) Trastuzumab. *Oncologist*, **16**, 800–810.
729  Brock,A. *et al.* (2021) High-Performance Large-Scale Image Recognition Without
730      Normalization. *arXiv:2102.06171 [cs, stat]*.
731  Bromley,J. *et al.* (1993) Signature verification using a "siamese" time delay neural network. *Int.*
732      *J. Patt. Recogn. Artif. Intell.*, **07**, 669–688.
733  Browne,F. *et al.* (2007) Supervised Statistical and Machine Learning Approaches to Inferring
734      Pairwise and Module-Based Protein Interaction Networks. In, *2007 IEEE 7th*
735      *International Symposium on BioInformatics and BioEngineering*. IEEE, Boston, MA,
736      USA, pp. 1365–1369.
737  Carpenter,E.P. *et al.* (2008) Overcoming the challenges of membrane protein crystallography.
738      *Current Opinion in Structural Biology*, **18**, 581–586.
739  Chen,H.H. *et al.* (2012) Novel Protein Therapeutics for Systolic Heart Failure: Chronic
740      Subcutaneous B-Type Natriuretic Peptide. *Journal of the American College of*
741      *Cardiology*, **60**, 2305–2312.
742  Chen,M. *et al.* (2019) Multifaceted protein–protein interaction prediction based on Siamese
743      residual RCNN. *Bioinformatics*, **35**, i305–i314.
744  Cong,Q. *et al.* (2019) Protein interaction networks revealed by proteome coevolution. *Science*,
745      **365**, 185–189.
746  Desta,I.T. *et al.* (2020) Performance and Its Limits in Rigid Body Protein-Protein Docking.
747      *Structure*, **28**, 1071-1081.e3.
748  Devlin,J. *et al.* (2019) BERT: Pre-training of Deep Bidirectional Transformers for Language
749      Understanding. *arXiv:1810.04805 [cs]*.
750  Dick,K. *et al.* (2020) PIPE4: Fast PPI Predictor for Comprehensive Inter- and Cross-Species
751      Interactomes. *Scientific Reports*, **10**, 1390.
752  Ding,Y. *et al.* (2016) Predicting protein-protein interactions via multivariate mutual information
753      of protein sequences. *BMC Bioinformatics*, **17**, 398.
754  Elnaggar,A. *et al.* (2021) ProtTrans: Towards Cracking the Language of Life's Code Through
755      Self-Supervised Learning.
756  Evans,R. *et al.* (2021) Protein complex prediction with AlphaFold-Multimer Bioinformatics.
757  Falcon,W. *et al.* (2020) PyTorchLightning/pytorch-lightning: 0.7.6 release Zenodo.
758  Gillis,J. *et al.* (2014) Bias tradeoffs in the creation and analysis of protein–protein interaction
759      networks. *Journal of Proteomics*, **100**, 44–54.
760  Henikoff,S. and Henikoff,J.G. (1992) Amino acid substitution matrices from protein blocks.
761      *Proceedings of the National Academy of Sciences*, **89**, 10915–10919.
762  Hochreiter,S. (1998) The Vanishing Gradient Problem During Learning Recurrent Neural Nets
763      and Problem Solutions. *Int. J. Unc. Fuzz. Knowl. Based Syst.*, **06**, 107–116.
764  Hochreiter,S. and Schmidhuber,J. (1997) Long Short-Term Memory. *Neural Computation*, **9**,
765      1735–1780.
766  Huttlin,E.L. *et al.* (2017) Architecture of the human interactome defines protein communities and
767      disease networks. *Nature*, **545**, 505–509.

768  Izmailov,P. *et al.* (2019) Averaging Weights Leads to Wider Optima and Better Generalization.
769       *arXiv:1803.05407 [cs, stat].*
770  Jansen,R. (2003) A Bayesian Networks Approach for Predicting Protein-Protein Interactions
771       from Genomic Data. *Science*, **302**, 449–453.
772  Jassal,B. *et al.* (2019) The reactome pathway knowledgebase. *Nucleic Acids Research*, gkz1031.
773  Jumper,J. *et al.* (2021) Highly accurate protein structure prediction with AlphaFold. *Nature*, **596**,
774       583–589.
775  Kanehisa,M. (2000) KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids*
776       *Research*, **28**, 27–30.
777  Kudo,T. (2018) Subword Regularization: Improving Neural Network Translation Models with
778       Multiple Subword Candidates. *arXiv:1804.10959 [cs].*
779  Kudo,T. and Richardson,J. (2018) SentencePiece: A simple and language independent subword
780       tokenizer and detokenizer for Neural Text Processing. In, *Proceedings of the 2018*
781       *Conference on Empirical Methods in Natural Language Processing: System*
782       *Demonstrations*. Association for Computational Linguistics, Brussels, Belgium, pp. 66–
783       71.
784  Li,Y. and Ilie,L. (2017) SPRINT: ultrafast protein-protein interaction prediction of the entire
785       human interactome. *BMC Bioinformatics*, **18**, 485.
786  Lipton,Z.C. *et al.* (2015) A Critical Review of Recurrent Neural Networks for Sequence
787       Learning. *arXiv:1506.00019 [cs].*
788  Loshchilov,I. and Hutter,F. (2019) Decoupled Weight Decay Regularization. In, *ICLR*.
789  Malenfant,S.J. *et al.* (2014) Pertuzumab: a new targeted therapy for HER2-positive metastatic
790       breast cancer. *Pharmacotherapy*, **34**, 60–71.
