

Large-Scale Protein-Protein Post-Translational Modification Extraction With Distant Supervision And Confidence Calibrated BioBERT

Aparna Elangovan

The University of Melbourne

Yuan Li

The University of Melbourne

Douglas E.V. Pires

The University of Melbourne

Melissa J. Davis

Walter and Eliza Hall Institute of Medical Research

Karin Verspoor (✉ karin.verspoor@rmit.edu.au)

RMIT University

Research Article

Keywords: Protein-protein interaction, Post-translational modifications, BioBERT, Natural language processing, Deep learning, Distant supervision

Posted Date: September 22nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-898489/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at BMC Bioinformatics on January 4th, 2022. See the published version at <https://doi.org/10.1186/s12859-021-04504-x>.

RESEARCH

Large-scale protein-protein post-translational modification extraction with distant supervision and confidence calibrated BioBERT

Aparna Elangovan¹, Yuan Li¹, Douglas E. V. Pires¹, Melissa J. Davis^{2,3} and Karin Verspoor^{1,4*}

*Correspondence:

karin.verspoor@rmit.edu.au

⁴School of Computing

Technologies RMIT University,

Melbourne, Australia

Full list of author information is available at the end of the article

Abstract

Motivation: Protein-protein interactions (PPIs) are critical to normal cellular function and are related to many disease pathways. A range of protein functions are mediated and regulated by protein interactions through post-translational modifications (PTM). However, only 4% of PPIs are annotated with PTMs in biological knowledge databases such as IntAct, mainly performed through manual curation, which is neither time- nor cost-effective. Here we aim to facilitate annotation by extracting PPIs along with their pairwise PTM from the literature by using distantly supervised training data using deep learning to aid human curation. We further assessed model generalisation in a real-world scenario, evaluating its performance on a randomly sampled subset of predictions from 18 million PubMed abstracts.

Method: We use the IntAct PPI database to create a distant supervised dataset annotated with interacting protein pairs, their corresponding PTM type, and associated abstracts from the PubMed database. We train an ensemble of BioBERT models – dubbed PPI-BioBERT-x10 – to improve confidence calibration. We extend the use of ensemble average confidence approach with confidence variation to counteract the effects of class imbalance to extract high confidence predictions.

Results and conclusion: The PPI-BioBERT-x10 model evaluated on the test set resulted in a modest F1-micro 41.3 (P=58.1, R=32.1). However, by combining high confidence and low variation to identify high quality predictions, tuning the predictions for precision, we retained 19% of the test predictions with 100% precision. We evaluated PPI-BioBERT-x10 on 18 million PubMed abstracts and extracted 1.6 million (546507 unique PTM-PPI triplets) PTM-PPI predictions, and filter \approx 5,700 (4584 unique) high confidence predictions. Of the 5700, human evaluation on a small randomly sampled subset shows that the precision drops to 33.7% despite confidence calibration and highlights the challenges of generalisability beyond the test set even with confidence calibration. We circumvent the problem by only including predictions associated with multiple papers, improving the precision to 58.8%. In this work, we highlight the benefits and challenges of deep learning-based text mining in practice, and the need for increased emphasis on confidence calibration to facilitate human curation efforts.

Keywords: Protein-protein interaction; Post-translational modifications; BioBERT; Natural language processing; Deep learning; Distant supervision

1 **Background**

2 Critical biological processes, such as signaling cascades and metabolism, are reg-
3 ulated by protein-protein interactions (PPIs) that modify other proteins in order
4 to modulate their stability or activity via post-translational modifications (PTMs).
5 PPIs are curated in large online repositories such as IntAct [1] and HPRD [2].
6 However, most PPIs are not annotated with a function, for example, we found the
7 IntAct database has over 100,000 human PPIs, but less than 4000 of these are
8 annotated with PTMs such as phosphorylation, acetylation or methylation. Un-
9 derstanding the nature of PTM between an interacting protein pair is critical for
10 researchers to determine the impact of network perturbations and downstream bio-
11 logical consequences. PPIs and PTMs in biological databases are usually manually
12 curated, which is time consuming and requires highly trained curators. [3]. Orchard
13 *et al.* [4] have also highlighted additional challenges in maintaining manually cu-
14 rated databases, ensuring they are up to date, as well as the economic aspects of
15 manual curation. Hence, the adoption of automated curation methods is essential
16 for sustainability of this work.

17 Here we extract PTMs by text mining PubMed abstracts, extracting protein pairs
18 along with their corresponding PTM. Given an input journal abstract, the output
19 is a triplet of the form [Protein1, PTM function, Protein2] where Protein1 and Pro-
20 tein2 are Uniprot identifiers [5] of the proteins. We also aim to aid human curation
21 of PTM-PPIs, hence we assess how well machine learning models generalise by ap-
22 plying them to 18 million PubMed abstracts to extract PTM-PPI triplets. In this
23 paper, we use confidence calibration [6, 7] as a mechanism to understand general-
24 isability to know when a prediction works to extract high quality predictions. We
25 believe our paper is the first to study the practical applicability and challenges of
26 large scale PTM-PPI extraction using NLP with deep learning and distant super-
27 vision.

28 We focus on extracting PTMs including phosphorylation, dephosphorylation,
29 methylation, ubiquitination, deubiquitination, and acetylation (these PTMs were
30 selected based on the availability of training data). We create a training dataset
31 using distant supervised approach [8, 9] using IntAct [1] as the source knowledge
32 base to extract PTM-PPI triplets from PubMed abstracts. We train an ensemble
33 of BioBERT [10] models to improve neural confidence calibration [6, 7]. We then
34 applied the trained model on 18 million PubMed abstracts to extract PPI pairs
35 along with their corresponding PTM function and attempt to ensure high quality
36 predictions using neural confidence calibration techniques [6, 7] to augment and
37 facilitate human curation efforts.

38 Related works in protein interaction extraction through deep learning

39 The datasets in PPI extraction such as AIMed [11] and BioInfer [12] used to evalu-
40 ate text mining approaches have remained the same for over a decade (since 2007)
41 and focus on extracting protein interactions but not the nature of the PTM in-
42 teraction between them. These datasets have also been used to evaluate the latest
43 machine learning approaches including deep learning [13, 14] in protein pair extrac-
44 tion. However, the latest deep learning trends do not seem to be widely popular
45 in PPI curation outside the limited context of benchmarking methods using the
46 AIMed [11] and BioInfer [12] datasets. Automated PPI curation attempts using
47 text mining and rule-based approaches seem more prevalent [15, 16, 17].

48 STRING v11 [18], one of the most popular PPI databases, uses text mining as a
49 curation method. Their text mining pipeline has largely remained the same since
50 STRING v9.1 [15]. STRING v9.1 uses a weighted PPI co-occurrence rule-based
51 approach, where the weights depend on whether a protein pair occurs together
52 within the same document, the same paragraph and/or the same sentence. Rule
53 based approaches can be quite effective [19, 15] even with limited training data,
54 depending on the task. For instance, Szklarczyk *et al.* define an interaction unit

55 in STRING v11 database [18] as a “*functional association, i.e. a link between two*
56 *proteins that both contribute jointly to a specific biological function*”. This definition
57 allows for co-occurrence rule-based approach to be quite effective, *i.e.* if a protein
58 pair co-occurs in text frequently then the pair is highly likely to be related.

59 iPTMnet [20] consolidates information about PPIs and PTMs from various man-
60 ually curated databases such as HPRD [2] and PhosphoSitePlus [21] as well as
61 text mining sources. For text mining, iPTMnet uses RLIMS-P [17] and eFIP [16]
62 to automatically curate enzyme-substrate-site relationships. These tools use rule-
63 based approaches using text patterns to extract proteins involved in PTM. The
64 iPTMnet statistics dated Nov 2019 ([https://research.bioinformatics.udel.](https://research.bioinformatics.udel.edu/iptmnet/stat)
65 [edu/iptmnet/stat](https://research.bioinformatics.udel.edu/iptmnet/stat)) indicate that the total number of enzyme-substrate pairs cu-
66 rated using RLIMS-P [17] is fewer than 1,000 pairs. This modest number highlights
67 the main challenge using text patterns: while they can extract relationships with
68 fairly high precision, they are not robust to variations in how PPI relations can be
69 described in text. We therefore explore machine learning-based methods.

