

1 CONSTRUCTION OF *IN SILICO* PROTEIN-PROTEIN
2 INTERACTION NETWORKS ACROSS DIFFERENT
3 TOPOLOGIES USING MACHINE LEARNING

4 Loïc Lannelongue^{1,2,3,*}, Michael Inouye^{1,2,3,4,5,6,*}

5 ¹Cambridge Baker Systems Genomics Initiative, Department of Public Health and Primary Care, University of Cambridge,
6 Cambridge, UK

7 ²British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of
8 Cambridge, Cambridge, UK

9 ³Health Data Research UK Cambridge, Wellcome Genome Campus and University of Cambridge, Cambridge, UK

10 ⁴Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, Melbourne, Victoria, Australia

11 ⁵British Heart Foundation Centre of Research Excellence, University of Cambridge, Cambridge, UK

12 ⁶The Alan Turing Institute, London, UK

13

14 * Correspondence: LL (LL582@medschl.cam.ac.uk) and MI (mi336@medschl.cam.ac.uk; minouye@baker.edu.au)

15 **ABSTRACT**

16 Protein-protein interactions (PPIs) are essential to understanding biological pathways as well as
17 their roles in development and disease. Computational tools have been successful at predicting
18 PPIs *in silico*, but the lack of consistent and reliable frameworks for this task has led to network
19 models that are difficult to compare and, overall, a low level of trust in the PPI predictions. To
20 better understand the underlying mechanisms that underpin these models, we designed B4PPI,
21 an open-source framework for benchmarking that accounts for a range of biological and statistical
22 pitfalls while facilitating reproducibility. We use B4PPI to shed light on the impact of network
23 topology and how different algorithms deal with highly connected proteins. By studying functional
24 genomics-based and sequence-based models (the two most popular approaches) on human PPIs,
25 we show their complementarity as the former performs best on lone proteins while the latter
26 specialises in interactions involving hubs. We also show that algorithm design has little impact on
27 performance with functional genomic data. We replicate our results between both human and *S.*
28 *cerevisiae* data and demonstrate that models using functional genomics are better suited to PPI
29 prediction across species. With rapidly increasing amounts of sequence and functional genomics
30 data, our study provides a systematic foundation for future construction, comparison and
31 application of PPI networks.

32

33 INTRODUCTION

34 Protein-protein interactions (PPIs) are central to protein function and inform a wide range of
35 biomedical applications, from mechanistic studies [1], [2] to drug development [3], [4]. Better
36 understanding these interactions is critical for successfully mapping biological pathways, but the
37 diversity of PPIs and the scale of the network make this a difficult task. Experimental methods to
38 map PPIs exist, but even when high-throughput tend to focus on proteins of interest.

39 Computational methods can address the issue of scalability and experimental bias. Given a pair of
40 proteins and some characteristics of each one, machine learning models can learn to predict the
41 likelihood of interaction. Numerous methods have been developed for this, using the full range of
42 machine learning models, from early work on *Saccharomyces cerevisiae* [5]–[8] to algorithms
43 dedicated to human PPIs [9]–[13]. Yet, despite a wealth of tools, the mechanics and consequences
44 of the underlying inference are still poorly understood, and it is unclear why models with similar
45 performance make vastly different predictions. Reported performance scores often cannot be
46 compared or replicated due to proprietary data and inconsistent or flawed assessment methods.
47 As a consequence, there are multiple issues for *in silico* PPIs: it is unclear what the state-of-the-art
48 is, analyses are difficult to reconcile, the development of new models is inefficient, follow-up
49 mechanisms studies are likely undermined and, ultimately, there are different versions of the
50 underlying molecular networks that describe protein function.

51 A unified framework for PPI inference would improve the development and reliable assessment
52 of new models, and would facilitate the overdue widespread adoption of PPI predictions for
53 downstream analysis. Replicable, trustworthy and generalisable high-performing models can
54 capture more causal biology and enhance many aspects of biological research such as
55 experimental designs and drug development.

56 In this work, we design a robust and standardised approach to *in silico* PPI prediction that accounts
57 for both biological and statistical pitfalls and leverages the strength of large, open-source and
58 professionally curated databases. We make publicly available benchmarking standards for human
59 and yeast PPIs to accelerate future discoveries and lay the foundations for similar datasets for
60 other organisms. Within this framework, we study and compare the main approaches to PPI
61 prediction in humans, based on functional genomic (FG) information or amino acid sequences
62 alone. We highlight why both perspectives are still relevant today and how each adapts to the PPI
63 network's topology. In particular, we show that the presence of highly connected proteins in the
64 networks has a drastic impact on prediction models and is an area where FG and sequence models
65 diverge. We also replicate these results between human and yeast (*S. cerevisiae*) and show which
66 tools are most suitable to cross-species predictions. This work provides robust foundations for
67 future developments in PPI prediction models, but also gives critical insight into which models can
68 and should be used in different situations.

70 B4PPI: A robust and open-source benchmark for PPI prediction

71 The lack of a consistent way to assess PPI prediction algorithms has hindered the development of
72 such algorithms and reduced their impact by making it difficult to reuse models for downstream
73 analysis [14]. Benchmarks are important for replicability [15], and when combined with carefully
74 curated datasets, they enable fast development through trial and error. Our Benchmarking
75 Pipeline for the Prediction of Protein-Protein Interactions in Humans (B4PPI-Human) includes both
76 carefully selected training and testing sets and a collection of input features to enable such trials.
77 Standard UniProt IDs are used throughout to easily combine these with other data sources.
78 Relevant metrics are selected with guidelines on how to share them. All this, alongside the pre-
79 processing steps and relevant guidelines, is made available online and can be downloaded easily
80 from <https://github.com/Llanelongue/B4PPI>. An example of a reporting sheet is in **Figure 1**.

81 The complexity of the underlying biological mechanisms of PPIs introduces pitfalls that need to be
82 considered when evaluating models. First, the way non-interacting proteins are selected for
83 training is important. While some efforts have used proteins known to be localised in different
84 parts of the cell [5], [12], [16], [17], this has been shown to be unreliable and a source of significant
85 bias that overestimates accuracy [18]. An alternative is to use a database of experimentally tested
86 non-interacting proteins, but leading resources such as Negatome have only ~1,300 pairs and thus
87 offer limited coverage [16], [19]. Considering the scarcity of PPIs, randomly sampling pairs of
88 proteins has a very low risk of false negative and limits selection bias (i.e. focusing on known
89 proteins of interest) [18], [20]. However, the impact of the associated imbalance between
90 interacting and non-interacting proteins should be taken into account when training models on
91 balanced datasets [21]. Lastly, each observation is itself a pair of proteins. Even when ensuring
92 that the two sets don't have pairs in common, there can be individual proteins present in both the
93 training and testing sets. This protein-level overlap, often overseen, has been shown to
94 significantly affect the performance of an algorithm and should therefore be properly assessed
95 [22]. Despite being documented in the literature, these pitfalls are still unevenly accounted for in
96 published works. This, alongside inconsistencies in the choice of testing sets and performance
97 metrics explains why, despite the number of algorithms released in recent years, there is still no
98 simple way to compare a new approach to the state-of-the-art, or even know what the state-of-
99 the-art is.

100 The essential aspects of training and assessment that should be systematically accounted for are
101 (1) the quality of the positive examples (i.e. the interacting proteins), (2) how non-interacting
102 proteins are selected for the gold standard, (3) a suitable split between training and testing sets,
103 in particular regarding individual proteins, and (4) the metrics to evaluate and compare models.
104 B4PPI seeks to address these four aspects of benchmarking.

105 When building a gold standard for machine learning algorithms, quality and representativity are
106 the most important aspects to consider, which makes IntAct [23] a database of choice for

107 interacting proteins. It aggregates reliable evidence of molecular interactions from over 20,000
108 publications, which are manually curated, and includes data from other interactions databases
109 such as the IMEx consortium [24]. We further limited the risk of false positives by removing low-
110 quality interactions, for example the ones based on spatial colocalisation only (**Methods**). The final
111 dataset comprised 78,229 interactions, covering 12,026 proteins (out of the 20,386 registered in
112 UniProt).

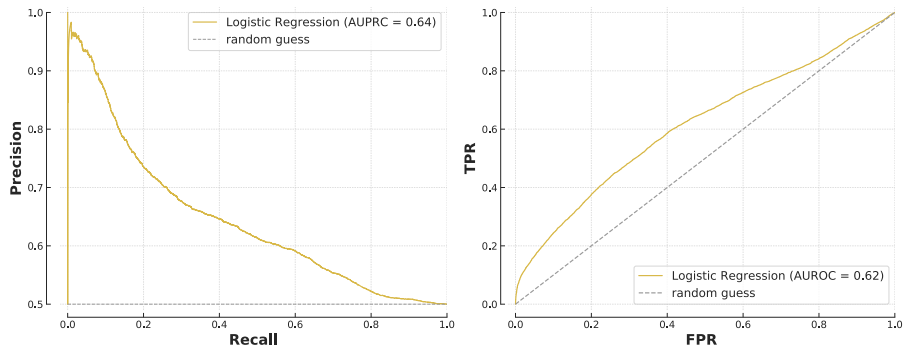
113 To select non-interacting proteins to serve as negative examples, randomly sampling protein pairs
114 is the approach with the lowest probability of error considering the scarcity of the PPI network
115 [25]. Non-interacting proteins can be sampled using a uniform distribution, i.e. all proteins have
116 an equal probability of being selected, which leads to an unbiased set, representative of the
117 general population of protein pairs. However, PPI networks are known to be similar to scale-free
118 networks, i.e. composed of a few highly connected nodes, called *hubs*, and numerous *lone*
119 proteins with few interactions [26] (**Supplementary Figure 1**). Consequently, hubs are over-
120 represented in a set of PPIs. For example in our curated set from IntAct, the top 20% of proteins
121 by number of interactions were involved in 94% of PPIs. But when uniformly sampling protein
122 pairs, the same top 20% were only involved in 37% of non-interacting proteins. Although expected,
123 this can be an issue for machine learning algorithms that would identify hubs and systematically
124 predict a positive interaction when hubs are involved. Such a strategy would maximise accuracy
125 on the training set but lead to a majority of false positives when making predictions on new pairs.
126 To mitigate this, a balanced sampling can be used [27], where the probability of sampling a protein
127 for the negative set is proportionate to its frequency in the positive set. It has been shown that
128 each strategy serves a different purpose [20]; balanced sampling is beneficial for training models
129 but shouldn't be used for evaluating them, as the induced bias makes metrics less meaningful.
130 This was the strategy implemented here, where non-interacting proteins were selected with
131 balanced sampling for the training set and uniform sampling for the testing sets.

132

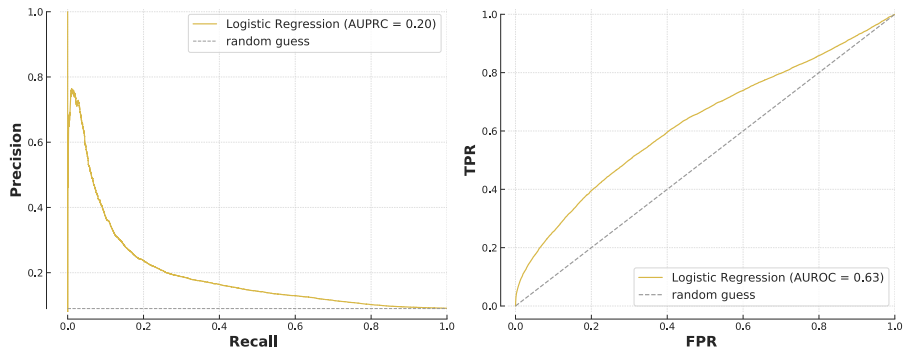
Reporting sheet B4PPI-Human: Logistic Regression

PR and ROC curves on T1* and T2**

Predictions on T1 (n=24,898, 50% +)



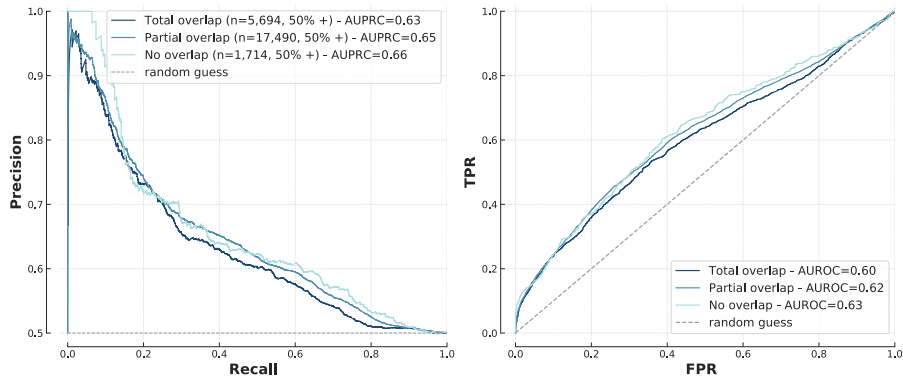
Predictions on T2 (n=136,939, 9% +)



* First testing set, used to compare models on an independent set and investigate protein-level overlap.

** Second testing set, used to assess generalisation on an imbalanced dataset (10 times more negative examples than positive ones).

Impact of protein-level overlap



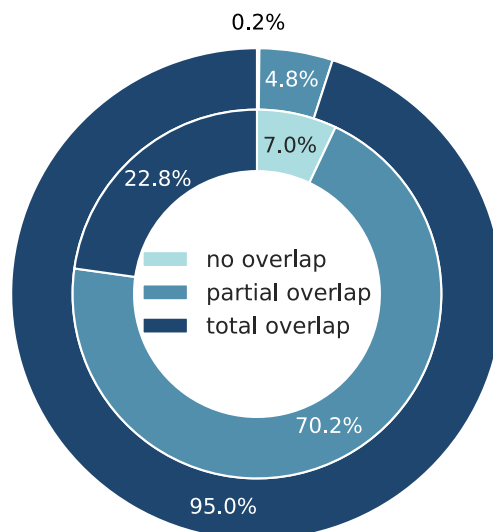
	Running time	Memory	Energy used	Carbon footprint (UK)
Training once	6s	Negligible	< 0.01 kWh	0.007 gCO ₂ e
Training incl. hyperparameters tuning	same*	Negligible	same	same
Inference	<1s	Negligible	~0	~0

* No hyperparameters to tune

133

134 Figure 1: Reported performance sheet of the logistic regression (FG-based) on B4PPI-Human.

135 In the presence of limited data, the division of the gold standard between training and testing sets
136 is critical to simultaneously optimise learning and obtain meaningful generalisation metrics. Here,
137 the testing set should achieve several objectives, (1) provide performance metrics on a new,
138 independent set, (2) measure the impact, or absence of impact, of protein-level overlap, (3)
139 demonstrate how the model can generalise to real-world data. Since a single set cannot achieve
140 simultaneously (2) and (3), as the careful selection of proteins to measure overlap biases the
141 dataset, we designed two testing sets $T1$ and $T2$ (**Methods**). $T1$ should be used to compare
142 different approaches with an independent set and investigate protein-level overlap, and $T2$ should
143 be used to assess generalisation. $T1$ was built by purposefully leaving some proteins out of the
144 training set; we demonstrated the importance of this as dividing the training and testing sets
145 conventionally (using, for example, the popular *scikit-learn* library) resulted in almost all pairs
146 (95%) in the testing set sharing at least one protein with the training set (**Figure 2**), which may lead
147 to overestimating performances [22]. $T2$, with ten times more negative examples than positive
148 ones (**Supplementary Table 1**), can then be used to assess how models perform in a more realistic
149 setting where positive interactions are rare compared to non-interacting proteins.