791  Merity,S. *et al.* (2017) Regularizing and Optimizing LSTM Language Models.
792       *arXiv:1708.02182 [cs].*
793  Misra,D. (2020) Mish: A Self Regularized Non-Monotonic Activation Function.
794       *arXiv:1908.08681 [cs, stat].*
795  Orchard,S. *et al.* (2012) Protein interaction data curation: the International Molecular Exchange
796       (IMEx) consortium. *Nature Methods*, **9**, 345–350.
797  Orchard,S. *et al.* (2014) The MIntAct project—IntAct as a common curation platform for 11
798       molecular interaction databases. *Nucl. Acids Res.*, **42**, D358–D363.
799  Oughtred,R. *et al.* (2020) The BioGRID database: A comprehensive biomedical resource of
800       curated protein, genetic, and chemical interactions. *Protein Science: A Publication of the*
801       *Protein Society*, **30**, 187–200.
802  Park,Y. and Marcotte,E.M. (2012) Flaws in evaluation schemes for pair-input computational
803       predictions. *Nat Methods*, **9**, 1134–1136.
804  Paszke,A. *et al.* (2019) PyTorch: An Imperative Style, High-Performance Deep Learning
805       Library. In, Wallach,H. *et al.* (eds), *Advances in Neural Information Processing Systems*
806       *32*. Curran Associates, Inc., pp. 8024–8035.
807  Richoux,F. *et al.* (2019) Comparing two deep learning sequence-based models for protein-
808       protein interaction prediction. *arXiv:1901.06268 [cs, q-bio, stat].*
809  Roux,K.J. *et al.* (2012) A promiscuous biotin ligase fusion protein identifies proximal and
810       interacting proteins in mammalian cells. *Journal of Cell Biology*, **196**, 801–810.

811  Salwinski,L. *et al.* (2004) The Database of Interacting Proteins: 2004 update. *Nucleic Acids Res.*,
812      **32**, D449-451.
813  Schuster,M. and Nakajima,K. (2012) Japanese and Korean voice search. In, *2012 IEEE*
814      *International Conference on Acoustics, Speech and Signal Processing (ICASSP).*, pp.
815      5149–5152.
816  Sikder,S. *et al.* (2019) Long-term delivery of protein and peptide therapeutics for cancer
817      therapies. *Expert opinion on drug delivery*.
818  Snider,J. *et al.* (2015) Fundamentals of protein interaction network mapping. *Mol Syst Biol*, **11**,
819      848.
820  Srivastava,N. *et al.* (2014) Dropout: a simple way to prevent neural networks from overfitting. *J.*
821      *Mach. Learn. Res.*, **15**, 1929–1958.
822  Szklarczyk,D. *et al.* (2019) STRING v11: protein–protein association networks with increased
823      coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic*
824      *Acids Res*, **47**, D607–D613.
825  Tabe-Bordbar,S. *et al.* (2018) A closer look at cross-validation for assessing the accuracy of gene
826      regulatory networks and models. *Sci Rep*, **8**, 6620.
827  Tsomaia,N. (2015) Peptide therapeutics: targeting the undruggable space. *European journal of*
828      *medicinal chemistry*.
829  Turner,B. *et al.* (2010) iRefWeb: interactive analysis of consolidated protein interaction data and
830      their supporting evidence. *Database*, **2010**, baq023–baq023.
831  Usmani,S. *et al.* (2017) THPdb: Database of FDA-approved peptide and protein therapeutics.
832      *PloS one*.
833  Vaswani,A. *et al.* (2017) Attention is all you need. In, *Proceedings of the 31st International*
834      *Conference on Neural Information Processing Systems*, NIPS'17. Curran Associates Inc.,
835      Red Hook, NY, USA, pp. 6000–6010.
836  Vidal,M. and Fields,S. (2014) The yeast two-hybrid assay: still finding connections after 25
837      years. *Nat Methods*, **11**, 1203–1206.
838  Wan,L. *et al.* (2013) Regularization of Neural Networks using DropConnect. In, Dasgupta,S. and
839      McAllester,D. (eds), *Proceedings of the 30th International Conference on Machine*
840      *Learning*, Proceedings of Machine Learning Research. PMLR, Atlanta, Georgia, USA,
841      pp. 1058–1066.
842  Wright,L. and Demeure,N. (2021) Ranger21: a synergistic deep learning optimizer.
843      *arXiv:2106.13731 [cs]*.
844  Yang,J. *et al.* (2012) BioLiP: a semi-manually curated database for biologically relevant ligand–
845      protein interactions. *Nucleic Acids Research*, **41**, D1096–D1103.
846  Yong,H. *et al.* (2020) Gradient Centralization: A New Optimization Technique for Deep Neural
847      Networks. In, *ECCV*.
848  Yosinski,J. *et al.* (2014) How transferable are features in deep neural networks? In,
849      Ghahramani,Z. *et al.* (eds), *Advances in Neural Information Processing Systems 27*.
850      Curran Associates, Inc., pp. 3320–3328.
851  Zaremba,W. *et al.* (2015) Recurrent Neural Network Regularization. *arXiv:1409.2329 [cs]*.
852  Zhang,M.R. *et al.* (2019) Lookahead Optimizer: k steps forward, 1 step back. In, *NeurIPS*.