70 Automatic extraction of PPIs using deep learning can be beneficial as it has the
71 potential to extract PPIs from a variety of text where the PPI relationships are
72 described in ways that cannot be easily captured by a manually crafted rule-based
73 system. However, deep learning requires a much larger volume of training data.
74 Generalisability of the model to ensure prediction quality is key to its widespread
75 adoption for automatic extraction of PPI relationships from text at scale. Enhanc-
76 ing quality of prediction at large scale needs to focus on reducing false positives to
77 minimise corrupting existing knowledge base entries (*i.e.* predicting that a relation-
78 ship exists when it doesn't) and, hence, confidence calibration approaches to reduce
79 poor quality predictions become a crucial step in large scale text mining. Confidence
80 calibration is the problem of predicting probability estimates representative of the
81 true correctness [7] and in this paper we use confidence calibration to know when

82 a prediction is likely correct and use it as a mechanism to improve generalisation.
83 The aspects of generalisability have been largely limited to the evaluation of a test
84 set and the limitations of using the test set performance as a proxy for real world
85 performance have been challenged in previous studies [22, 23].

86 Creating gold standard training data with fine-grained annotations is a manual,
87 labor-intensive task and is a limiting factor in applying machine learning to new
88 domains or tasks. Being able to leverage one or more existing data sources is key to
89 using machine learning in new domains or for new tasks. Distant supervision [8, 9]
90 exploits existing knowledge bases, such as IntAct [1], instead of annotating a new
91 dataset. However there are two main limitations to the use of distant supervision
92 datasets: **(a)** noisy labels require noise reduction techniques to improve label quality
93 [24] **(b)** they require negative samples to be generated as the databases usually only
94 contain positive examples of a relationship.

95 Deep learning architectures such as BiLSTM and BioBERT have been previously
96 used to benchmark methods for protein relation extraction using Natural Language
97 Processing (NLP) and the AIMed dataset [13, 25, 22]. However, these works [13,
98 25, 22] do not measure the ability of these models to calibrate confidence scores. We
99 chose a state-of-the-art deep learning approach, BioBERT [10], train an ensemble
100 to enhance confidence calibration [6] and use confidence variation to counteract the
101 effects of class imbalance during confidence calibration.

102 **Methods**

103 **Dataset**

104 We obtain the dataset of human PTM-PPI interactions from the IntAct database
105 [1], a database that is part of the International Molecular Exchange (IMEX) [26]
106 Consortium. The database contains Uniprot identifiers, protein aliases, interaction
107 type (where available), and the PubMed identifier of the paper describing the inter-
108 action. Of the over 100,000 human PPIs in IntAct [1], only 3,381 PPIs describe PTM

109 interactions. Hence, we start with an initial dataset containing these 3,381 PPIs. We
 110 then remove duplicate entries, *i.e.* the interacting proteins, the PubMedId and the
 111 interaction type are the same resulting in 2797 samples. We then pre-process this
 112 data to remove self-relations (where participant protein1 is the same as protein2).

Table 1 Illustration of data preparation. The IntAct database has the PubMedId, the Uniprot identifiers of the participating proteins and the interaction type. More than 1 pair of interacting proteins can be annotated against a given PubMed Id and not all of these interactions are described in the abstract. This forms the noisily labelled training data.

IntAct	Pubmed id: 24291004 , Uniprot protein pair: P04150, P31749, Interaction type: Phosphorylation
Corresponding abstract	<i>Glucocorticoid resistance ...[truncated display]...</i> we identify the AKT1 kinase as a major negative regulator of the NR3C1 glucocorticoid receptor protein activity driving glucocorticoid resistance in T-ALL. Mechanistically, AKT1 impairs glucocorticoid-induced gene expression by direct phosphorylation of NR3C1 at position S134 and blocking glucocorticoid-induced NR3C1 translocation to the nucleus. Moreover, we demonstrate that loss of PTEN and consequent AKT1 activation can effectively block glucocorticoid-induced apoptosis and induce resistance to glucocorticoid therapy. Conversely, pharmacologic inhibition of AKT with MK2206 effectively restores glucocorticoid-induced NR3C1 translocation <i>in vitro and in vivo.</i>
Genes and normalised NCBI ids	<ul style="list-style-type: none"> ● Start-End, Gene Mention, NCBI gene Id ● 137 - 141, AKT1, 207 ● 186 - 191, NR3C1, 2908 ● 294 - 298, AKT1, 207 ● 375 - 380, NR3C1, 2908 ● 434 - 439, NR3C1, 2908 ● 508 - 512, PTEN, 5728 ● 528 - 532, AKT1, 207 ● 748 - 753, NR3C1, 2908
NCBI to Uniprot map	<ul style="list-style-type: none"> ● NCBI gene: Uniprot identifiers ● 2908: P04150, E5KQF5, E5KQF6, F1D8N4, B7Z7I2 ● 207: P31749, B0LPE5, B3KVVH4 ● 5728: P60484, F6KD01
Normalised abstract for pair P04150, P31749	<i>Glucocorticoid resistance ...[truncated display]...</i> we identify the P31749 kinase as a major negative regulator of the P04150 glucocorticoid receptor protein activity driving glucocorticoid resistance in T-ALL. Mechanistically, P31749 impairs glucocorticoid-induced gene expression by direct phosphorylation of P04150 at position S134 and blocking glucocorticoid-induced P04150 translocation to the nucleus. Moreover, we demonstrate that loss of P60484 and consequent P31749 activation can effectively block glucocorticoid-induced apoptosis and induce resistance to glucocorticoid therapy. Conversely, pharmacologic inhibition of AKT with MK2206 effectively restores glucocorticoid-induced P04150 translocation <i>in vitro and in vivo.</i>
Negative sample pairs	<ul style="list-style-type: none"> ● P04150, P60484 ● P31749, P60484

113 Once we pre-process the data, the steps involved in dataset transformation is
 114 illustrated in Table 1. We first identify all the gene mentions in an abstract and their
 115 corresponding Uniprot protein identifier [5], as detailed in section *Gene mentions*
 116 *and protein identifiers*. We apply noise reduction techniques to remove samples
 117 where the abstract may not describe the PPI relationship, as detailed in section
 118 *Noise reduction*. We then split the dataset into train, test and validation sets such
 119 that they are stratified by interaction type and have unique PubMed ids in each
 120 set to avoid test set leakage resulting from random splits identified in the AIMed
 121 dataset [22]. Negative training samples are generated as detailed in section *Negative*
 122 *sample generation*.

123 Gene mentions and protein identifiers

124 To prepare the training data, the Uniprot identifiers mentioned in IntAct need to
 125 match the Uniprot identifiers that can be inferred from the abstract text by ap-
 126 plying named entity recognition and normalisation NLP techniques, *gene mention*
 127 \rightarrow *Uniprot accession number*. GNormPlus[27] identifies gene names and normalises
 128 them to the NCBI gene identifier [28]. The NCBI gene identifier is a gene-level
 129 identifier that needs to be converted to a Uniprot accession code which is a protein
 130 identifier. A single NCBI gene can be associated with multiple Uniprot accession
 131 numbers either due to biological reasons such as alternative splicing where a sin-
 132 gle gene results in different protein isoforms or due to technical reasons related to
 133 Uniprot accession number changes [5]. For a given abstract A , and Uniprot iden-
 134 tifiers of the two PPI participants IU_1, IU_2 annotated in the IntAct database as
 135 $\langle A, IU_1, IU_2 \rangle$, we perform the following:

- 136 1 Use GNormPlus to identify gene name mentions $[g_1, g_2..g_m]$ in abstract
 137 A , and to normalise gene names to corresponding NCBI gene ids, $NG =$
 138 $[ng_1, ng_2, \dots, ng_m]$.
- 139 2 Given a NCBI gene id, $ng_i \in NG$, we obtain a list of corresponding Uniprot
 140 accession numbers $U = [u_1, u_2, \dots, u_n]$.
- 141 3 If IU_1 or IU_2 exists in U , then we have found a matching IntAct Uniprot
 142 identifier IU_x annotated against the abstract, where $IU_x \in [IU_1, IU_2]$. Map
 143 NCBI gene id ng_i to Uniprot identifier IU_x . If neither IU_1 nor IU_2 exist in
 144 U , *i.e.* the Uniprot identified in the abstract is not annotated against the
 145 abstract, then simply return the first Uniprot identifier u_1 .

146 This gene mention identification to Uniprot identifier mapping process is illus-
 147 trated through an example in Table 1.