150
151 **Figure 2: The impact of train/test splitting strategies on protein-level overlap. The common splitting strategy is to**
152 **allocate pairs randomly (outer ring) while here we set aside proteins for testing (inner ring).**

153 The choice of metrics is a crucial element of a benchmark, and summarising the results by a single
154 number, such as accuracy or AUROC, is often misleading [28]. We report both the Receiver
155 Operating Characteristic (ROC) and the Precision-Recall (PR) curves which highlight nuanced and
156 complementary aspects of PPI models. In addition, to address the environmental impact of
157 bioinformatic tools [29], we also reported the carbon footprint of training models, measured using
158 the Green Algorithms calculator [30].

159 The elements described above represent the minimum needed for reproducible benchmarks and
160 researchers who wish to use their own input features can evaluate their models on these
161 partitions. However, to rapidly test a new model, it is useful to have access to carefully selected
162 and highly accurate protein properties. The two main categories of features used are amino acid
163 sequences and functional genomics annotations, such as subcellular localisation and biological
164 functions. These are available with B4PPI, from the professionally curated databases UniProt, the
165 Human Protein Atlas (HPA) [31], [32] and Bgee [33] (**Methods** and **Table 1** for the full list of
166 features).

167 With B4PPI, we could compare different models in a consistent manner to better understand what
168 aspects of the underlying biology are captured by each method. We focused here on FG-based
169 and sequence-based models as they have been widely used and rarely compared, despite
170 attempts at combining them.

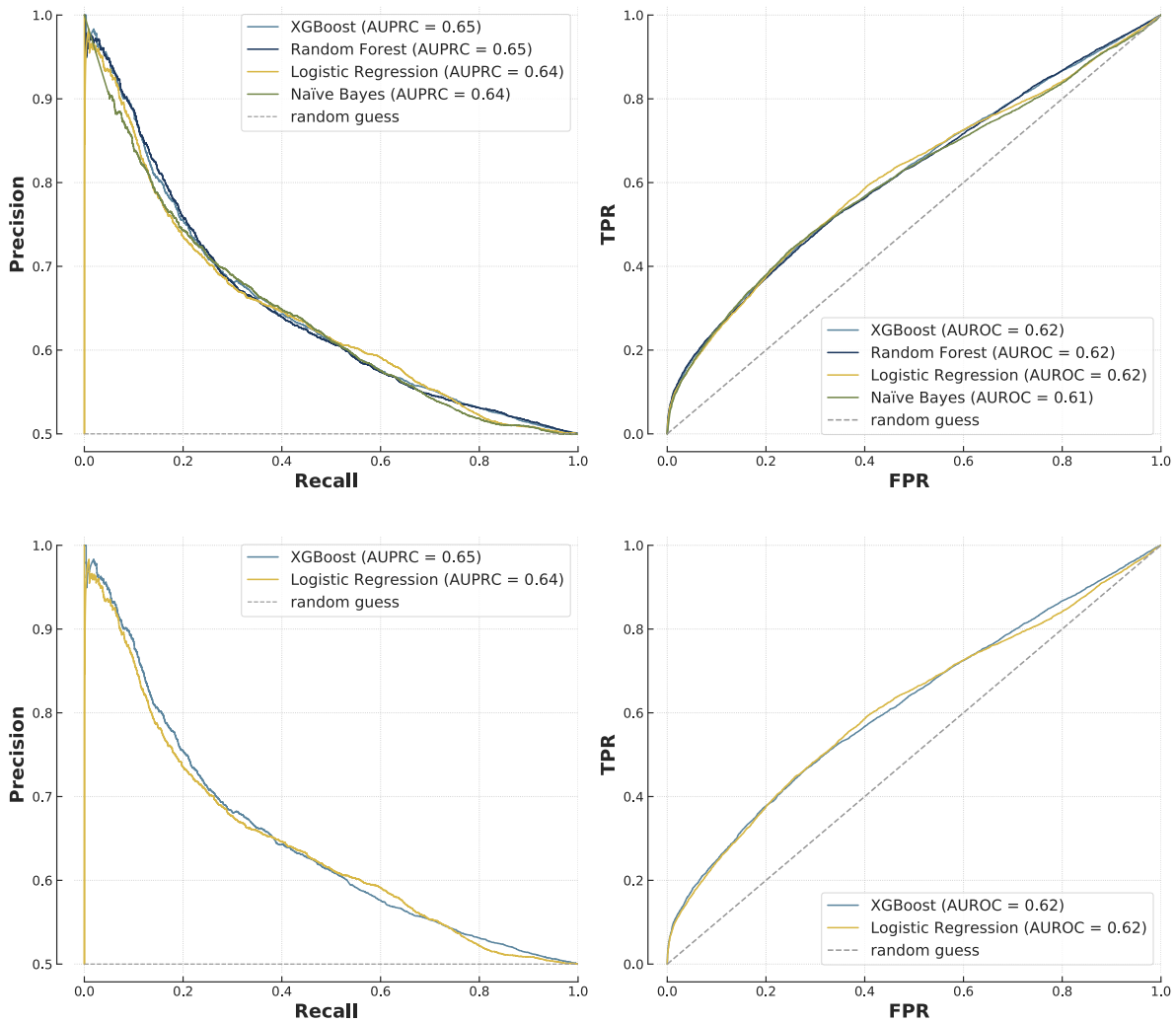
171 **FG-based linear models achieved top performance**

172 In FG-based models, FG annotations are pre-processed to compute similarity measures, such as
173 colocalisation, between proteins (**Methods**). The low dimensionality of the transformed problem
174 explains the success of standard machine learning algorithms; in particular, Naïve Bayes Classifiers
175 [34], decision trees [35] and Random Forests [36] have been the most popular choices [5]–[7], [37].
176 Despite the proven track record of such tools, the more recent XGBoost algorithm [38] has been
177 shown to outperform them in other situations like kidney disease diagnostic [39], which motivated
178 its inclusion in this analysis.

179 Using logistic regression as a baseline, we reported PR and ROC curves on the two test sets (**Figure**
180 **1**). A list of coordinates for these two curves was made available so that future models can be
181 compared without unnecessary re-training. We also reported the training time, 6 seconds, and
182 the carbon footprint, close to 0 gCO₂e. We then compared other models to this baseline and
183 produced similar performance sheets (**Supplementary Figure 2**).

184 We found that more complex algorithms brought little improvement over logistic regression, as
185 most models performed similarly on *T1* (**Figure 3** and **Supplementary Figure 3**). XGBoost and
186 Random Forest showed minor improvement in AUROC and AUPRC, but the difference between
187 the ROC curves of the logistic regression and XGBoost is non-significant ($p=0.27$) (**Methods**).
188 Moreover, XGBoost was more efficient than Random Forest as it had nearly half the runtime (30s
189 vs 54s). When studying the coefficients of the linear regression, we found that most decisions are
190 based on common biological processes, co-localisation (cellular compartment) and common
191 domains, all three coefficients being significant ($p<0.001$) (**Supplementary Figure 4**).

192 The reporting standard also enabled us to look at finer performance metrics, broken down by
193 protein-level overlap (i.e. individual proteins common to the training and testing sets). Comparing
194 PR and ROC curves showed that both logistic regression and XGBoost were unaffected by the level
195 of overlap (**Figure 1** and **Supplementary Figure 2**), and can therefore transfer effectively to new
196 proteins.



197

198

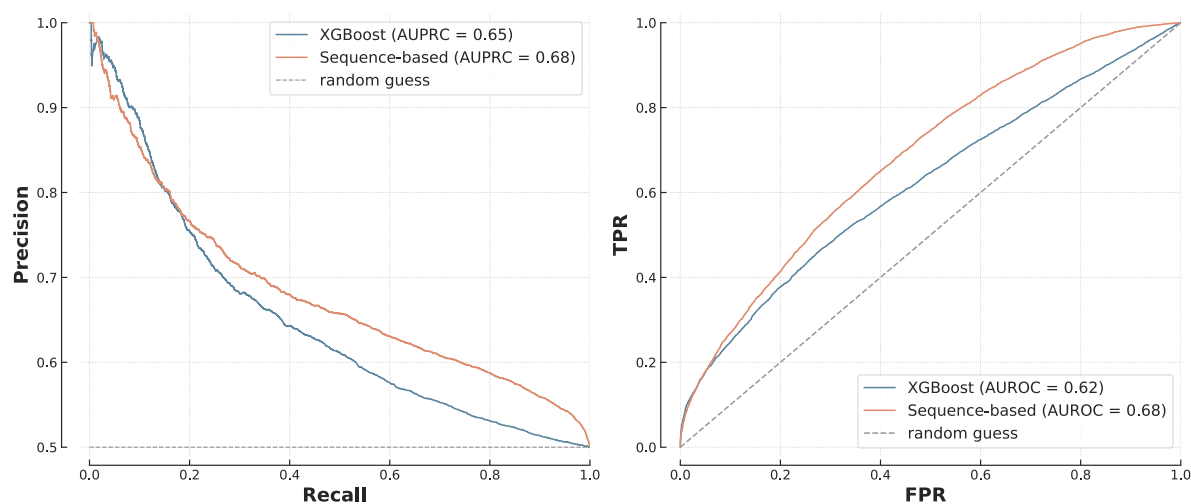
199 **Figure 3: Comparison of FG-based models on *T1* (n=24,898, 50% positive), with PR curves (left) and ROC curves**
 200 **(right) for all the models tested (top) and then only XGBoost and Logistic Regression for clarity (bottom).**

201 **Sequence models outperformed FG-based algorithms on known**
 202 **proteins**

203 The alternative to FG-based models is to use amino acids sequences as the input for a PPI
 204 prediction algorithm. We compared several deep learning architectures and reported the
 205 performance of an optimised Siamese neural network (**Methods, Figure 9 and Supplementary**
 206 **Figure 5**). Despite having access to no functional information about the proteins, the sequence
 207 model outperformed the best performing FG-based model, XGBoost, except at low recall and high
 208 precision (AUPRC=0.68 vs 0.65 and ROC curves significantly different, $p=9 \times 10^{-47}$) (**Figure 4**).
 209 However, while XGBoost was trained in only 30s with less than 0.01 kWh of energy, the deep
 210 learning approach trained for 1h10 with 0.62 kWh, emitting 22,000 times more greenhouse gases
 211 (GHGs). In addition, the performance of the sequence model was heavily affected by the choice
 212 of deep learning architecture and its hyper-parameters, such as number of layers or learning rate.
 213 These require extensive (and expensive) optimisation. Protein-level overlap had a significant

214 impact on these results. The model had an AUPRC of 0.68 on average, but 0.75 when restricted to
215 proteins present in both training and testing sets, and only 0.62 when there was no overlap
216 (**Supplementary Figure 5**). This demonstrates that (1) in the absence of specific adjustments, such
217 deep learning models are poorly suited to make predictions on previously unseen proteins and (2)
218 in-depth benchmarks like B4PPI are important to reliably measure performances. While this
219 comparison of FG-based and sequence-based models could indicate that deep learning is the best
220 approach to PPI prediction, and support the numerous similar claims in the literature, it could also
221 be the consequence of unaccounted-for biological properties of PPIs.

222



223

224 **Figure 4: Siamese network vs XGBoost on T1. The difference between the ROC curves was statistically significant**
225 **($p = 9 \times 10^{-47}$).**

226 **The role of network hubs is essential to PPI prediction**

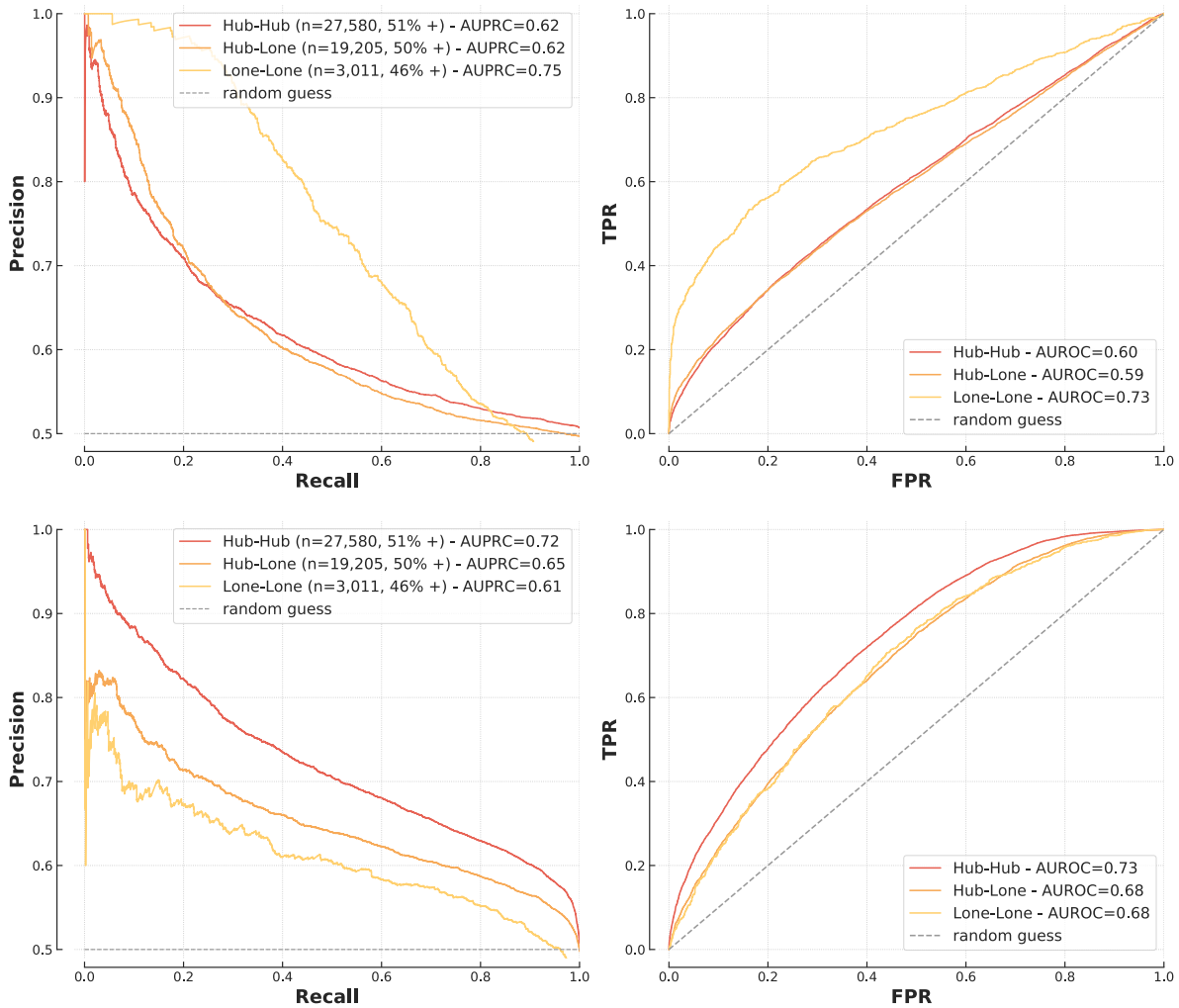
227 A scale-free topology has important biological implications [26] so we hypothesised that a one-
228 fits-all approach for hubs and lone proteins is unlikely to be optimal. In assessing interactions
229 between protein hubs (hub-hub), between a protein hub and a lone protein (hub-lone) and
230 between lone proteins (lone-lone) (**Methods**), we found a distinct pattern whereby FG-based
231 models had greater AUPRC and AUROC for interactions involving only lone proteins while
232 sequence-based models performed better for hubs (**Figure 5**).

233 These findings can be explained by the pre-processing of similarity measures for FG models.
234 Because of their central role in biological pathways, hubs are highly studied and therefore
235 annotated for many processes and localisations. For example in the training set, hubs have on
236 average 11.6 annotations for biological processes (significant feature in the logistic regression
237 model discussed above) while non-hubs only have 5.8 (median 6 vs 3). The same phenomenon is
238 observed for cellular compartments (6.2 vs 3.6 annotations on average). Because the similarity
239 measures used by the FG models quantify overlaps in annotations, hubs annotated for a large

240 number of processes provide little information about the probability of interaction, which can
241 explain why FG-based models perform best when hubs are not involved.

242 These results provide insight into the strengths of each approach and, importantly, show that a
243 PPI approach should be context specific, particularly with respect to the network topologies of
244 interest. Indeed, the apparent superiority of the deep learning model shown on **Figure 4** is largely
245 due to the composition of *T1*, made up of 70% of hub-hub or hub-lone interactions.

246



247

248

249 **Figure 5: Performance of XGBoost (top) and sequence-based model (bottom) on hubs and lone proteins.**

250 **Cross-species validation of PPI prediction models and relative** 251 **performances**

252 *S. cerevisiae* is a well-studied model organism with a known interactome and has been used
253 extensively for *in silico* PPI predictions [8], [37], [40]. We replicated the analyses presented above
254 on *S. cerevisiae* proteins and found that our findings regarding network topology and models'
255 relative performances were robust across species. The data was selected similarly to previously,

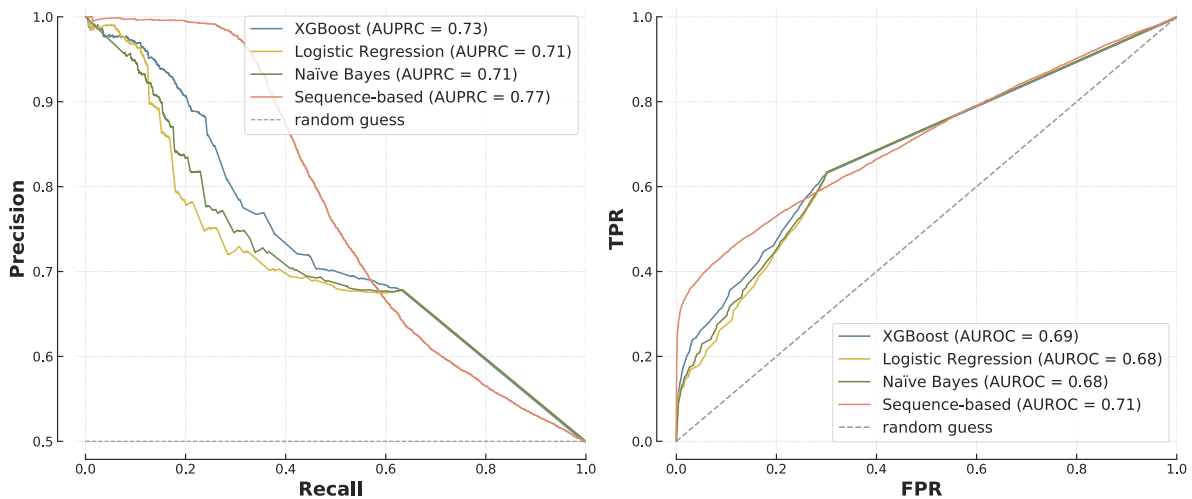
256 extracted and curated from IntAct and UniProt, but without data from HPA and Bgee as these
257 databases do not curate yeast (**Methods**).

258 As shown previously, all FG-based models had similar performances with AUPRC between 0.71
259 and 0.73 (**Figure 6**); however, in this analysis, the differences between XGBoost and other models
260 were statistically significant ($p = 2 \times 10^{-12}$ for Naïve Bayes and $p = 9 \times 10^{-31}$ for logistic regression).
261 The sequence model outperformed FG-based models in most cases ($p = 2 \times 10^{-6}$), except at high
262 recall (**Figure 6**). Second, similar to humans, FG-based models were not sensitive to protein-level
263 overlap while sequence-based models had different performances depending on the level of
264 overlap (**Supplementary Figure 6**). Finally, we found consistency regarding the role of network
265 hubs; FG-based models were better able to predict lone-lone protein interactions while the
266 sequence-based model was better at predicting interactions amongst protein hubs
267 (**Supplementary Figure 7**).

268 While experimental data on PPIs is readily available for humans and *S. cerevisiae*, many non-model
269 organisms lack data despite their biological relevance [41]. For these, cross-species predictions –
270 i.e. training a model on a species to make predictions on another – are of particular interest. We
271 showed that FG-based models are generally more suitable than sequence-based ones for this task.

272 We investigated whether models trained on yeast could be used to predict human PPIs, finding
273 that the yeast-trained FG-based models (logistic regression and XGBoost) achieved similar AUPRC
274 and AUROC as those which were human-trained to predict human PPIs ($p = 0.26$ for XGBoost)
275 (**Figure 7**). Conversely, the yeast-trained sequence model was unable to predict human PPIs
276 (AUPRC = 0.52 vs 0.68, $p = 3 \times 10^{-272}$). We observed the same phenomenon when using human-
277 trained models to make predictions on yeast (**Supplementary Figure 8**).

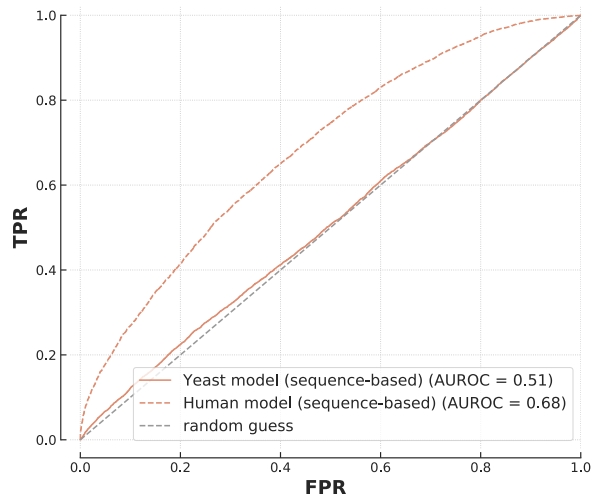
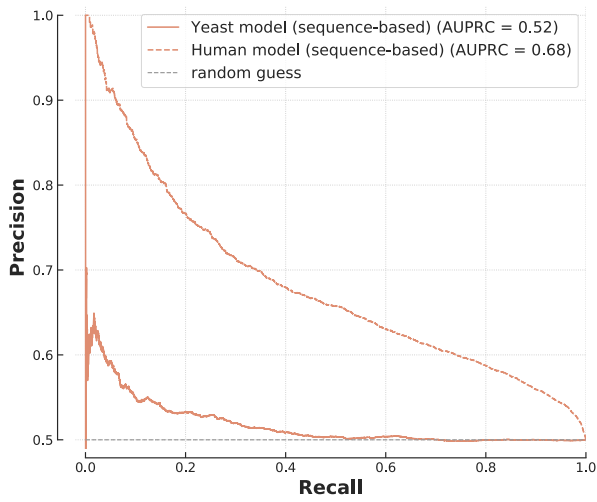
278



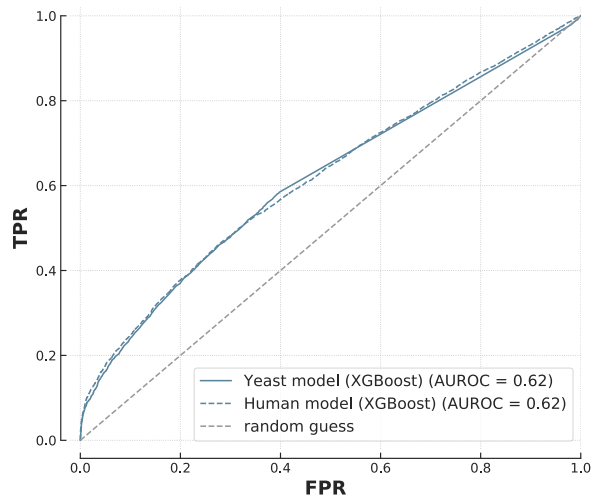
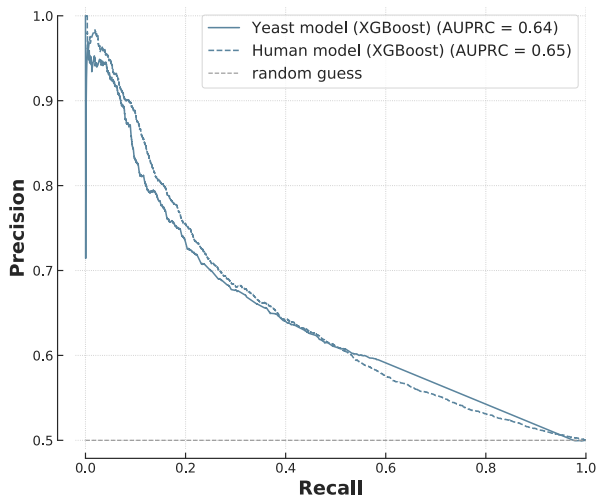
279

280 **Figure 6: Comparison of a range of models on the yeast testing set.**

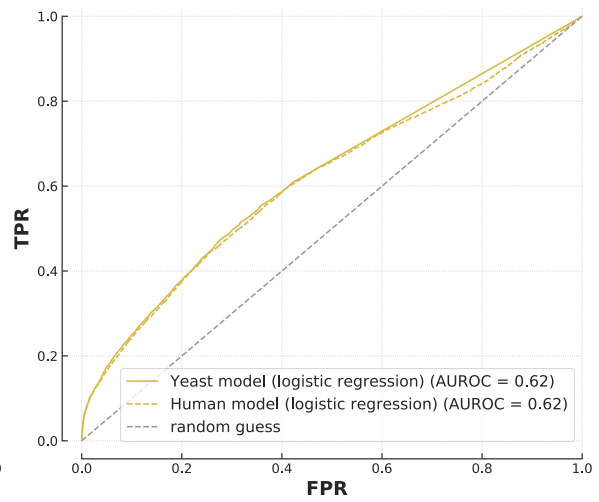
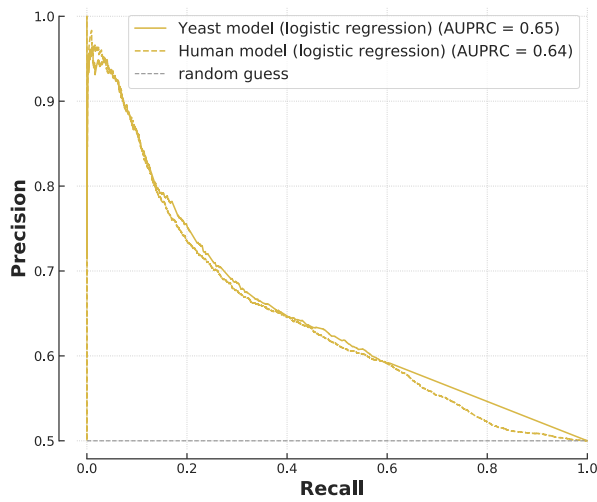
281



282



283



284

285 **Figure 7: Cross-species predictions. Models trained on human PPIs (dotted lines) and yeast PPIs (solid lines) were**
 286 **used to make predictions on the human testing set. The top plot is the sequence-based model, the other ones are**
 287 **FG-based (XGBoost in the middle, logistic regression at the bottom).**

289 In this work, we sought to identify and explain the strengths and weaknesses of a wide selection
290 of approaches to PPI prediction, and thereby provide the community with both a benchmarking
291 resource, insight into which PPI approach to select or trust in a particular scenario, and finally with
292 a set of well-studied PPI models which have also been validated across species.

293 The FG-based and sequence-based models are common for PPI prediction but are rarely directly
294 compared. In particular, it was unclear where the differences lie and if one approach should be
295 preferred today. We found that when using FG annotations, the choice of algorithm has little
296 impact on the predictions and a logistic regression performs close to the state-of-the-art while
297 providing clear insight into the decision-making process; here, colocalisation, common biological
298 processes and shared domains are the main indicators of interaction. The fact that a highly flexible
299 and non-linear model such as XGBoost performs similarly to logistic regression, making identical
300 predictions in 93% of cases, shows that performance is likely driven by the quality and the pre-
301 processing of the FG annotations instead of the modelling; once the similarity measures have been
302 calculated, there are limited non-linearities and a simple logistic regression achieves top
303 performance. Sequence-based models on the other hand need specifically optimised
304 architectures but achieve similarly high performances, if not higher in some settings, without any
305 biological information apart from amino acid sequences.

306 We found that the two approaches adapt to the presence of hubs and lone proteins differently
307 and show complementary strengths. While sequence-based models are mostly useful when hubs
308 are involved, FG-based models perform well for interactions between lone proteins. This simple
309 result offers important insight into the specificities of each approach and explains discrepancies in
310 reported performances in the literature, as the topology of the testing set has a large impact on
311 metrics. These results are not specific to human PPIs as the same conclusions were drawn from
312 analysis on *S. cerevisiae*. Cross-species predictions are instrumental to study non-model
313 organisms, and we showed that FG-based decision rules translate well to new species while
314 sequence-based models do not.