148 Noise reduction

149 In order to extract PTM-PPI relationships mentioned in the abstract, we aim to
150 retain only those training samples where the PTM-PPI is described in the abstract.
151 Hence we exclude IntAct PTM-PPI triplets from the training data if either par-
152 ticipant’s Uniprot identifiers does not exist in the normalized abstract. This is to
153 minimize the false positive noise in the distant supervised dataset, *i.e.* removing
154 samples where the relationship is not described in the abstract, based on the as-
155 sumption that, if the proteins are not explicitly mentioned in the abstract then it
156 is highly likely (unless protein names are not recognized by the NER tool) that
157 the abstract does not describe the relationship between those proteins. We also re-
158 move PPIs where the stemmed interaction type (phosphoryl, dephosphoryl, methyl,
159 acetyl, deubiquitin, ubiquitin, demethyl) is not mentioned in the abstract.

160 Negative sample generation

161 Since knowledge bases focus on relationships that exist, distant supervision-based
162 approaches require negative training examples to be created. In our dataset, nega-
163 tive samples are protein pairs that are mentioned in the abstract but do not have
164 a function referencing that paper annotated against the pair in Intact. In order
165 to generate negative samples, we identify protein mentions from the abstracts us-
166 ing GNormPlus [27], which normalizes mentions to NCBI gene IDs which are then
167 converted to Uniprot identifiers as described in section *Gene mentions and protein*
168 *identifiers*. If a given protein pair $\langle p1, p2 \rangle$, where $p1$ and $p2$ are mentioned in the
169 abstract, but is not annotated against any of the 7 types of PPI relationship within
170 the abstract according to the Intact database, then it is assumed that $\langle p1, p2 \rangle$
171 form negative samples for that abstract, see example in Table 1.

172 It is important to emphasize that a negative sample does not mean that a given
173 PPI relationship does not exist, but rather that the abstract does not describe such a

174 relationship. It could also be a noisy negative sample, *i.e.* the abstract describes the
175 functional relationship between pair, but it is simply not captured in the annotation.

176 Training BioBERT for PPI extraction

177 We fine-tuned the pretrained BioBERT v1.1 [10], based on BERT-base-Cased, which
178 is pretrained on a large collection of PubMed abstracts. We applied fine-tuning to
179 adapt BioBERT to the PTM-PPI extraction task, as a multi-class classification
180 problem, and assumed there is at most one type of PTM relationship between a
181 protein pair. For full details on training configuration, see supplementary section A.1

182 The original BioBERT is trained using a sentence as a training sample for the
183 language modelling task, while we work with complete abstracts. To utilize an entire
184 abstract as a single training sample, we feed all the sentences within the abstract
185 including the full stop (end-of-sentence marker, the period “.”) separating each
186 sentence. BioBERT has a maximum input sequence length of 512 units, this includes
187 the sub word units to fit the entire abstract. Sub word units are concatenated to
188 represent words and limit the vocabulary size, *e.g.*, a word such as Immunoglobulin
189 is tokenized into 7 units (I, ##mm, ##uno, ##g, ##lo, ##bul, ##in) where
190 ## are special symbols to indicate continuation [29]. Therefore, an abstract with
191 250 words can result in a much longer sequence once it is tokenised. In our training
192 data, post tokenisation, 90% of the normalised abstracts are under the maximum
193 limit of 512 sub units. For more detailed distribution see supplementary Table A10.
194 Therefore, we simply truncate sequences longer than 510 (to accommodate the
195 mandatory marker [CLS] and [SEP] tokens that BERT requires [29]) and feed the
196 input abstracts.

197 Uncertainty estimation

In order to improve the probability estimate associated with each prediction, we
use an ensemble of 10 PPI-BioBERT models (referred to as PPI-BioBERT-x10),
all trained with exactly the same hyperparameters and training data but capturing

slightly different models due to Bernoulli (binary) dropout [30] layers in the network. The use of ensembles of models to improve uncertainty estimates, rather than just improving overall accuracy, has been shown to be effective in the computer vision task of image classification Lakshminarayanan *et al.* [6], hence we follow that approach. Using the ensemble of the 10 models, the predicted confidence \hat{p}_j and predicted class \hat{y}_j for the input x_j is:

$$\hat{p}_j = p_j(y = c) = \frac{1}{M} \sum_{i=1}^M p(y_{\theta_i} = c | x_j, \theta_i) \quad (1)$$

$$\hat{y}_j = \operatorname{argmax}_c p(y = c) = \frac{1}{M} \sum_{i=1}^M p(y_{\theta_i} = c | x_j, \theta_i) \quad (2)$$

In addition to averaging confidence, we also capture the standard deviation of the predicted confidence across the ensemble as:

$$\operatorname{std}(\hat{p}_j) = \sqrt{\frac{1}{M} \sum_{i=1}^M (p(y_{\theta_i} = c | x_j, \theta_i) - \hat{p}_j)^2} \quad (3)$$

198 where M is the number of models in the ensemble; θ_i is model i; y_{θ_i} is the output
199 predicted by model i; x is the input; c is the predicted class.

We measure the ability of the ensemble in confidence calibration using Expected Calibration Error (ECE) used previously by Guo *et al.* [7] in computer vision. ECE is defined as:

$$\operatorname{accu}(B_k) = \frac{1}{|B_k|} \sum_{j \in B_k} 1(\hat{y}_j = y_j) \quad (4)$$

$$\operatorname{conf}(B_k) = \frac{1}{|B_k|} \sum_{j \in B_k} \hat{p}_j \quad (5)$$

$$ECE = \sum_{k=1}^K \frac{|B_k|}{n} |\operatorname{accu}(B_k) - \operatorname{conf}(B_k)| \quad (6)$$

200 where the confidence score is divided into K equally spaced bins B_k , n in the number
201 of samples, $\operatorname{accu}(B_k)$ is the accuracy of predictions whose confidence fall into bin

202 B_k and $conf(B_k)$ is the average confidence of predictions that confidence fall into
203 bin B_k . The best calibrated model will have zero ECE.

204 Large scale PPI function extraction from PubMed abstracts

205 We extracted the complete collection of 18 million abstracts available in PubMed as
206 of April 2019. We applied GNormPlus [27] to recognize the proteins and normalize
207 them to UniProt identifiers. We then apply the PPI-BioBERT-x10 model to extract
208 PTM-PPI relationships from the entire PubMed abstracts. In order to extract high
209 quality PTM-PPI relationships, we only select predictions that have low variation
210 in the confidence score and high prediction confidence.

211 Results

212 IntAct-based dataset

213 The final dataset includes the Uniprot identifiers of the participants, the normal-
214 ized Abstract (the protein names normalised to Uniprot identifiers), all the proteins
215 identified through NER in the abstract and the class label (the type of PTM inter-
216 action). After pre-processing and applying noise reduction to the original dataset
217 with approximately over 3,000 interactions in IntAct, we are left with a total of
218 279 positive and 1579 negative samples across training, test and validation set. The
219 distribution of interaction types and positive/negative samples is shown in Table 2,
220 with phosphorylation forming almost 75% of the positive samples in the training
221 set. The positive sample rate is between 14% - 17% in the train, test and validation
222 set. Interaction types such as such acetylation, deubiquitination and ubiquitination,
223 have less than 10 positive samples in total across training, test and validation.

224 Confidence calibration on the test and validation sets

225 The performance on the test and validation sets using the ensemble BioBERT model
226 (PPI-BioBERT-x10) is detailed in Table 3. The test and validation sets have an F1-

Table 2 Train/Test/ Val Positive samples for each interaction type. NOTE The negative samples count against a interaction type is solely to indicate how many were derived from abstracts describing the interaction type. All negative samples belong to a single class "other / Negative".

	Train		Val		Test		Total	
	-	+	-	+	-	+	-	+
acetylation	31	5	2	1	9	1	42	7
dephosphorylation	167	28	36	10	33	6	236	44
deubiquitination	4	2	0	0	0	0	4	2
methylation	22	10	5	1	22	4	49	15
phosphorylation	862	139	118	21	227	44	1207	204
ubiquitination	30	5	5	1	5	1	40	7
Total	1116	189	166	34	296	56	1578	279

227 micro score of 41.3 and 58.6, respectively, across all interactions. The confusion
 228 matrix is shown in Figure 1.

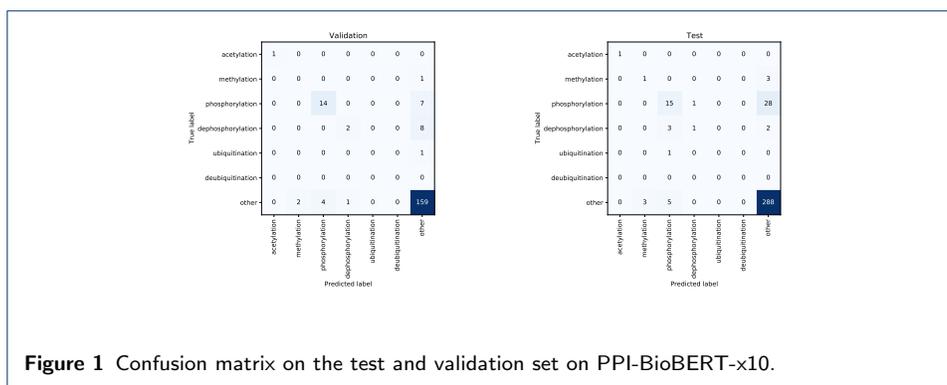


Figure 1 Confusion matrix on the test and validation set on PPI-BioBERT-x10.

229 We visualise the ensemble model’s confidence calibration using reliability dia-
 230 grams, similar to Guo *et al.* [7], as shown in Figure 2. We observe that the pre-
 231 dicted confidence range for each interaction type is proportional to the percentage
 232 of training samples, *i.e.*, the negative samples have the highest proportion and over
 233 80% of the negative sample predictions have a confidence score between 0.9 and 1.0,
 234 whereas acetylation, which has less than 1% of the training samples, has prediction
 235 confidence of less than 0.6. This relationship between the distribution of predicted
 236 confidence and the proportion of training samples for a given class label (interac-
 237 tion type) is consistent in the train, validation and test set as shown in Figure 2.
 238 This highlights 2 main limitations of measuring calibration error using ECE: **(a)** it
 239 penalizes model calibration regardless of the class imbalance in the training dataset

Table 3 The performance of ensemble PPI-BioBERT-x10 on the test and validation set. **ECE** is the expected calibration error. **SD** denotes the average standard deviation within the ensemble.

Dataset	Interaction	P	R	F1	ECE	SD	Support
Test	acetylation	100.00	100.00	100.00	0.49	0.25	1
Test	dephosphorylation	50.00	16.67	25.00	0.67	0.40	6
Test	methylation	25.00	25.00	25.00	0.60	0.28	4
Test	phosphorylation	62.50	34.09	44.12	0.79	0.26	44
Test	ubiquitination	0.00	0.00	0.00	-	-	1
Test	ECE	-	-	-	0.75	-	31
Test	average SD	-	-	-	-	0.28	31
Test	macro avg	47.50	35.15	38.82	-	-	56
Test	micro avg	58.06	32.14	41.38	-	-	56
Val	acetylation	100.00	100.00	100.00	0.53	0.16	1
Val	dephosphorylation	66.67	20.00	30.77	0.61	0.37	10
Val	methylation	0.00	0.00	0.00	0.53	0.29	1
Val	phosphorylation	77.78	66.67	71.79	0.78	0.26	21
Val	ubiquitination	0.00	0.00	0.00	-	-	1
Val	ECE	-	-	-	0.73	-	24
Val	average SD	-	-	-	-	0.28	24
Val	macro avg	48.89	37.33	40.51	-	-	34
Val	micro avg	70.83	50.00	58.62	-	-	34

240 when high quality predictions can be available at lower confidence scores for classes
 241 with lower proportion of samples (**b**) it penalises both over- and under-calibration
 242 equally. Hence, we find it challenging to rely on ECE and/or average ensemble con-
 243 fidence to extract high quality predictions, especially for interaction types that have
 244 a much lower proportion of training samples.

245 We therefore inspected the variation in the confidence score predicted by each
 246 model for a given sample as shown in Figure 3. Intuitively, if for a given input,
 247 all the models within the ensemble consistently predict with similar confidence
 248 then we can rely on the ensemble confidence better compared to high variation
 249 confidence scores. When we compare the variation in confidence scores of the correct
 250 *vs.* incorrect predictions, the incorrect predictions do not seem to have low variation
 251 in predicted confidence, as shown in Figure 3, except for 3 incorrect phosphorylation
 252 predictions. We apply a heuristic to select interaction-wise confidence and standard
 253 deviation thresholds where we select the 50th percentile as the cut-off for both

254 confidence and confidence standard deviation based on the training set. This results
 255 in retaining 19% of positive predictions (6 out of 31) in the test set. Since our
 256 training data is noisy, we manually verify all predictions that have low variation in
 257 the test and validation set. On manual verification, we find all the predictions using
 258 this heuristic are in fact correct, including 3 incorrectly labelled ones, as shown in
 259 Table 4, indicating that a combination of low variation and relative high confidence
 260 has the potential to be robust against noisy labels.

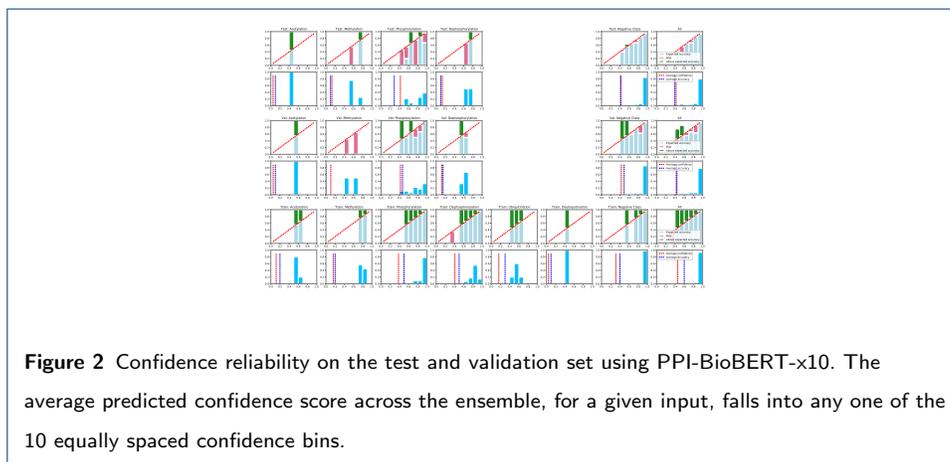
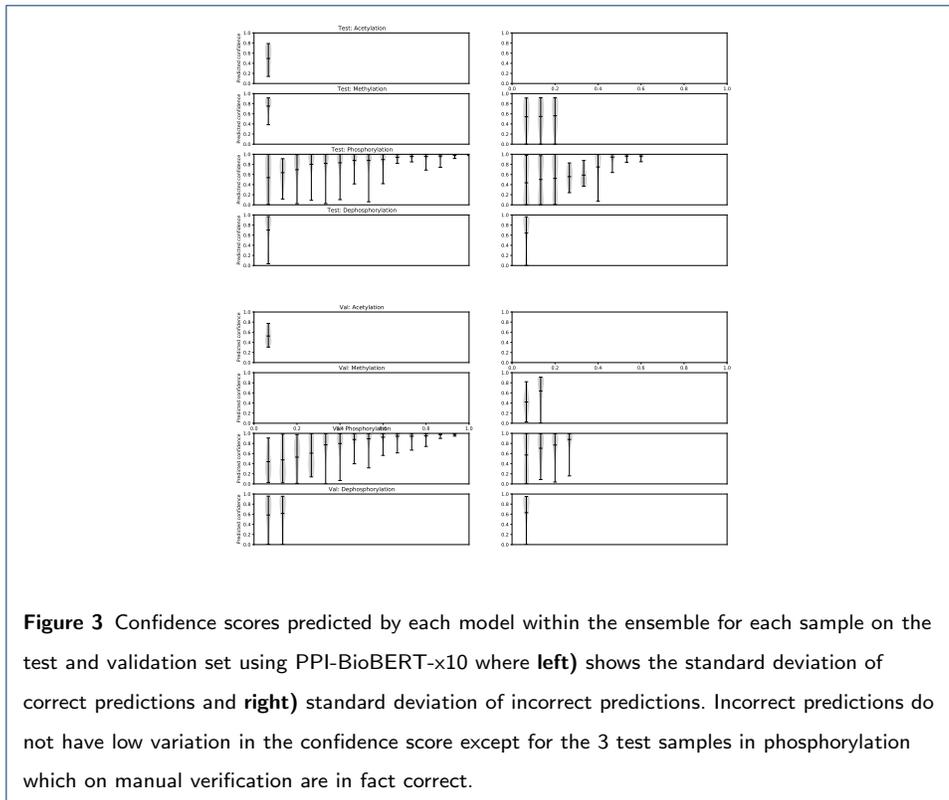