315 These observations are consistent with the way each algorithm learns. FG-based models make
316 predictions based on general, but less complex, rules about PPIs which translate well to new
317 proteins and new species. This is particularly useful considering that many proteins are still not
318 represented in interaction databases. On the other hand, sequence-based models have millions
319 of parameters which give them the flexibility to recognise individual proteins and learn specific
320 interaction patterns. Although this enables such models to make predictions without functional
321 information, it also limits high performance to proteins present in the training set. This likely
322 explains the poor results of sequence models on previously unseen proteins and cross-species
323 datasets. It is also consistent with the high performance of these models on network hubs, which
324 are overrepresented in training sets and therefore well captured by the models.

325 These analysis and results required a robust and reliable benchmarking pipeline. We designed the
326 open-source B4PPI, which accounts for a range of biological and statistical pitfalls. By being freely
327 accessible and using standard identifiers for proteins, B4PPI can be used by any researcher working
328 on *in silico* PPI prediction to assess performances and compare their approaches to the state-of-
329 the-art. An example reporting sheet is presented that includes relevant metrics, from PR and ROC
330 curves to runtime and carbon footprint, to ensure the models released can be trusted and
331 encourage wider use of PPI imputation for downstream analysis. B4PPI also comes with pre-
332 processed features to enable rapid development of new approaches.

333 Our study has limitations. We focused on the two most widely used approaches to PPI prediction,
334 namely FG-based and sequence-based; however, some alternative approaches have also been
335 proposed, using, for example, higher-level protein structures [37], phylogeny [42] and the
336 topology of existing networks [43]–[45], but the latter depends heavily on the quality of the existing
337 PPI networks. Most FG annotations are from gene ontologies which have a hierarchical structure
338 which we do not account for here, contrary to Armean et al. [46] for example. Moreover, we
339 analysed two common interactomes, human and yeast, yet there are many more. As
340 demonstrated though with the yeast dataset, similar benchmarks and analysis can be transferred
341 to other model organisms in a relatively straightforward manner.

342 We showed here the limits of classic sequence-based deep learning models for cross-species
343 predictions, but it is worth noting some recent deep learning models that have been successfully
344 used for cross-species predictions [47], [48] by including biological and chemical information about
345 amino acids as well as structural knowledge. The results presented in this work can hopefully guide
346 similar future work and help move this area further.

347 While a benchmarking standard for PPI prediction is needed, it is important to remember the
348 downsides of benchmarks, as demonstrated in computer vision or natural language processing. A
349 fixed set of metrics can motivate the community to overly focus on those, at the expense of
350 applicability and usefulness. To limit this, B4PPI includes a range of metrics but the relevant
351 indicators for each use-case should nonetheless be carefully considered.

352 The size and complexity of the PPI network makes *in silico* prediction tools indispensable, but it is
353 important to ensure that the models developed are reliable and readily available to the community
354 for downstream analysis and to give insights into biological pathways. For this, consistent and
355 reliable evaluation pipelines are necessary as well as a better understanding of what machine
356 learning models learn. The results presented here make key progress in both areas and facilitate
357 the development, evaluation and reliability of future PPI models.

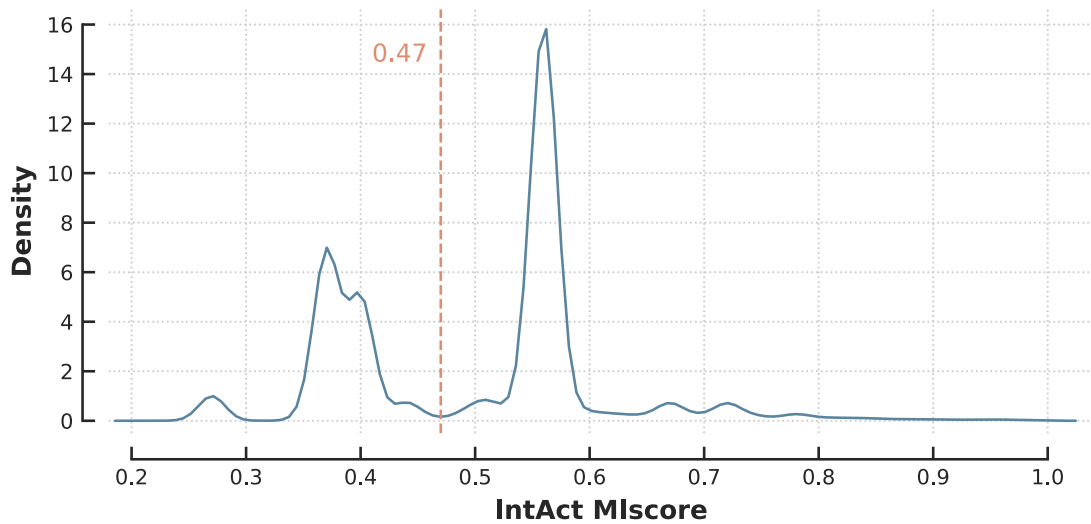
358 **METHODS**

359 **B4PPI-Human**

360 The data was obtained from large and professionally curated databases. This limits experimental
361 bias, as each interaction is based on several experiments, and leverages experts' knowledge in the
362 curation process. Standard UniProt IDs are used throughout to ensure maximum compatibilities.
363 Most of the manipulations were done in Python [49] with Jupyter Notebooks [50] using the Pandas
364 library [51], [52] and Numpy [53]. The plots were drawn using Matplotlib [54], Seaborn [55] and
365 the MetBrewer colour palettes [56]. All the code and final data are available on GitHub
366 (<https://github.com/Llannelongue/B4PPI>); some intermediary pre-processed datasets are not
367 available online due to file size limits but they can be recreated using the code available. Data is
368 available under Creative Commons Attribution (CC BY 4.0) License.

369 *Protein-protein interaction data*

370 The train machine learning algorithms, the quality of the gold standard is paramount. Data on PPIs
371 was obtained from IntAct [23] and downloaded from the EMBLE-EBI FTP server (timestamp:
372 15/10/2021). We restricted the data to human protein-protein interactions with UniProt IDs. To
373 reduce the risk of false positives, we removed spoke complex expansions (where the pairwise
374 interactions within a complex are unreliable) and interactions based on colocalisation only. This
375 quality control step leaves 128,790 PPIs, covering 15,506 proteins (out of 20,386 in UniProt). Based
376 on this dataset, we created an index of the number of recorded interactions per protein and made
377 a list of hubs (highly connected proteins). In line with the literature, hubs are defined as the 20%
378 of proteins with the most interactions [57], which here is equivalent to proteins with more than
379 21 partners. The quality of the interactions is assessed further by looking at the MIscore [58], a
380 quality score based on the manual curation of the interactions that takes into account the
381 detection method, the interaction type and the number of publications reporting it. In case of
382 duplicated PPIs, the highest MIscore was used. When looking at the distribution of the MIscores
383 in the dataset (**Figure 8**), a natural threshold of 0.47 is visible, which restrict the dataset to 78,229
384 interactions, covering 12,026 proteins.



385

386 **Figure 8: Distribution of the MIscore in IntAct.**

387 *Functional genomics annotations and amino acids sequences*

388 Protein sequences in humans are well documented and can be obtained from UniProt, but FG
 389 features can be more challenging as they should be diverse (i.e. cover a wide range of properties),
 390 of high-quality and have high coverage (i.e. few missing proteins). For the same reasons as
 391 described above, aggregated, manually curated and professionally reviewed databases are
 392 preferred. Based on features that have been successfully used for the task before, it is relevant to
 393 include information about cellular and tissue localisation, biological functions and gene expression
 394 patterns [5], [8], [10], [37].

395 One of the main databases on proteins is UniProt [59] and in particular its knowledgebase
 396 UniProtKB. Swiss-Prot, the section of UniProtKB that is reviewed and manually curated, is used in
 397 this work to ensure optimal quality. The data from Swiss-Prot is downloaded through their API by
 398 restricting to reviewed, non-obsolete, human proteins (last download is 09/11/2021). The
 399 different columns are then cleaned to extract the information of interest in a standardised format,
 400 and we use UniProt IDs throughout. There is information for 20,386 proteins and more details
 401 about each feature are in **Table 1**. UniProt’s API is also used to map UniProt IDs between different
 402 databases and to map outdated IDs. In particular, we extract amino acids sequences for each
 403 protein, more than 95% of which come from the translation of coding sequences submitted to the
 404 International Nucleotide Sequence Database Collaboration [60]. Annotated domains and motifs
 405 are also included in the database. Additionally, we extract gene ontology (GO) annotations of
 406 biological processes, cellular components and molecular functions. For each protein, each of the
 407 FG features is represented as a bag-of-words, i.e. a sparse vector of length the number of
 408 annotations in the database.

409 When working with gene expression data, both biological and technical noise need to be
 410 accounted for correctly. The Bgee public repository [33] does that by regrouping curated healthy
 411 wild-type standardised gene expression patterns. The human data is mainly from GTEx v6

412 (phs000424.v6.p1), with an added layer of manual curation to remove unhealthy subjects. For a
 413 gene, the final data provides binary calls of presence or absence of expression for each
 414 combination of anatomical entity and developmental stage. We downloaded the database from
 415 their FTP server (version 14.2) and obtained information for 59,777 genes, 320 anatomical entities
 416 and 33 developmental stages, which leads to 1,147 stage/entity combinations. The Bgee entries
 417 are matched to the UniProt IDs using UniProt’s own mapping table.

418 The Human Protein Atlas (HPA) [31], [32] provides data mapping human proteins to tissues and
 419 cells. In particular, we used the Tissue Atlas [31] that presents the distribution of proteins in tissues
 420 and cell types and the Cell Atlas [32] that contains the distribution across subcellular locations.
 421 The Tissue Atlas contains data similar to Bgee, but the overlap is likely to be limited as the two
 422 databases only share GTEx RNA-seq data. While Bgee has a more thorough curation process, HPA
 423 contains a lot of original in-house experimental results, which justifies the inclusion of both data
 424 sources. We downloaded the HPA data from their website (release 20.1, Ensembl version 92.38).
 425 Despite its name, the data in HPA is identified by Ensembl gene IDs, which are mapped to UniProt
 426 IDs using UniProt’s API. We restricted the dataset to the reviewed proteins present in Swiss-Prot
 427 and to ensure the quality of annotations, we discarded the entries HPA annotated as “uncertain”.
 428 For the tissue IHC data, we mapped expression levels to numerical values (high=3, medium=2,
 429 low=1 and “not detected”=-1) with untested tissues being mapped to 0. Similar pre-processing
 430 was used for the consensus RNA-seq data and the subcellular location.

431 **Table 1: Features used to train human models (GO = Gene Ontology).**

Feature	Number of different annotations	Missing values (/20,386)	Source
Biological processes (GO)	12,248	3,338	UniProt [59]
Cellular components (GO)	1,754	1,765	UniProt
Molecular functions (GO)	4,346	4,552	UniProt
Domains	2,313	11,815	UniProt
Motifs	819	18,103	UniProt
Sequence	N/A	0	UniProt
Gene expression profile	1,147	1,296	Bgee [33]
Tissue IHC data	62	9,536	HPA [31], [32]
Tissue and cell type	189	9,536	HPA
RNA-seq	61	1,448	HPA
Subcellular location	33	7,820	HPA

432

433

434 *Pre-processing to measure features similarity*

435 For a protein, each FG feature was represented as a vector, of length the number of annotations.
436 To measure the feature-specific similarity between two proteins, we compared the two vectors
437 using cosine similarity [61], a popular tool widely used for similar tasks in Natural Language
438 Processing. For two vectors $A = (A_i)$ and $B = (B_i)$, their cosine similarity $CS(A, B)$ is:

$$439 \quad CS(A, B) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_i A_i B_i}{\sqrt{\sum_i A_i^2} \sqrt{\sum_i B_i^2}}$$

440 As a result, for each of the 207,784,305 possible pairs of proteins, we obtained 12 similarity
441 features: biological processes, cell components, molecular function, domains and motifs from
442 UniProt, gene expression from Bgee, tissue/cell expression, tissue expression, RNA-seq expression
443 and subcellular locations from the Human Protein Atlas (**Table 1**).

444 *Creation of the gold standard*

445 The PPIs obtained from IntAct are divided between a training set and two testing sets. First, a set
446 of 1,562 proteins (13%) was randomly set aside to ensure some unseen proteins are present in
447 the testing set; the necessity of this is shown in **Figure 2**. The dataset was then randomly divided
448 under this constraint and included 53,331 PPIs in the training set, 12,449 in $T1$ and the same in $T2$
449 (**Supplementary Table 1**).

450 The negative examples (i.e. non interacting proteins) are obtained using random sampling among
451 all the possible pairs, excluding any pair that has been observed experimentally to limit the risk of
452 false negative. For the training set, balanced sampling is used [27] to favour learning, which means
453 that the probably of sampling a protein for the negative set is proportionate to its frequency in
454 the positive set. For $T1$ and $T2$, we used uniform sampling (all proteins have the same probability
455 of sampling) to limit the risk of bias. The training set and $T1$ both have 50% of positive examples,
456 while $T2$ has ten times more non-interacting proteins than interacting ones (**Supplementary Table**
457 **1**).

458 To investigate how models deal with different network topologies, especially hubs and lone
459 proteins, we had to create a separate testing set to ensure sufficient sample size in each category
460 (hub-hub, hub-lone and lone-lone interactions). We do so by aggregating PPIs from $T1$ and $T2$, and
461 using balanced sampling for the non-interacting proteins. This results in 49,796 pairs (50%
462 positive) (**Supplementary Table 2**).

464 *S. cerevisiae data*

465 The pipeline describe above was also followed for the *S. cerevisiae* data. UniProt lists 6,721 yeast
466 proteins and the same information as for humans (

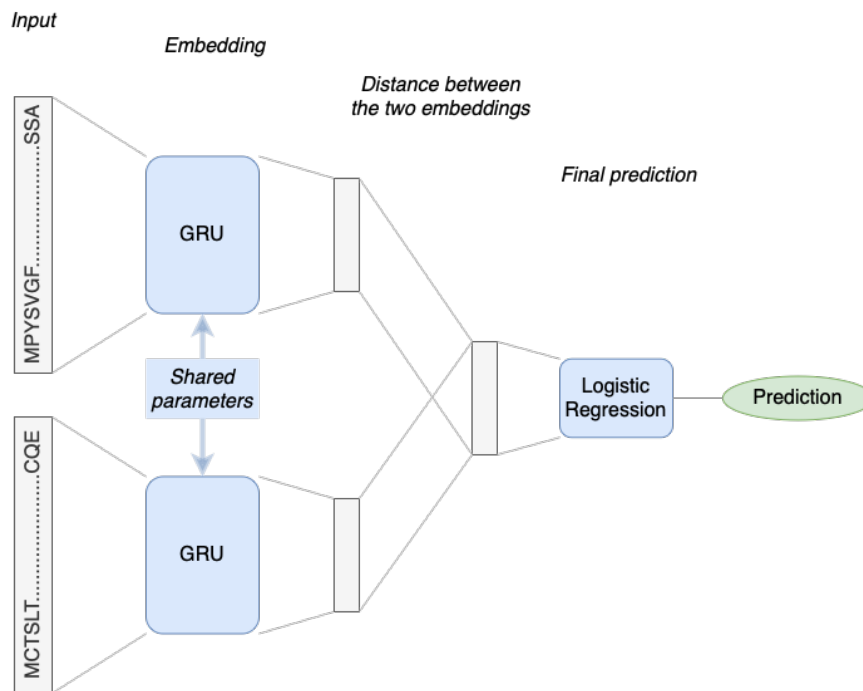
467 **Supplementary Table 3**) but HPA and Bgee do not include data for this organism. PPIs were
468 obtained from IntAct following the same procedure, although no selection based on MIScores was
469 made considering the absence of an obvious choice when looking at the distribution
470 (**Supplementary Figure 9**). The final PPI dataset comprised 43,068 interactions covering 5,679
471 proteins.