Table 4 Manually verified relationships in the validation and test set with ensemble prediction standard deviation less than interaction-wise threshold and predicted confidence greater than interaction-wise threshold, taking the precision to 100.0

PubMed	Phrases describing relationships in the abstract	Label	Prediction
Validation			
12150926	mTOR (P42345) -catalyzed phosphorylation of 4EBP1 (Q13541) in vitro	phosphorylation	phosphorylation
15733869	SGK (O00141) physically associates with CREB (P16220) and SGK (O00141) phosphorylates it on serine 133	phosphorylation	phosphorylation
15557335	Src (P12931) phosphorylates Alix (Q8WUM4) at a C-terminal	phosphorylation	phosphorylation
15527798	Phosphorylation of MDM2 (Q00987) by the protein kinase AKT (P31749)	phosphorylation	phosphorylation
10864201	radiation-induced phosphorylation of p53 (P04637) protein at serine 15, largely mediated by ATM (Q13315) kinase	phosphorylation	phosphorylation
24548923	Akt (P31749) -mediated phosphorylation of Carma1 (Q9BXL7)	phosphorylation	phosphorylation
19407811	BubR1 (O60566) forms a complex with PCAF (Q92831) and is acetylated at lysine 250	acetylation	acetylation
Test			
21920476	cofilin (P38432) is phosphorylated in Ser184 by both VRK1 (Q99986)	phosphorylation	phosphorylation
19424295	Previous studies of cofilin (P23528) have shown that it is phosphorylated primarily by the LIM domain kinases Limk1 (P53667)	phosphorylation	phosphorylation
22726438	FGFR2 (P21802) phosphorylates tyrosine residues on Grb2 (P62993)	phosphorylation	phosphorylation
11154276	Akt (P31749) decreased ASK1 (Q99683) kinase activity stimulated by both oxidative stress and overexpression in 293 cells by phosphorylating a consensus Akt (P31749) site at serine 83 of ASK1 (Q99683)	phosphorylation	phosphorylation
25605758	phosphorylation of Rab5b (P61020) by LRRK2 (Q55007) also exhibited binding of AKT (P31749) (tail region) to Vim (P08670) (head region)	phosphorylation	phosphorylation
20856200	results in Vim (P08670) Ser39 phosphorylation	phosphorylation	phosphorylation
21986944	MAK (P20794) associates with CDH1 (P12830) (FZR1 (Q9UM11), fizzy/cell division cycle 20 related 1) and phosphorylates CDH1 (P12830)	Negative	phosphorylation
15862297	HDM2 (Q00987) phosphorylation by Chk2 (O96017) doubles in the presence of p53 (P04637)	Negative	phosphorylation
21887822	Hsp70 (P34932) is phosphorylated by Pik1 (P53350)	Negative	phosphorylation



261 Large scale PTM-PPI extraction from PubMed abstracts

262 We extracted approximately 1.6 million PPIs (546507 unique PTM-PPI triplets)
 263 from 18 million abstracts using the PPI-BioBERT-x10 model, see Table 5. The
 264 PTM-wise prediction confidence range during large scale extraction is shown in
 265 Figure 4, and confidence range is similar to training confidence ranges, where
 266 PTMs with lower training samples have lower confidence ranges. After applying
 267 interaction-wise probability and confidence standard deviation thresholds, we re-
 268 tain 5708 PPIs across 5 interaction types. From the 5708 predictions (4584 unique
 269 PTM-PPI triplets), we randomly select 30 predictions per PTM type for human
 270 verification. Despite the high precision (100.0%; see Table 4) on the test set after
 271 applying thresholds we find that during the large scale prediction the precision on
 272 the randomly sampled subset (subset in supplementary Table A14) drops to 33.7 as
 273 shown in Table 6. With the estimated precision of 33.7 %, of the 5708 \approx 1900 have
 274 the potential to be correct PTM-PPI triplets.

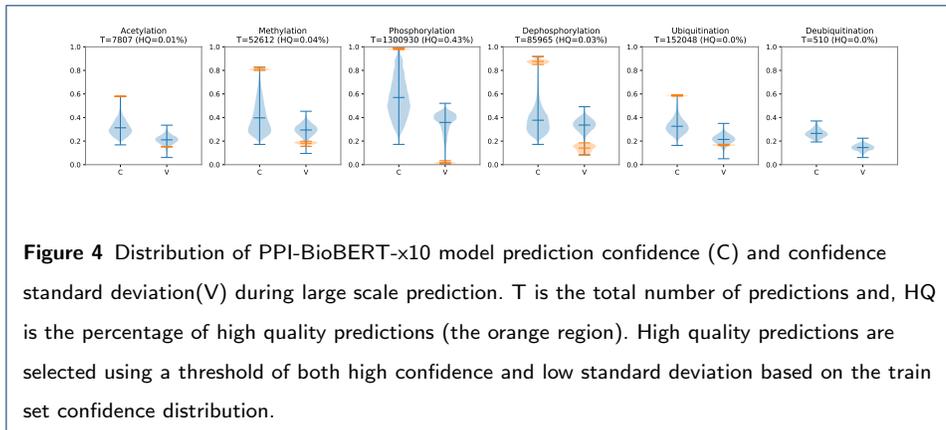


Table 5 The results of large scale prediction. **All** predictions indicate all the predictions from PPI-BioBERT-x10. **Unique** represents unique PTM-PPI triplet predictions. **HQ** represents high quality PTM-PPI after thresholding. **HQ MA** represents high quality PTM-PPI available in multiple abstracts.

PTM	All	All (U)	HQ	HQ (U)	HQ MA	HQ MA (U)
acetylation	7807	6113	1	1	0	0
dephosphorylation	85965	50004	29	29	1	1
deubiquitination	510	460	0	0	0	0
methylation	52612	29914	20	18	4	2
phosphorylation	1300930	381157	5654	4532	1659	537
ubiquitination	152048	78859	4	4	0	0
Total	1599872	546507	5708	4584	1664	540

275 We find that percentage of precision error per PTM does not seem to be correlated
 276 to the number of training samples. For instance, methylation only has 10 training
 277 samples but has the highest precision (57%, 11 out of 19), whereas phosphorylation
 278 has the highest number of training samples but only 20% (6 out of 30) are correct.

Table 6 Human evaluation on randomly sampled subset (30 interactions per PTM, unless there are fewer predictions) selected after thresholding on average confidence and standard deviation.

	acety.	dephosph.	methy.	phosph.	ubiquit.	Total
Correct	0	11	11	6	0	28
Incorrect - DNA Methylation	0	0	2	0	0	2
Incorrect - NER	0	2	1	3	0	6
Incorrect - No trigger word	0	1	0	2	4	7
Incorrect - Opposite type	0	1	0	0	0	1
Incorrect - Relationship not described	0	14	4	19	0	37
Not - sure	1	0	1	0	0	2
Total	1	29	19	30	4	83

279 We analyse the incorrect predictions further and categories them into:

- 280 • **Incorrect - DNA Methylation:** The abstract describes DNA methylation
281 instead of protein methylation. This type of error mainly affects methyla-
282 tion. For instance from the input abstract PubMed 19386523, the prediction
283 [Q01196 (RUNX1), Methylation , Q06455(ETO)] whereas the abstract de-
284 scribes “*RUNX1-ETO fusion gene on DNA methylation*”
- 285 • **Incorrect - NER:** NER has either not identified proteins mentions or not
286 normalised them correctly.
- 287 • **Incorrect - No trigger word:** The abstract does not even mention the trig-
288 ger word describing the predicted PTM. For instance,the prediction [P48431
289 (Sox2), Ubiquitination, Q01860 (Oct3/4)] where the input abstract PubMed
290 22732500 does not mention the trigger *ubiquit.*, but the abstract mainly de-
291 scribes “*knockdown of Sox2 and Oct3/4 gene expression in HCC cells can*
292 *reduce carboplatin-mediated increases in sphere formation and increase cellu-*
293 *lar sensitivity to chemotherapy*”.
- 294 • **Incorrect - Opposite PTM:** The abstract mentions the complementary
295 PTM, *e.g.* (phosphorylation, dephosphorylation), (ubiquitination, deubiqui-
296 tination). For instance, the prediction [Q15746 (myosin light chain kinase),
297 dephosphorylation, Q7Z406 (myosin)] from the abstract PubMed 2967285,
298 the abstract describes phosphorylation - “*phosphorylation of the dephospho-*
299 *rylated brain myosin with myosin light chain kinase and casein kinase II*”
- 300 • **Incorrect - Relationship not described:** This is effectively any other type
301 of error that doesn’t fall into any of the above categories. This category has
302 the largest percentage of error (70%, 37 out of 53 incorrect predictions). An
303 example prediction in this category is [P0870 (IL3), Phosphorylation, P36888
304 (FLT3)] from the PubMed abstract 10720129 doesn’t describe phosphoryla-
305 tion PTM between IL3 and FLT3 but rather “*Somatic mutation of the FLT3*
306 *gene, in which the juxtamembrane domain has an internal tandem duplica-*

307 *tion, is found in 20% of human acute myeloid leukemias and causes constitu-*
 308 *tive tyrosine phosphorylation of the products. In this study, we observed that*
 309 *the transfection of mutant FLT3 gene into an IL3-dependent murine cell line,*
 310 *32D, abrogated the IL3-dependency. ”*

- 311 • **Incorrect - Not related to PPI:** Here prediction is from an abstract that is
 312 not even related to PPI. This is the worst form of false positive error. Example
 313 of this is a prediction [P16070 (CD44), phosphorylation, P60568(IL2)] from
 314 PubMed abstract 7763733.
- 315 • **Unsure:** In this scenario, the human reviewer is unable to decide if the ab-
 316 stract describes the PTM-PPI by solely reading the abstract. Example of this
 317 is a prediction [P01019 (Angiotensin II), phosphorylation, P30556(AT1)] from
 318 PubMed abstract 9347311.

319 In order to reduce the false positive predictions, we further filter the high confi-
 320 dence and low variation predictions to a subset of predictions that are available in
 321 at least 2 abstracts to select 1659 predictions (537 unique PTM-PPI) and verify a
 322 randomly selected subset (30 per PTM) of predictions. We find that the precision on
 323 a randomly sampled subset (available as supplementary Table A13) now improves
 324 to 58.8% as shown in Table 7. The intuition behind this improvement is if the same
 325 PTM-PPI prediction can be inferred from multiple papers, the probability of the
 prediction being right is higher.

Table 7 Include multiple abstracts filter: Human evaluation on randomly sampled subset selected after thresholding on average confidence and standard deviation, with the additional condition that these predictions are present in multiple abstracts. We select 30 interactions per PTM, unless there are fewer predictions.

	methyL.	phosphoryl.	Total
Correct	4	16	20
Incorrect - NER	0	2	2
Incorrect - Not related to PPI	0	1	1
Incorrect - relationship not described	0	7	7
Not - sure	0	4	4
Total	4	30	34

327 We also compare the recall of the predictions with those PTM-PPI consolidated
 328 in iPTMnet [20]^[1] as shown in Table 8. Using PPI-BioBERT-x10, we have identified
 329 37.1% (3270 out of 8805) of iPTMnet PTM-PPIs, with 815 of these found in the high
 330 confidence region. Only a total of 358 PTM-PPI in iPTMnet were sourced directly
 331 from literature using text mining, specifically the rule-based text mining method
 332 RLIMS+[17], compared to the 3270 (815 after confidence thresholding) identified by
 333 PPI-BioBERT-x10. Hence, the recall of PPI-BioBERT-x10 in relation to the PTMs
 334 captured in iPTMnet is substantially higher than the RLIMS+ method, by over
 335 200% after confidence thresholding, demonstrating the robustness of our machine
 336 learning-based approach.

Table 8 Comparison of PPI-BioBERT-x10 predictions with iPTMnet. Of all PTM-PPI entries in iPTMnet (*iP Total*), *iP Unique* represents the subset of unique entries. Of the unique PTM-PPIs the subset that has associated UniProt identifiers is in column *iP Uniprot*s. *iP RLIMS* is the number of unique PPI-PTM sourced from RLIMS+ . The number of all the PPI-BioBERT-x10 predictions that can be recalled in iPTMnet is in *Ours*. *Ours HQ* represents the PPI-BioBERT-x10 predictions after confidence thresholding.

PTM	iP Total	iP Unique	iP Uniprot	Ours	Ours HQ	iP RLIMS
acetylation	141	73	12	0	0	0
methylation	7	4	4	0	0	0
phosphorylation	21050	8949	8805	3270	815	358
ubiquitination	2	1	0	0	0	0

337 Discussion

338 Noise in distant supervision

339 There are two types of noise in our distantly supervised dataset: **(a)** false posi-
 340 tive noise and **(b)** false negative noise. False positive noise occurs where a given
 341 PPI relationship [*Protein1*, *PTM function*, *Protein2*] may not be described in the
 342 abstract but is labelled as describing the relationship. False negative noise occurs
 343 where the relationship [*Protein1*, *PTM function*, *Protein2*] is in fact described in
 344 the abstract but is not labelled as describing the relationship.

345 Our simple heuristic removes noisy training samples, more specifically false posi-
 346 tive noise, *i.e.*, if the normalised Uniprot identifiers of the participating entities

^[1]https://research.bioinformatics.udel.edu/iptmnet_data/files/current/ptm.txt

Table 9 Results of noise levels, after noise reduction, in training data to verify if the PPI relationship is described in the abstract. The training data is randomly sampled (10 samples per interaction type unless the number of available training samples is lower) and verified by a human annotator.

Interaction Type	Correct	Not - sure	Total
acetylation	4	1	5
dephosphorylation	6	4	10
deubiquitination	1	1	2
methylation	4	6	10
phosphorylation	6	4	10
ubiquitination	2	3	5
Total	23	19	42

347 are not mentioned in the abstract and the abstract is not likely to describe the
 348 relationship (unless NER has failed to detect the protein names and normalise
 349 them). We manually inspected the quality of training data after noise reduction
 350 using a randomly sampled subset, as detailed in Table 9. While the noise reduc-
 351 tion heuristic has been fairly effective in removing false positive noise, in 45%
 352 of the cases the human reviewer (author A.E. with basic knowledge on PPIs
 353 and PTMs) was unable to decide whether the PPI relationship is described or
 354 not in the abstract. An example of a sample marked as unsure, is phosphory-
 355 lation between c-Jun N-terminal kinase (JNK) and c-Jun (IntAct entry <https://www.ebi.ac.uk/intact/interaction/EBI-7057279>) in the abstract PubMed
 356 19527717. This shows the complexity of the task beyond language, the need for
 357 prior context or assumed domain knowledge on how PPIs interact to interpret the
 358 abstracts for manual annotation. In the test set, we found that a combination of
 359 high confidence and low variation in confidence has been able detect some of the
 360 incomplete annotation (false negatives) as shown in Table 4.
 361

362 One of the main challenges is that it is difficult to quantify noise without manual
 363 verification. Hence, machine learning approaches that rely on distant supervised
 364 data need to manually verify label quality, at the very least, on randomly subset of
 365 the data. This requires a user interface (UI) that makes it easy for humans to review
 366 the samples effectively. For this work we built a UI (see Figure 5) using Amazon

387 gram word overlap between train and test folds resulting in inflated performance on
388 the test set. The substantially lower real world performance compared to the test
389 set, specifically in the context of BERT models, has also been reported in previous
390 papers [32]. In addition, previous studies in computer vision have also shown neural
391 networks can provide high confidence prediction even when the model is wrong
392 despite robust performance on the test set [7, 33]. In particular a type of model
393 uncertainty, epistemic uncertainty [33], where the model doesn't have sufficient
394 information to make a decision yet has made overconfident predictions continues
395 to be a challenge despite applying machine learning methods. More specifically,
396 in the context of this paper the need for diverse and effective negative samples
397 is pertinent, given the number of false positive predictions despite relatively high
398 confidence. The creation of effective negative samples and generalisability remain
399 open research questions in need of further study.

400 Human augmentation of PTM-PPI extraction

401 The prediction quality can be improved by providing more training samples [34, 33],
402 which requires manual curation. However, manual curation is difficult, time consum-
403 ing and not cost effective [4, 3]. This becomes a chicken and egg problem: automation
404 is meant to speed up curation and reduce manual curation effort, however we need
405 large amounts of manually curated training data so that the models can produce
406 sufficiently reliable predictions.