472 The split between training and testing sets was done similarly by setting aside 737 proteins for
473 testing and then randomly allocating PPIs to keep 30,369 PPIs for training. Because there is fewer
474 data on yeast, and only one testing set is needed to replicate the analysis conducted on humans,
475 dividing the remaining 12,699 further between *T1* and *T2* is not suitable here. But if the goal was
476 to measure generalisability of a yeast model, this could be easily done.

477 **Training**

478 FG-based machine learning models were trained using the scikit-learn library [62]. For models that
479 cannot deal with missing data, mean imputation was used (**Supplementary Table 4**). Hyper-
480 parameter search was done using Weight-and-Bias's Bayesian method [63] to find the optimal
481 settings of each algorithm in a reasonable time. All hyperparameter choices are in **Supplementary**
482 **Table 4** and **Supplementary Table 5**.

483 Deep learning models were trained using PyTorch Lightning [64], [65]. The Siamese architecture
484 [66], [67] was composed of a bidirectional Gated Recurrent Unit (GRU) [68] followed by a linear
485 output (**Figure 9**). Long-Short Term Memory networks (LSTM) [69] and Convolutional Neural
486 Networks (CNN) [70] were also tested, but GRU was preferred because of runtime efficiency and
487 its ability to account for proteins of various lengths. Full parameters are in **Supplementary Table**
488 **4**, **Supplementary Table 5** and in the open-source code.



489

490 Figure 9: Diagram of the deep learning architecture used to predict interactions from a pair of protein sequences.

491 Evaluation

492 The Receiver Operating Characteristic (ROC) and the Precision-Recall (PR) curves are
 493 complementary options for PPI prediction. While the ROC curve is unaffected by the prevalence
 494 of interacting proteins, a benefit as the true prevalence of PPIs is mostly unknown, it also means
 495 that both classes are considered equally, whereas often, PPIs are more interesting than non-
 496 interacting proteins. This is addressed by the PR curve where precision puts an emphasis on
 497 positive examples.

498 Both curves are reported alongside their respective Areas Under the Curve (AUC). To statistically
 499 compare ROC curves for a same testing set, we used a DeLong nonparametric test [71] and
 500 reported the p-value. We corrected for multiple testing by using a conservative significance
 501 threshold of 5×10^{-4} , corresponding to a Bonferroni correction for 100 pairwise comparisons [72].

502 Carbon footprint of this project

503 We used the Green Algorithms calculator (v2.1) [30] and estimated that the carbon footprint of
 504 this project was 51 kgCO₂e, which corresponds to 4.7 tree-years. We did our best to minimise
 505 greenhouse gas emissions in the first place, and as a commitment to the reduction of the carbon
 506 footprint of computational research, we funded tree planting in the east of England region through
 507 carbonfootprint.com. These trees are estimated to sequester 1 tonne of CO₂ in their lifetime,
 508 almost 20 times the emissions of this study.

509 **ACKNOWLEDGMENTS**

510 L.L. was supported by the University of Cambridge MRC DTP (MR/S502443/1). M.I. was supported
511 by the Munz Chair of Cardiovascular Prediction and Prevention and the NIHR Cambridge
512 Biomedical Research Centre (BRC-1215-20014)[*]. This work was supported by core funding from
513 the British Heart Foundation (RG/13/13/30194; RG/18/13/33946). This work was also supported
514 by Health Data Research UK, which is funded by the UK Medical Research Council, Engineering and
515 Physical Sciences Research Council, Economic and Social Research Council, Department of Health
516 and Social Care (England), Chief Scientist Office of the Scottish Government Health and Social Care
517 Directorates, Health and Social Care Research and Development Division (Welsh Government),
518 Public Health Agency (Northern Ireland), British Heart Foundation and Wellcome. This work was
519 performed using resources provided by the Cambridge Service for Data Driven Discovery (CSD3)
520 operated by the University of Cambridge Research Computing Service (www.csd3.cam.ac.uk),
521 provided by Dell EMC and Intel using Tier-2 funding from the Engineering and Physical Sciences
522 Research Council (capital grant EP/P020259/1), and DiRAC funding from the Science and
523 Technology Facilities Council (www.dirac.ac.uk).

524 *The views expressed are those of the author(s) and not necessarily those of the NIHR or the
525 Department of Health and Social Care.

526

527

- 529 [1] H. Goehler *et al.*, 'A Protein Interaction Network Links GIT1, an Enhancer of Huntingtin
530 Aggregation, to Huntington's Disease', *Molecular Cell*, vol. 15, no. 6, pp. 853–865, Sep. 2004, doi:
531 10.1016/j.molcel.2004.09.016.
- 532 [2] A. Vinayagam *et al.*, 'A Directed Protein Interaction Network for Investigating Intracellular
533 Signal Transduction', *Science Signaling*, vol. 4, no. 189, pp. rs8–rs8, Sep. 2011, doi:
534 10.1126/scisignal.2001699.
- 535 [3] M. Bakail and F. Ochsenbein, 'Targeting protein–protein interactions, a wide open field for
536 drug design', *Comptes Rendus Chimie*, vol. 19, no. 1–2, pp. 19–27, Jan. 2016, doi:
537 10.1016/j.crci.2015.12.004.
- 538 [4] S. Rapposelli, E. Gaudio, F. Bertozzi, and S. Gul, 'Editorial: Protein–Protein Interactions:
539 Drug Discovery for the Future', *Front. Chem.*, vol. 9, p. 811190, Nov. 2021, doi:
540 10.3389/fchem.2021.811190.
- 541 [5] R. Jansen, 'A Bayesian Networks Approach for Predicting Protein-Protein Interactions from
542 Genomic Data', *Science*, vol. 302, no. 5644, pp. 449–453, Oct. 2003, doi:
543 10.1126/science.1087361.
- 544 [6] L. V. Zhang, S. L. Wong, O. D. King, and F. P. Roth, 'Predicting co-complexed protein pairs
545 using genomic and proteomic data integration', *BMC Bioinformatics*, p. 15, 2004.
- 546 [7] X.-W. Chen and M. Liu, 'Prediction of protein–protein interactions using random decision
547 forest framework', *Bioinformatics*, vol. 21, no. 24, pp. 4394–4400, Dec. 2005, doi:
548 10.1093/bioinformatics/bti721.
- 549 [8] A. Ben-Hur and W. S. Noble, 'Kernel methods for predicting protein-protein interactions',
550 *Bioinformatics*, vol. 21, no. Suppl 1, pp. i38–i46, Jun. 2005, doi: 10.1093/bioinformatics/bti1016.
- 551 [9] M. S. Scott and G. J. Barton, 'Probabilistic prediction and ranking of human protein-protein
552 interactions', *BMC Bioinformatics*, vol. 8, no. 1, p. 239, Jul. 2007, doi: 10.1186/1471-2105-8-239.
- 553 [10] M. Kotlyar *et al.*, 'In silico prediction of physical protein interactions and characterization
554 of interactome orphans', *Nature Methods*, vol. 12, no. 1, pp. 79–84, Jan. 2015, doi:
555 10.1038/nmeth.3178.
- 556 [11] Y. Murakami and K. Mizuguchi, 'Homology-based prediction of interactions between
557 proteins using Averaged One-Dependence Estimators', *BMC Bioinformatics*, vol. 15, no. 1, p. 213,
558 Jun. 2014, doi: 10.1186/1471-2105-15-213.
- 559 [12] Z.-H. You, M. Zhou, X. Luo, and S. Li, 'Highly Efficient Framework for Predicting Interactions
560 Between Proteins', *IEEE Transactions on Cybernetics*, vol. 47, no. 3, pp. 731–743, Mar. 2017, doi:
561 10.1109/TCYB.2016.2524994.

- 562 [13] M. Chen *et al.*, 'Multifaceted protein–protein interaction prediction based on Siamese
563 residual RCNN', *Bioinformatics*, vol. 35, no. 14, pp. i305–i314, Jul. 2019, doi:
564 10.1093/bioinformatics/btz328.
- 565 [14] M. Kotlyar, A. E. M. Rossos, and I. Jurisica, 'Prediction of Protein-Protein Interactions',
566 *Current Protocols in Bioinformatics*, vol. 60, no. 1, p. 8.2.1-8.2.14, 2017, doi: 10.1002/cpbi.38.
- 567 [15] S. Mangul *et al.*, 'Systematic benchmarking of omics computational tools', *Nat Commun*,
568 vol. 10, no. 1, p. 1393, Dec. 2019, doi: 10.1038/s41467-019-09406-4.
- 569 [16] T. Sun, B. Zhou, L. Lai, and J. Pei, 'Sequence-based prediction of protein protein interaction
570 using a deep-learning algorithm', *BMC Bioinformatics*, vol. 18, no. 1, p. 277, May 2017, doi:
571 10.1186/s12859-017-1700-2.
- 572 [17] F. Li, F. Zhu, X. Ling, and Q. Liu, 'Protein Interaction Network Reconstruction Through
573 Ensemble Deep Learning With Attention Mechanism', *Front Bioeng Biotechnol*, vol. 8, p. 390, May
574 2020, doi: 10.3389/fbioe.2020.00390.
- 575 [18] A. Ben-Hur and W. Noble, 'Choosing negative examples for the prediction of protein-
576 protein interactions', *BMC Bioinformatics*, vol. 7, no. Suppl 1, p. S2, 2006, doi: 10.1186/1471-2105-
577 7-S1-S2.
- 578 [19] P. Blohm *et al.*, 'Negatome 2.0: a database of non-interacting proteins derived by literature
579 mining, manual annotation and protein structure analysis', *Nucleic Acids Res*, vol. 42, no. Database
580 issue, pp. D396–D400, Jan. 2014, doi: 10.1093/nar/gkt1079.
- 581 [20] Y. Park and E. M. Marcotte, 'Revisiting the negative example sampling problem for
582 predicting protein–protein interactions', *Bioinformatics*, vol. 27, no. 21, pp. 3024–3028, Nov.
583 2011, doi: 10.1093/bioinformatics/btr514.
- 584 [21] L. Hu, X. Wang, Y.-A. Huang, P. Hu, and Z.-H. You, 'A survey on computational models for
585 predicting protein–protein interactions', *Briefings in Bioinformatics*, vol. 22, no. 5, Sep. 2021, doi:
586 10.1093/bib/bbab036.
- 587 [22] Y. Park and E. M. Marcotte, 'Flaws in evaluation schemes for pair-input computational
588 predictions', *Nature Methods*, vol. 9, no. 12, pp. 1134–1136, Dec. 2012, doi: 10.1038/nmeth.2259.
- 589 [23] S. Orchard *et al.*, 'The MIntAct project--IntAct as a common curation platform for 11
590 molecular interaction databases.', *Nucleic Acids Res*, vol. 42, no. Database issue, pp. D358-63, Jan.
591 2014, doi: 10.1093/nar/gkt1115.
- 592 [24] S. Orchard *et al.*, 'Protein interaction data curation: the International Molecular Exchange
593 (IMEx) consortium', *Nature Methods*, vol. 9, no. 4, pp. 345–350, Apr. 2012, doi:
594 10.1038/nmeth.1931.
- 595 [25] M. Agrawal, M. Zitnik, and J. Leskovec, 'Large-scale analysis of disease pathways in the
596 human interactome', *Pac Symp Biocomput*, vol. 23, pp. 111–122, 2018.

- 597 [26] EMBL-EBI, 'Properties of PPINs: scale-free networks | Network analysis of protein
598 interaction data'. [https://www.ebi.ac.uk/training/online/courses/network-analysis-of-protein-
interaction-data-an-introduction/protein-protein-interaction-networks/properties-of-ppins-
scale-free-networks/](https://www.ebi.ac.uk/training/online/courses/network-analysis-of-protein-
599 interaction-data-an-introduction/protein-protein-interaction-networks/properties-of-ppins-
600 scale-free-networks/) (accessed Nov. 29, 2021).
- 601 [27] J. Yu, M. Guo, C. J. Needham, Y. Huang, L. Cai, and D. R. Westhead, 'Simple sequence-based
602 kernels do not predict protein-protein interactions', *Bioinformatics*, vol. 26, no. 20, pp. 2610–
603 2614, Oct. 2010, doi: 10.1093/bioinformatics/btq483.
- 604 [28] F. J. Provost, T. Fawcett, and R. Kohavi, 'The Case against Accuracy Estimation for
605 Comparing Induction Algorithms', in *Proceedings of the Fifteenth International Conference on
606 Machine Learning*, San Francisco, CA, USA, Jul. 1998, pp. 445–453.
- 607 [29] J. Grealey *et al.*, 'The carbon footprint of bioinformatics', *Bioinformatics*, preprint, Mar.
608 2021. doi: 10.1101/2021.03.08.434372.
- 609 [30] L. Lannelongue, J. Grealey, and M. Inouye, 'Green Algorithms: Quantifying the Carbon
610 Footprint of Computation', *Advanced Science*, vol. 8, no. 12, p. 2100707, 2021, doi:
611 10.1002/advs.202100707.
- 612 [31] M. Uhlén *et al.*, 'Tissue-based map of the human proteome', *Science*, vol. 347, no. 6220,
613 Jan. 2015, doi: 10.1126/science.1260419.
- 614 [32] P. J. Thul *et al.*, 'A subcellular map of the human proteome', *Science*, vol. 356, no. 6340,
615 May 2017, doi: 10.1126/science.aal3321.
- 616 [33] F. B. Bastian *et al.*, 'The Bgee suite: integrated curated expression atlas and comparative
617 transcriptomics in animals', *Nucleic Acids Research*, vol. 49, no. D1, pp. D831–D847, Jan. 2021, doi:
618 10.1093/nar/gkaa793.
- 619 [34] M. E. Maron, 'Automatic Indexing: An Experimental Inquiry', *J. ACM*, vol. 8, no. 3, pp. 404–
620 417, Jul. 1961, doi: 10.1145/321075.321084.
- 621 [35] L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen, *Classification and Regression Trees*.
622 Taylor & Francis, 1984. [Online]. Available: <https://books.google.fr/books?id=JwQx-WOmSyQC>
- 623 [36] L. Breiman, 'Random Forests', *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001, doi:
624 10.1023/A:1010933404324.
- 625 [37] Q. C. Zhang *et al.*, 'Structure-based prediction of protein–protein interactions on a
626 genome-wide scale', *Nature*, vol. 490, no. 7421, pp. 556–560, Oct. 2012, doi:
627 10.1038/nature11503.
- 628 [38] T. Chen and C. Guestrin, 'XGBoost: A Scalable Tree Boosting System', in *Proceedings of the
629 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '16*,
630 San Francisco, California, USA, 2016, pp. 785–794. doi: 10.1145/2939672.2939785.