407 Triaging journals alone takes up approximately 15% of the curation time [35] to
408 select relevant papers, while the rest of the tasks such as named entity recognition
409 and normalisation, detecting relationships take up the rest 85%. In our research,
410 we surface PTM-PPI triplets where the entities are normalised. Hence, our solution
411 supports human augmentation [36], *i.e.* using machine learning to make it much
412 easier for humans to curate. This is one of the major advantages of machine learning
413 compared to rule based systems, as adding more data and retraining can improve

414 the model without manually rewriting the rules. The second advantage of using
415 confidence thresholding is the ability to adjust the thresholds per PTM enabling
416 faster curation of PTM-PPIs with reduced representation, such as for relatively
417 rare interaction types like deubiquitination. In order to take advantage of machine
418 learning, we need to ensure human annotators have: **(a)** convenient user interface
419 so the curator can quickly accept or reject the prediction and **(b)** provide a high
420 hit ratio of correct PPIs using confidence thresholding so that humans can add
421 more PPIs to the knowledge bases faster. This will enable continuous retraining the
422 model with more diverse training data (both positive and negative samples) so the
423 model can evolve and further reduce manual curation effort over time.

424 **Conclusion**

425 We created a distant supervised training dataset for extracting PTM-PPIs, in-
426 cluding annotation for phosphorylation, dephosphorylation, methylation, demethy-
427 lation, ubiquitination, and acetylation from PubMed abstracts by leveraging the
428 IntAct database.

429 Using this dataset we trained an ensemble model, PPI-BioBERT-x10, and applied
430 an approach by Lakshminarayanan *et al.* [6] to improve confidence calibration of the
431 neural network by averaging the confidence scores predicted by the ensemble. We
432 find that the predicted confidence range for each PTM is proportional to the num-
433 ber of training samples available and the predicted confidence range is consistent
434 between train, test validation and large scale predictions. As a result, we extended
435 the work by Lakshminarayanan *et al.* [6] using a threshold per PTM by combining
436 average confidence with confidence standard deviation to improve confidence cali-
437 bration and counteract the effects of class imbalance during calibration, resulting
438 in 100% precision (retaining 19% of the positive predictions) in the test set.

439 We applied PPI-BioBERT-x10 to text-mine 18 million PubMed abstracts, extract-
440 ing 1.6 million PTM-PPI predictions ($\approx 540,000$ unique PTM-PPI) and shortlisting

441 approximately 5,700 (\approx 4,500 unique PTM-PPI) PTM-PPIs with relative high con-
442 fidence and low variation. However, on manual review of a randomly sampled subset
443 we find that the precision of prediction drops to 33.7% and generalisability aspect of
444 knowing when a prediction is correct remains a challenge. By selecting predictions
445 that appear in multiple abstracts (\approx 1600 triplets), the precision on a randomly
446 sampled subset improves to 58.8%. We also find that PPI-BioBERT-x10 is able to
447 identify over 200% more PTM-PPIs compared to RLIMS+[17] a rule based system,
448 revealing the advantage of PPI-BioBERT-x10 with regards to recall compared to
449 a manually crafted rule based system. We also propose that to decrease manual
450 curation effort and to improve the prediction quality over time requires using an
451 effective user interface along with confidence thresholding to allow humans curators
452 to easily accept or reject predictions and continuously retrain the model.

453 In this work, we studied the advantages and challenges of adopting deep learning-
454 based text mining for a new task, PTM-PPI triplet extraction, using distantly
455 supervised data. This work highlights the need for effective confidence calibration,
456 the importance of generalisability beyond test set performance and how real world
457 performance can vary substantially compared to the test set. Further study in-
458 cludes the use of effective adversarial samples to improve the robustness of machine
459 learning models for practical use.

460 **Appendix**

461 A.1 BioBERT Training details

462 We use the following settings to train BioBERT v1.1 based on BERT Base cased
463 using PyTorch 1.4.0. We use gradient accumulation to train with an effective batch
464 size of 64. We train on a Amazon SageMaker P3 instance using a single GPU for
465 approximately 1 hours with early stopping patience.

- 466 • Loss function - cross entropy loss
- 467 • Optimiser - Adam

- 468 • Optimiser Learning rate - .00001,
- 469 • Optimiser Weight decay - 0.01,
- 470 • Batch size - 8
- 471 • Gradient accumulation steps - 8
- 472 • Epochs - 1000
- 473 • Early stopping patience epochs - 20

474 The full source code is here <https://github.com/elangovana/large-scale-ptm-ppi>

475 A.2 BioBERT Token lengths

476 Here we present the distribution of length of the tokens, as a result of using BERT
477 tokeniser, to tokenise the abstract, see Table A10.

Table A10 The distribution of abstract lengths as a result of using BioBERT tokeniser. Approximately 90% of the unique normalised abstracts is under the max limit of 512.

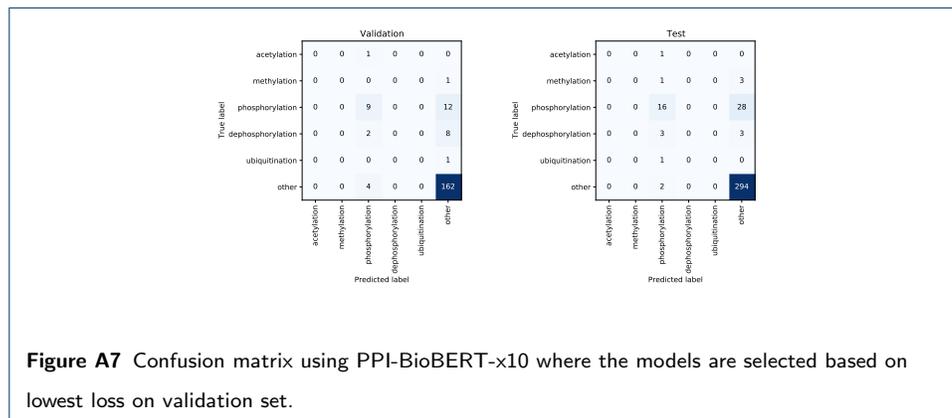
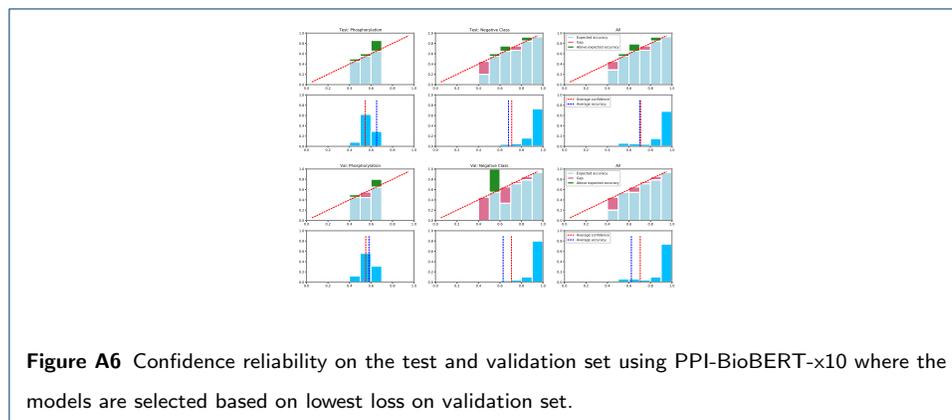
	Train	Test	Validation
count	1305.00	352.00	200.00
mean	378.21	353.21	375.32
std	95.53	102.59	78.16
min	174.00	175.00	194.00
0%	174.00	175.00	194.00
10%	256.00	209.20	241.90
20%	285.00	250.20	323.00
30%	323.00	296.30	335.00
40%	344.60	323.00	361.80
50%	377.00	348.00	379.00
60%	402.00	376.00	389.00
70%	426.00	414.70	420.90
80%	464.00	447.00	467.00
90%	511.00	496.90	472.00
max	612.00	538.00	496.00

478 A.3 Performance using lowest loss

479 The performance on test and validation set where the models in the ensemble are
480 selected based on lowest loss on the validation set. The Table A11 shows the F-score,
481 the reliability diagrams are in Figure A6 and the Figure A7 shows the corresponding
482 confusion matrix.

Table A11 The performance of ensemble PPI-BioBERT-x10 on the test and validation set where the models are selected based on lowest loss on validation set.