- 631 [39] A. Ogunleye and Q.-G. Wang, 'XGBoost Model for Chronic Kidney Disease Diagnosis',
632 *IEEE/ACM Trans. Comput. Biol. and Bioinf.*, vol. 17, no. 6, pp. 2131–2140, Nov. 2020, doi:
633 10.1109/TCBB.2019.2911071.
- 634 [40] E. Sprinzak and H. Margalit, 'Correlated sequence-signatures as markers of protein-protein
635 interaction' 1Edited by G. von Heijne', *Journal of Molecular Biology*, vol. 311, no. 4, pp. 681–692,
636 Aug. 2001, doi: 10.1006/jmbi.2001.4920.
- 637 [41] J. J. Russell *et al.*, 'Non-model model organisms', *BMC Biol*, vol. 15, no. 1, pp. 55, s12915-
638 017-0391–5, Dec. 2017, doi: 10.1186/s12915-017-0391-5.
- 639 [42] G. Marmier, M. Weigt, and A.-F. Bitbol, 'Phylogenetic correlations can suffice to infer
640 protein partners from sequences', *PLoS Comput Biol*, vol. 15, no. 10, p. e1007179, Oct. 2019, doi:
641 10.1371/journal.pcbi.1007179.
- 642 [43] X. Zhong and J. C. Rajapakse, 'Graph embeddings on gene ontology annotations for
643 protein–protein interaction prediction', *BMC Bioinformatics*, vol. 21, no. Suppl 16, p. 560, Dec.
644 2020, doi: 10.1186/s12859-020-03816-8.
- 645 [44] L. Becchetti, A. Fazzone, and L. Martini, 'Network and Sequence-Based Prediction of
646 Protein-Protein Interactions', *arXiv:2107.03694 [cs, q-bio]*, Jul. 2021, Accessed: Aug. 16, 2021.
647 [Online]. Available: <http://arxiv.org/abs/2107.03694>
- 648 [45] I. A. Kovács *et al.*, 'Network-based prediction of protein interactions', *Nat Commun*, vol.
649 10, no. 1, p. 1240, Dec. 2019, doi: 10.1038/s41467-019-09177-y.
- 650 [46] I. M. Armean, K. S. Lilley, M. W. B. Trotter, N. C. V. Pilkington, and S. B. Holden, 'Co-complex
651 protein membership evaluation using Maximum Entropy on GO ontology and InterPro
652 annotation', *Bioinformatics*, vol. 34, no. 11, pp. 1884–1892, Jun. 2018, doi:
653 10.1093/bioinformatics/btx803.
- 654 [47] S. Mahapatra and S. S. Sahu, 'Improved prediction of protein–protein interaction using a
655 hybrid of functional-link Siamese neural network and gradient boosting machines', *Briefings in*
656 *Bioinformatics*, no. bbab255, Jul. 2021, doi: 10.1093/bib/bbab255.
- 657 [48] S. Sledzieski, R. Singh, L. Cowen, and B. Berger, 'D-SCRIPT translates genome to phenome
658 with sequence-based, structure-aware, genome-scale predictions of protein-protein interactions',
659 *cels*, vol. 0, no. 0, Sep. 2021, doi: 10.1016/j.cels.2021.08.010.
- 660 [49] 'The Python Language Reference — Python 3.10.1 documentation'.
661 <https://docs.python.org/3/reference/> (accessed Jan. 10, 2022).
- 662 [50] 'Jupyter Project Documentation — Jupyter Documentation 4.1.1 alpha documentation'.
663 <https://docs.jupyter.org/en/latest/> (accessed Jan. 10, 2022).
- 664 [51] J. Reback *et al.*, *pandas-dev/pandas: Pandas 1.0.3*. Zenodo, 2020. doi:
665 10.5281/zenodo.3715232.

- 666 [52] W. McKinney, 'Data Structures for Statistical Computing in Python', Austin, Texas, 2010,
667 pp. 56–61. doi: 10.25080/Majora-92bf1922-00a.
- 668 [53] C. R. Harris *et al.*, 'Array programming with NumPy', *Nature*, vol. 585, no. 7825, pp. 357–
669 362, Sep. 2020, doi: 10.1038/s41586-020-2649-2.
- 670 [54] J. D. Hunter, 'Matplotlib: A 2D Graphics Environment', *Comput. Sci. Eng.*, vol. 9, no. 3, pp.
671 90–95, 2007, doi: 10.1109/MCSE.2007.55.
- 672 [55] M. Waskom, 'seaborn: statistical data visualization', *JOSS*, vol. 6, no. 60, p. 3021, Apr. 2021,
673 doi: 10.21105/joss.03021.
- 674 [56] 'GitHub - BlakeRMills/MetBrewer: Color palette package in R inspired by works at the
675 Metropolitan Museum of Art in New York'. <https://github.com/BlakeRMills/MetBrewer> (accessed
676 Jan. 10, 2022).
- 677 [57] G. Jin, S. Zhang, X.-S. Zhang, and L. Chen, 'Hubs with Network Motifs Organize Modularity
678 Dynamically in the Protein-Protein Interaction Network of Yeast', *PLoS ONE*, vol. 2, no. 11, p.
679 e1207, Nov. 2007, doi: 10.1371/journal.pone.0001207.
- 680 [58] B. Aranda *et al.*, 'PSICQUIC and PSIScore: accessing and scoring molecular interactions',
681 *Nat Methods*, vol. 8, no. 7, pp. 528–529, Jun. 2011, doi: 10.1038/nmeth.1637.
- 682 [59] The UniProt Consortium, 'UniProt: the universal protein knowledgebase in 2021', *Nucleic
683 Acids Research*, vol. 49, no. D1, pp. D480–D489, Jan. 2021, doi: 10.1093/nar/gkaa1100.
- 684 [60] M. Arita, I. Karsch-Mizrachi, G. Cochrane, and on behalf of the International Nucleotide
685 Sequence Database Collaboration, 'The international nucleotide sequence database
686 collaboration', *Nucleic Acids Research*, vol. 49, no. D1, pp. D121–D124, Jan. 2021, doi:
687 10.1093/nar/gkaa967.
- 688 [61] A. Singhal, 'Modern Information Retrieval: A Brief Overview', p. 9, 2001.
- 689 [62] F. Pedregosa *et al.*, 'Scikit-learn: Machine Learning in Python', *J. Mach. Learn. Res.*, vol. 12,
690 no. null, pp. 2825–2830, Nov. 2011.
- 691 [63] L. Biewald, 'Experiment Tracking with Weights and Biases'. 2020. [Online]. Available:
692 <https://www.wandb.com/>
- 693 [64] W. Falcon and The PyTorch Lightning team, *PyTorch Lightning*. 2019. doi:
694 10.5281/zenodo.3828935.
- 695 [65] A. Paszke *et al.*, 'PyTorch: An Imperative Style, High-Performance Deep Learning Library',
696 in *Advances in Neural Information Processing Systems*, 2019, vol. 32. Accessed: Jan. 10, 2022.
697 [Online]. Available:
698 <https://proceedings.neurips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740->
699 [Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html)
- 700 [66] J. Bromley, I. Guyon, Y. LeCun, E. Säcker, and R. Shah, 'Signature Verification using a
701 "Siamese" Time Delay Neural Network', in *Advances in Neural Information Processing Systems 6*,

702 J. D. Cowan, G. Tesauro, and J. Alspecter, Eds. Morgan-Kaufmann, 1994, pp. 737–744. Accessed:
703 Oct. 21, 2019. [Online]. Available: [http://papers.nips.cc/paper/769-signature-verification-using-a-](http://papers.nips.cc/paper/769-signature-verification-using-a-siamese-time-delay-neural-network.pdf)
704 [siamese-time-delay-neural-network.pdf](http://papers.nips.cc/paper/769-signature-verification-using-a-siamese-time-delay-neural-network.pdf)

705 [67] S. Chopra, R. Hadsell, and Y. LeCun, ‘Learning a Similarity Metric Discriminatively, with
706 Application to Face Verification’, in *2005 IEEE Computer Society Conference on Computer Vision*
707 *and Pattern Recognition (CVPR’05)*, San Diego, CA, USA, 2005, vol. 1, pp. 539–546. doi:
708 10.1109/CVPR.2005.202.

709 [68] K. Cho *et al.*, ‘Learning Phrase Representations using RNN Encoder-Decoder for Statistical
710 Machine Translation’, *arXiv:1406.1078 [cs, stat]*, Sep. 2014, Accessed: Nov. 29, 2021. [Online].
711 Available: <http://arxiv.org/abs/1406.1078>

712 [69] S. Hochreiter and J. Schmidhuber, ‘Long Short-Term Memory’, *Neural Computation*, vol. 9,
713 no. 8, pp. 1735–1780, Nov. 1997, doi: 10.1162/neco.1997.9.8.1735.

714 [70] Y. LeCun *et al.*, ‘Handwritten Digit Recognition with a Back-Propagation Network’, in
715 *Advances in Neural Information Processing Systems*, 1990, vol. 2. Accessed: Feb. 04, 2022.
716 [Online]. Available:
717 <https://papers.nips.cc/paper/1989/hash/53c3bce66e43be4f209556518c2fcb54-Abstract.html>

718 [71] E. R. DeLong, D. M. DeLong, and D. L. Clarke-Pearson, ‘Comparing the Areas under Two or
719 More Correlated Receiver Operating Characteristic Curves: A Nonparametric Approach’,
720 *Biometrics*, vol. 44, no. 3, p. 837, Sep. 1988, doi: 10.2307/2531595.

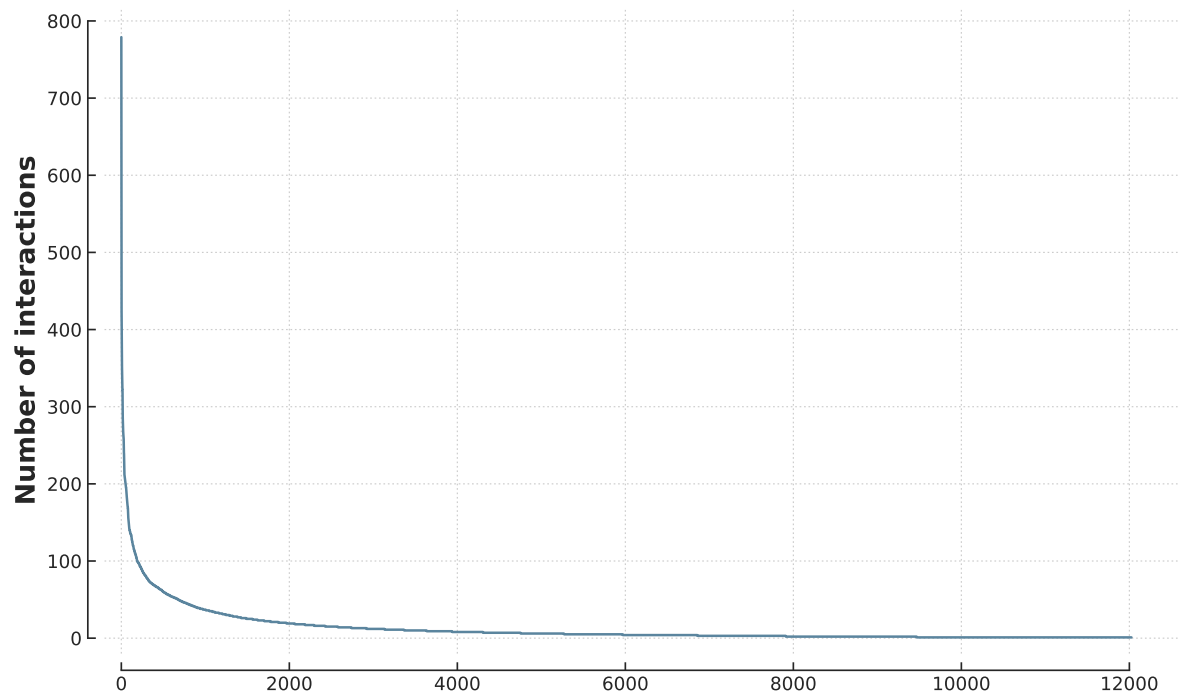
721 [72] J. Neyman and E. S. Pearson, ‘On the Use and Interpretation of Certain Test Criteria for
722 Purposes of Statistical Inference: Part I’, *Biometrika*, vol. 20A, no. 1/2, pp. 175–240, 1928, doi:
723 10.2307/2331945.

724

725

726 **SUPPLEMENTARY MATERIAL**

727 **Supplementary Figures**



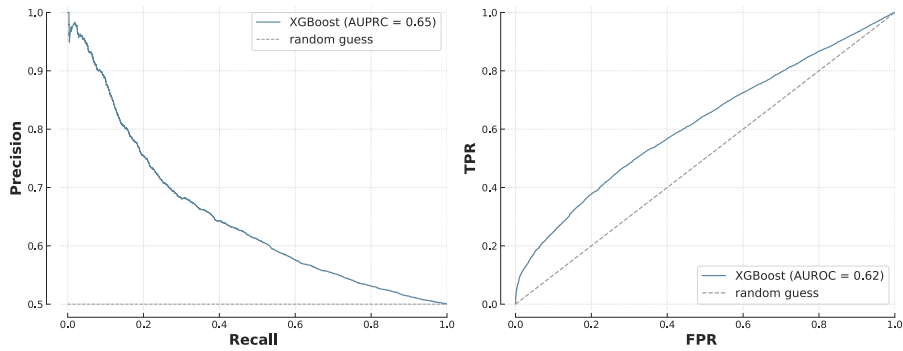
728
729 **Supplementary Figure 1: Distribution of proteins' degree in IntAct. The exponential decrease is characteristic of a**
730 **scale-free network.**

731

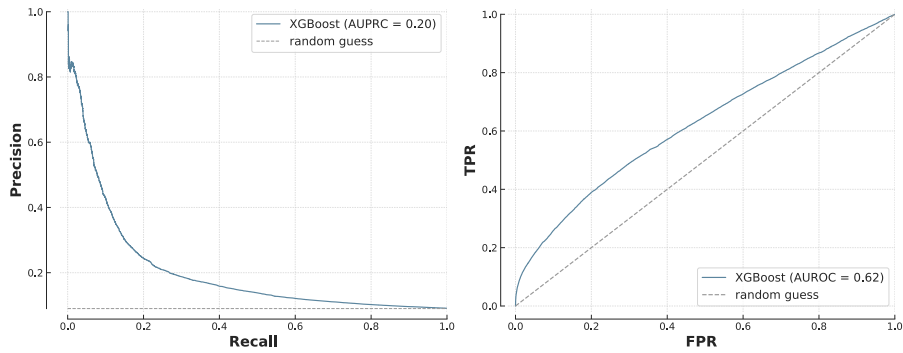
Reporting sheet B4PPI-Human: XGBoost

PR and ROC curves on T1* and T2**

Predictions on T1 (n=24,898, 50% +)



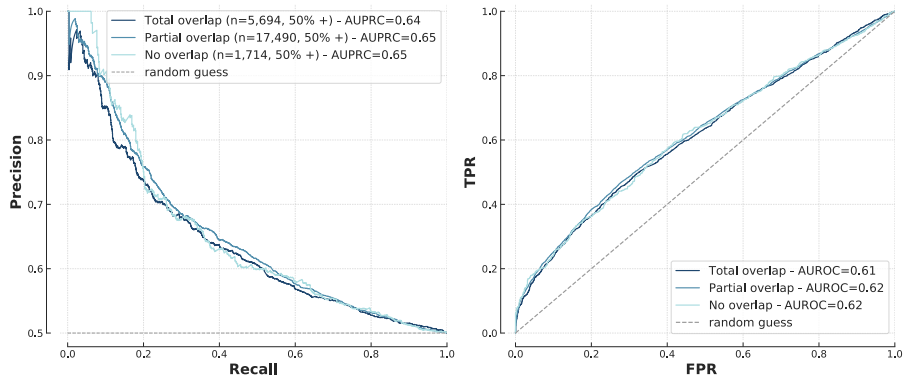
Predictions on T2 (n=136,939, 9% +)



* First testing set, used to compare models on an independent set and investigate protein-level overlap.

** Second testing set, used to assess generalisation on an imbalanced dataset (10 times more negative examples than positive ones).

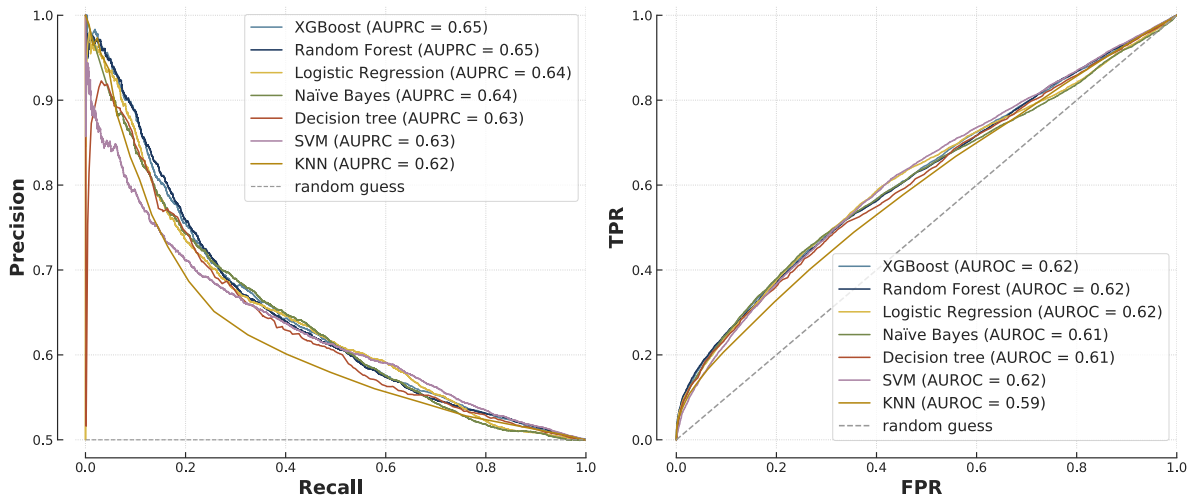
Impact of protein-level overlap



	Running time	Memory	Energy used	Carbon footprint (UK)
Training once	30s	Negligible	< 0.01 kWh	0.002 gCO ₂ e
Training incl. hyperparameters tuning	22min	Negligible	< 0.01 kWh	1 gCO ₂ e
Inference	<1s	Negligible	~0	~0

732

733 Supplementary Figure 2: Performance sheet of XGBoost on B4PPI-Human.



734

735 Supplementary Figure 3: Comparison of a broader range of FG-based models.

736

```

=====
              coef      std err          z      P>|z|      [0.025   0.975]
-----
RNAseqHPA          0.0196      0.008      2.526      0.012      0.004   0.035
tissueHPA         -0.0380      0.024     -1.576      0.115     -0.085   0.009
tissueCellHPA     0.0617      0.024      2.614      0.009      0.015   0.108
subcellularLocationHPA 0.0144      0.006      2.235      0.025      0.002   0.027
bioProcessUniprot 0.3038      0.011     27.377      0.000      0.282   0.326
cellCompUniprot   0.2650      0.007     37.125      0.000      0.251   0.279
molFuncUniprot    0.0329      0.008      4.353      0.000      0.018   0.048
domainUniprot     0.1696      0.012     14.470      0.000      0.147   0.193
motifUniprot      0.0235      0.007      3.249      0.001      0.009   0.038
Bgee              0.0196      0.007      2.735      0.006      0.006   0.034
=====

```

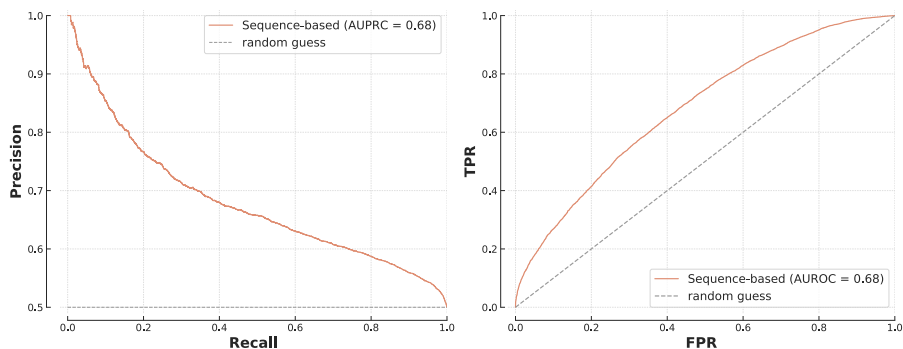
737

738 Supplementary Figure 4: Output of the logistic regression on the training set.

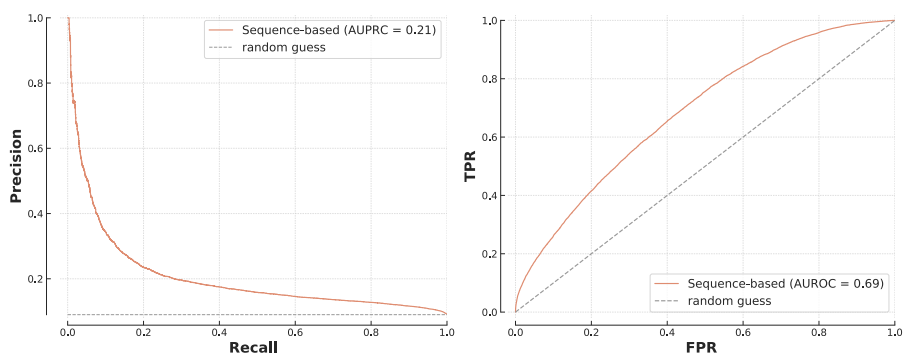
Reporting sheet B4PPI-Human: Sequence-based

PR and ROC curves on T1* and T2**

Predictions on T1 (n=24,898, 50% +)



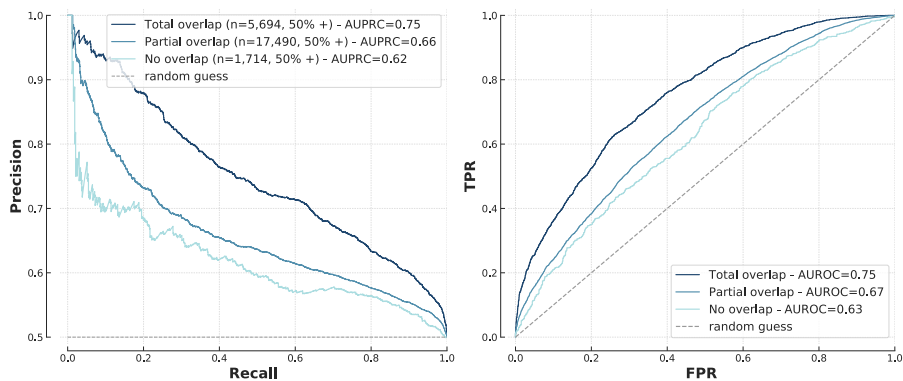
Predictions on T2 (n=136,939, 9% +)



* First testing set, used to compare models on an independent set and investigate protein-level overlap.

** Second testing set, used to assess generalisation on an imbalanced dataset (10 times more negative examples than positive ones).

Impact of protein-level overlap

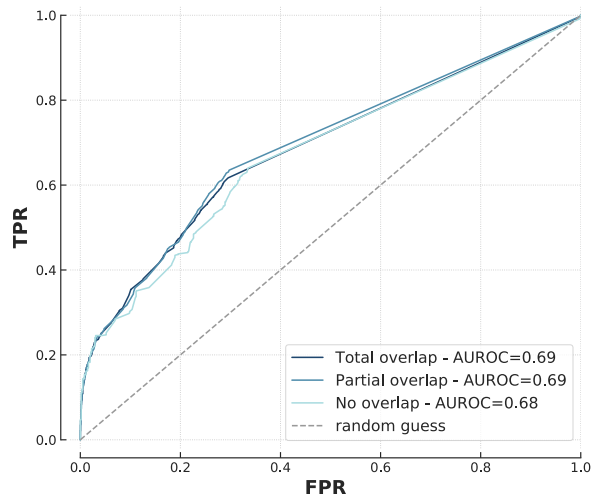
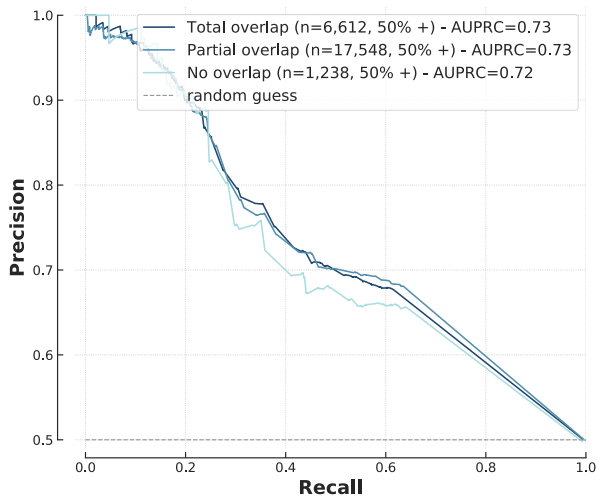


	Running time	Memory	Energy used	Carbon footprint (UK)
Training once	1h10	15 GB	0.62 kWh	158 gCO ₂ e
Training incl. hyperparameters tuning	>100h	>1.5 TB	> 62 kWh	> 15 kgCO ₂ e
Inference	1min46	6 GB	0.01 kWh	4 gCO ₂ e

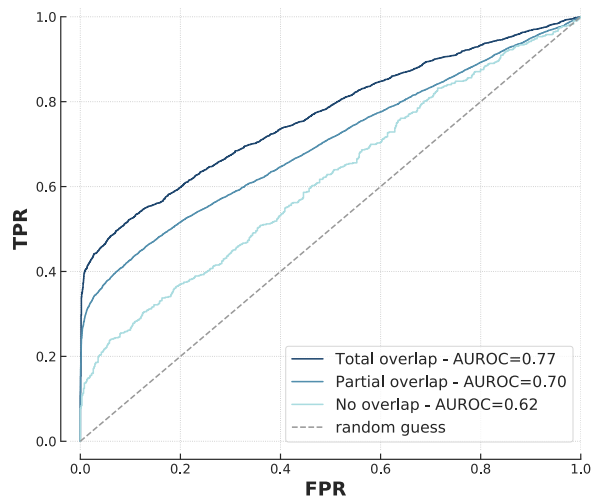
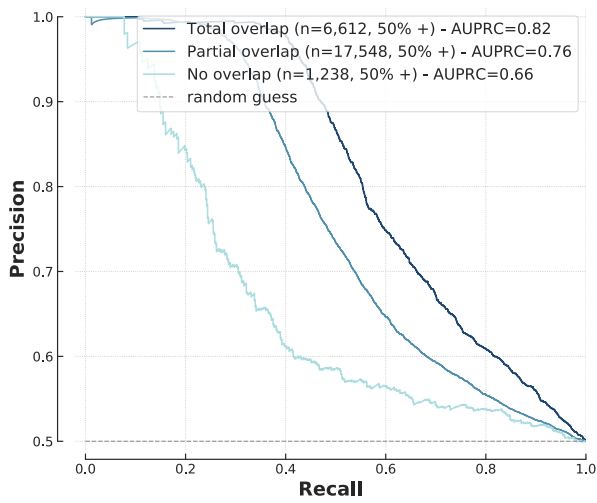
Number of (trainable) parameters: 1.6m

739

740 Supplementary Figure 5: Performance sheet of the sequence-based model.

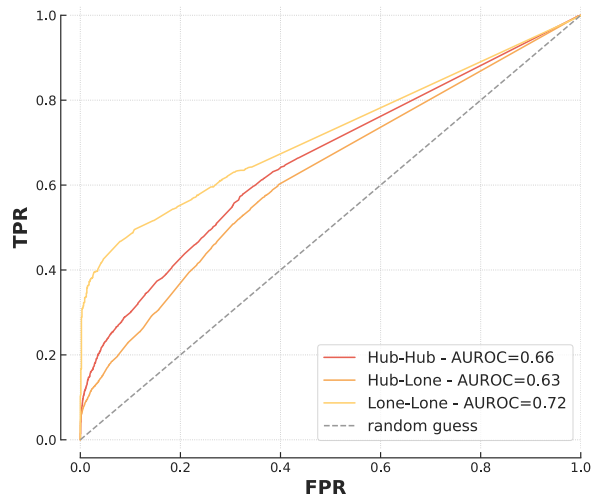
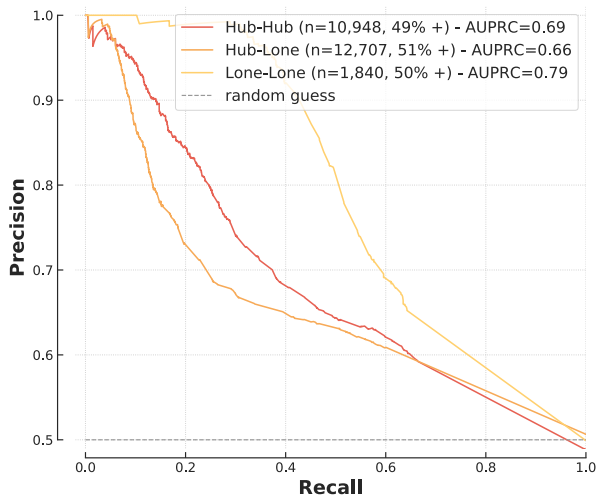


741

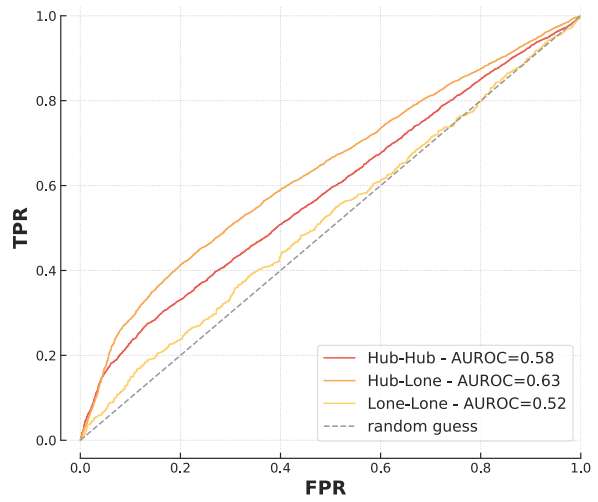
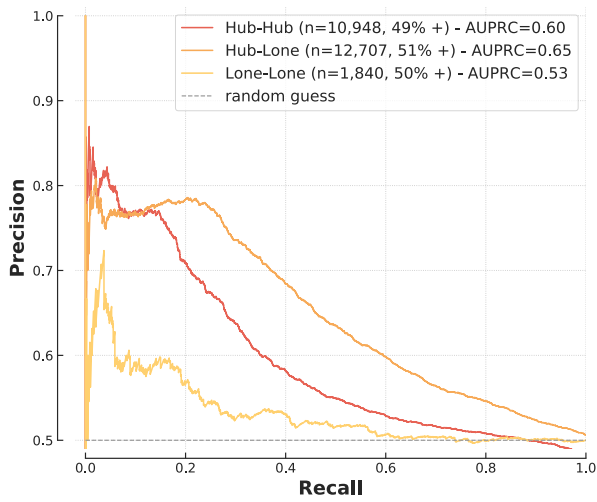


742

743 Supplementary Figure 6: Impact of protein-level overlap on the yeast dataset for XGBoost (top, FG-based) and the
 744 sequence-based model (bottom).



745



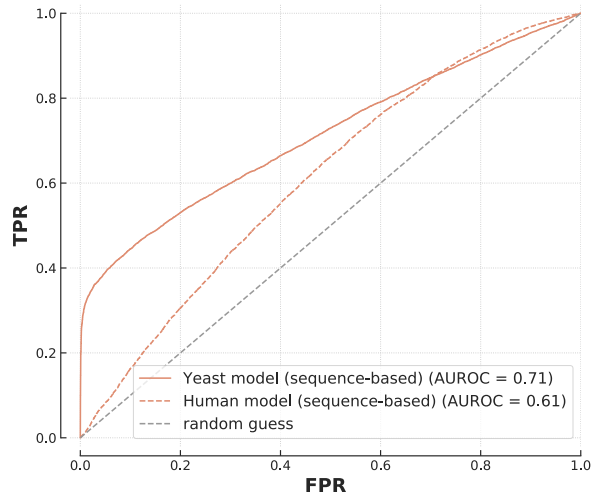
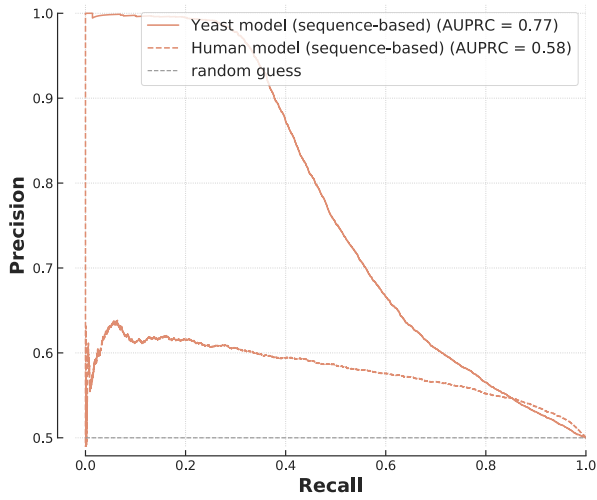
746

747

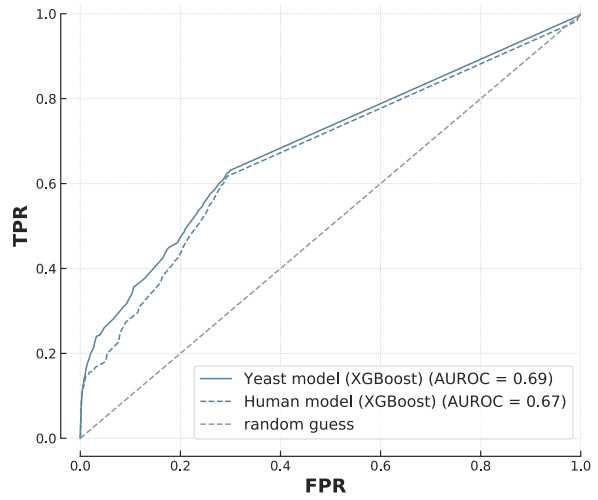
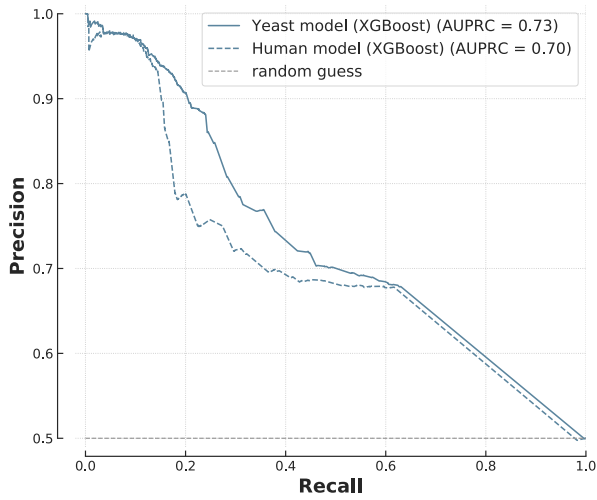
748

Supplementary Figure 7: Impact of hubs for FG-based (XGBoost, top) and sequence-based (bottom) model on yeast interactions.

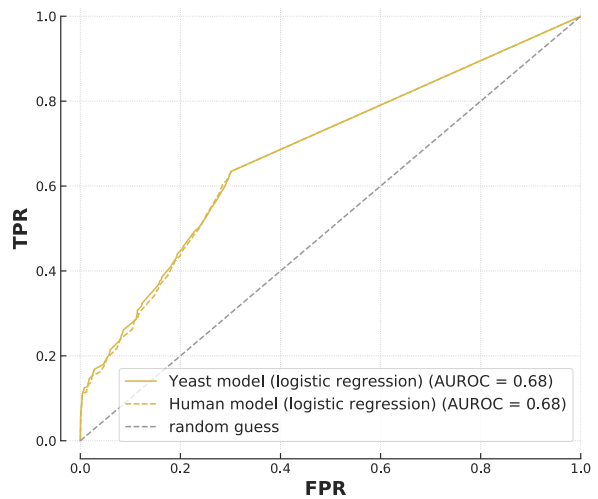
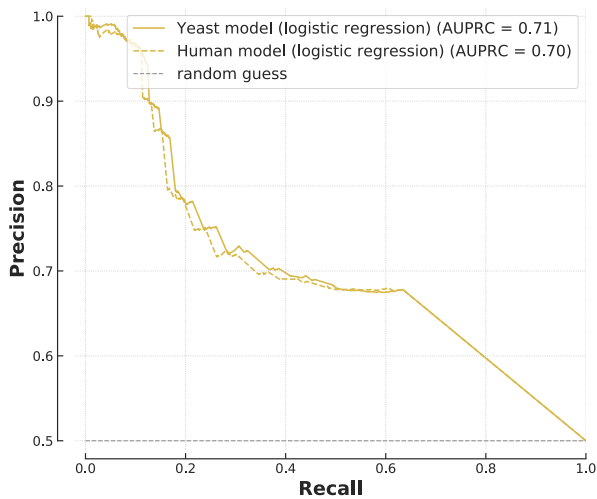
749



750



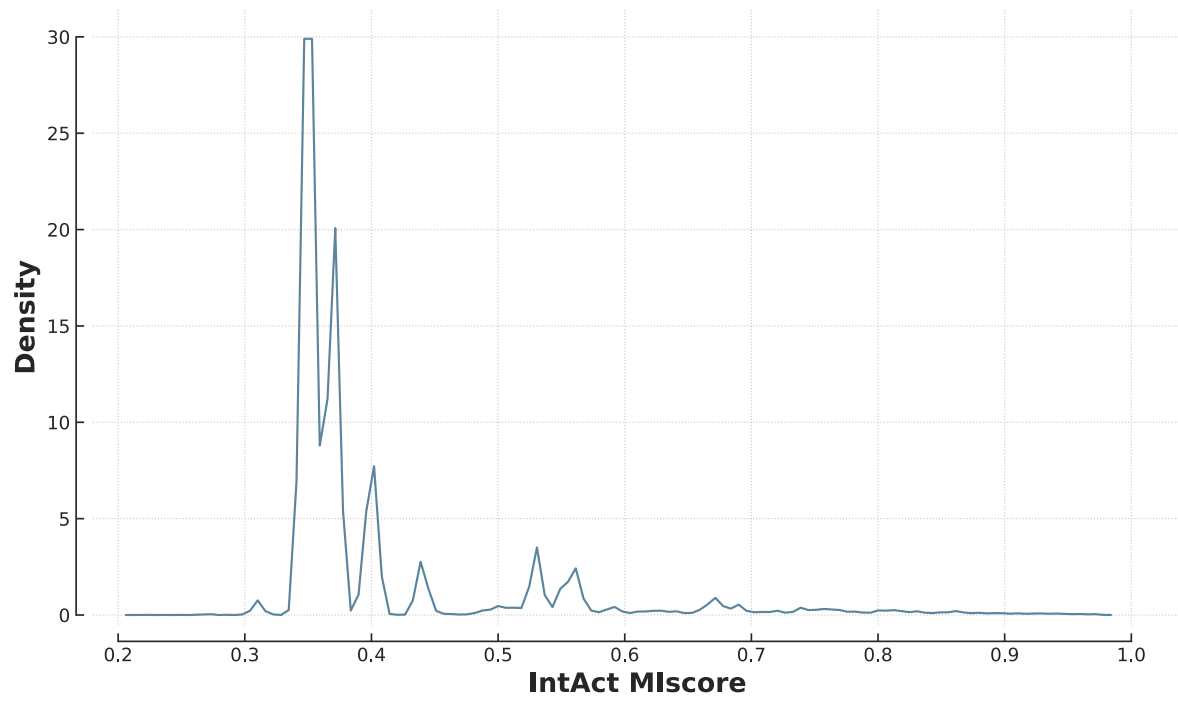
751



752

753 **Supplementary Figure 8: Cross-species predictions.** Models trained on human PPIs (dotted lines) and yeast PPIs
 754 (solid lines) were used to make predictions on the yeast testing set (n=25,398, 50% positive). The top plot is the
 755 sequence-based method, the others the FG-based ones (XGBoost in the middle, logistic regression at the bottom).

756



757

758 Supplementary Figure 9: distribution of IntAct's MIscore in the yeast dataset.

759

760 **Supplementary Tables**

761 **Supplementary Table 1: Sample size in B4PPI.**

Set	Number of examples (% of positive)	
	B4PPI-Human	Yeast dataset
Training	106,662 (50%)	60,738 (50%)
T1	24,898 (50%)	25,398 (50%)
T2	136,939 (9%)	N/A

762

763 **Supplementary Table 2: Sample size in each category in the dataset used to investigate networks topology.**

Type of interaction	Number of pairs	% of PPIs
Hub-hub	27,580	50.69 %
Hub-lone	19,205	49.70 %
Lone-lone	3,011	45.67 %

764

765 **Supplementary Table 3: Details of the features used for B4PPI-Yeast.**

Feature	Number of different annotations	Missing values (/6,721)	Source
Biological processes (GO)	3,114	1,510	UniProt [59]
Cellular components (GO)	820	839	UniProt
Molecular functions (GO)	2,079	2,242	UniProt
Domains	606	5,135	UniProt
Motifs	181	6,273	UniProt
Sequence	N/A	0	UniProt

766

767

768 **Supplementary Table 4: Optimal parameters for the models trained on B4PPI-Human**

Algorithm	Missing data imputation	Scaling	Optimal hyperparameters
Logistic Regression	Yes (mean)	Yes	Penalty = none, tol = 0.0001
XGBoost	No	No	colsample_bytree = 0.8059, learning_rate = 0.00002186, max_depth = 29, min_child_weight = 25, n_estimators = 116, subsample = 0.4595
Decision Tree	Yes (mean)	No	Criterion = entropy, min_samples_split = 895, splitter = random
SVM	Yes (mean)	No	C = 1, degree = 3, gamma = scale, kernel = rbf (default values were used due to long runtime)
Random Forest	Yes (mean)	No	Criterion = gini, max_features = log2, min_sample_split = 487, n_estimators = 336
KNN	Yes (mean)	No	Algorithm = brute, leaf_size = 53, n_neighbors = 35, p = 2, weights = uniform
Naïve Bayes	Yes (mean)	No	N/A
Sequence-based Siamese architecture	N/A	N/A	Batch size = 200, gradient_clip_val = 10, RNN = bidirectional GRU, output = linear, hidden size = 512, n_layers = 1, learning rate = 0.001 (GRU) and 0.0001 (output)

769

770 **Supplementary Table 5: Optimal parameters for the models trained on the yeast dataset.**

Algorithm	Missing data imputation	Scaling	Optimal hyperparameters
Logistic Regression	Yes (mean)	Yes	Penalty = none, tol = 0.0001
XGBoost	No	No	colsample_bytree = 0.7087, learning_rate = 0.00001129, max_depth = 26, min_child_weight = 3, n_estimators = 244, subsample = 0.9318
Naïve Bayes	Yes (mean)	No	N/A
Sequence-based Siamese architecture	N/A	N/A	Batch size = 200, gradient_clip_val = 10, RNN = bidirectional GRU, output = linear, hidden size = 512, n_layers = 1, learning rate = 0.001 (GRU) and 0.0001 (output)

771

772

773