InteractionType	Test				Validation			
	Precision	Recall	F1-score	support	Precision	Recall	F1-score	support
acetylation	0.00	0.00	0.00	1	0.00	0.00	0.00	1
dephosphorylation	0.00	0.00	0.00	6	0.00	0.00	0.00	10
methylation	0.00	0.00	0.00	4	0.00	0.00	0.00	1
phosphorylation	66.67	36.36	47.06	44	56.25	42.86	48.65	21
ubiquitination	0.00	0.00	0.00	1	0.00	0.00	0.00	1
macro avg	13.33	7.27	9.41	56	11.25	8.57	9.73	34
micro avg	66.67	28.57	40.00	56	56.25	26.47	36.00	34



483 **A.4 Human verified results**

484 **Acknowledgements**

485 **Funding**

486 This work is supported by Australian Research Council grants DP190101350 (KV) and LP160101469 (KV,YL).

487 **Abbreviations**

488 NER: Named entity recognition

489 PPI: Protein Protein Interaction

Table A12 The distribution of unique normalised protein counts in the abstracts in the test set.

Interaction	count	mean	std	min	25%	50%	75%	max
acetylation	6	3.50	1.38	2.00	3.00	3.00	3.75	6.00
demethylation	2	5.00	0.00	5.00	5.00	5.00	5.00	5.00
dephosphorylation	21	3.62	1.69	2.00	2.00	3.00	5.00	6.00
deubiquitination	1	2.00	nan	2.00	2.00	2.00	2.00	2.00
methylation	11	2.82	0.40	2.00	3.00	3.00	3.00	3.00
phosphorylation	125	3.56	1.95	1.00	2.00	3.00	5.00	9.00
ubiquitination	2	5.00	1.41	4.00	4.50	5.00	5.50	6.00

Table A13 Human verification of randomly sampled subset of 30 PTM-PPI per PTM type. The samples were chosen after confidence calibration and the predictions were present in multiple abstracts within the confidence and standard deviation thresholds.

PubmedId	Participant1 Uniprot	Participant1 Name	Participant2 Uniprot	Participant2 Name	Prediction	Human result
19877273	P23443	SGK1	P42345	mTOR	phosphorylation	Correct
28539327	P01375	TNF	Q13546	RIPK1	phosphorylation	Incorrect - relationship not described
9148953	P06493	Cdc2	P30291	Wee1	phosphorylation	Correct
18028023	P31749	Akt	P42345	mammalian target of rapamycin	phosphorylation	Incorrect - relationship not described
14572154	P31749	PKB/Akt	P60484	PTEN	phosphorylation	Not - sure
1327869	O00757	FBPase-2	Q16875	PFK-2	phosphorylation	Incorrect - relationship not described
22855742	O14757	Checkpoint kinase 1	Q13535	ATR	phosphorylation	Correct
16477614	P23560	Brain-derived neurotrophic factor	P28482	extracellular-signal regulated kinase	phosphorylation	Incorrect - relationship not described
15289331	P31749	Akt	P49841	glycogen synthase kinase-3beta	phosphorylation	Correct
23399841	P01583	IL-1	P10145	CXCL8	phosphorylation	Not - sure
16873552	P12931	Src	P41240	C-terminal Src kinase	phosphorylation	Correct
15994312	Q05397	focal adhesion kinase	Q14289	proline-rich tyrosine kinase-2	phosphorylation	Correct
17190911	P35568	insulin receptor substrate-1	P42336	PI3-K	phosphorylation	Correct
1840422	P17677	B-50	Q92686	neurogranin	phosphorylation	Incorrect - NER
27633668	P20042	eukaryotic translation initiation factor 2	Q9BQ13	heme-regulated inhibitor	phosphorylation	Incorrect - relationship not described
12112022	P31749	PKBalpha	Q14289	Protein kinase B	phosphorylation	Incorrect - relationship not described
23337506	P01138	nerve growth factor	P04629	tropomyosin receptor kinase receptor	phosphorylation	Correct
12759443	O14920	IKK beta	Q04206	p65	phosphorylation	Not - sure
29162743	P16220	CREB	P40763	signal transducer and activator of transcripti...	phosphorylation	Incorrect - NER
9347311	P01019	Angiotensin II	P30556	AT1	phosphorylation	Not - sure
21325496	P24941	cyclin-dependent kinase 2	P38936	p21	phosphorylation	Correct
7763733	P16070	CD44	P60568	IL2	phosphorylation	Incorrect - Not related to PPI
26659448	Q02763	Tek	Q15389	angiotensin-1	phosphorylation	Correct
20110283	P17676	CCAAT/enhancer binding protein beta	P49841	glycogen synthase kinase 3beta	phosphorylation	Correct
15970650	Q15746	myosin light chain kinase	Q7Z406	myosin	phosphorylation	Correct
29496905	P23443	SGK1	P35568	insulin receptor substrate-1	phosphorylation	Correct
17402366	P01562	IFN-alpha	P42224	STAT1	phosphorylation	Correct
18198129	P12931	Src	Q05397	FAK	phosphorylation	Correct
17053882	P31749	Akt	P42336	phosphatidylinositol 3-kinase	phosphorylation	Incorrect - relationship not described
24360952	P12931	Src	Q05397	Focal adhesion kinase	phosphorylation	Correct
27411844	P04637	p53	Q9NQR1	SETD8	methylation	Correct
20614940	P04637	p53	Q96KQ7	G9a	methylation	Correct
24151879	P04637	p53	Q96KQ7	G9a	methylation	Correct
25554733	P04637	p53	Q9NQR1	SETD8	methylation	Correct

490 PTM: Post Translational Modification

491 ECE: Expected Calibration Error

492 NLP: Natural Language Processing

493 BERT: Bidirectional Encoder Representations from Transformers

494 **Availability of data and materials**495 We publish our source code, training, test and validation dataset and the predicted ≈ 5700 PPIs on496 <https://github.com/elangovana/large-scale-ptm-ppi>.

497 The Amazon SageMaker Ground Truth UI code is available at

498 <https://github.com/elangovana/ppi-sagemaker-groundtruth-verification>.499 **Ethics approval and consent to participate**

500 Not applicable

501 Competing interests

502 We certify that there are no financial and non-financial competing interests regarding the materials discussed in the
503 manuscript.

504 Consent for publication

505 Not applicable

506 Authors' contributions

507 AE carried out computational experiments and was responsible for the design, data analysis, data interpretation,
508 developing software, discussion and drafting and editing the manuscript. AE, MD, KV conceptualised the study.
509 MD, YL, DP and KV supervised AE, were involved in the design and execution of the study, supported data
510 interpretation, and contributed to writing the manuscript. All authors read and approved the final manuscript.

511 Author details

512 ¹School of Computing and Information Systems, The University of Melbourne, Melbourne, Australia. ²The Walter
513 and Eliza Hall Institute of Medical Research, Melbourne, Australia. ³Department of Clinical Pathology, Faculty of
514 Medicine, Dentistry & Health Sciences The University of Melbourne, Melbourne, Australia. ⁴School of Computing
515 Technologies RMIT University, Melbourne, Australia.

516 References

- 517 1. Orchard, S., Ammari, M., Aranda, B., Breuza, L., Briganti, L., Broackes-Carter, F., Campbell, N.H., Chavali,
518 G., Chen, C., Del-Toro, N.: The mintact project—intact as a common curation platform for 11 molecular
519 interaction databases. *Nucleic Acids Research* **42**(D1), 358–363 (2013)
- 520 2. Mishra, G.R., Suresh, M., Kumaran, K., Kannabiran, N., Suresh, S., Bala, P., Shivakumar, K., Anuradha, N.,
521 Reddy, R., Raghavan, T.M.: Human protein reference database—2006 update. *Nucleic Acids Research* **34**,
522 411–414 (2006)
- 523 3. Orchard, S., Ammari, M., Aranda, B., Breuza, L., Briganti, L., Broackes-Carter, F., Campbell, N.H., Chavali,
524 G., Chen, C., del-Toro, N., Duesbury, M., Dumousseau, M., Galeota, E., Hinz, U., Iannuccelli, M.,
525 Jagannathan, S., Jimenez, R., Khadake, J., Lagreid, A., Licata, L., Lovering, R.C., Meldal, B., Melidoni, A.N.,
526 Milagros, M., Peluso, D., Peretto, L., Porras, P., Raghunath, A., Ricard-Blum, S., Roechert, B., Stutz, A.,
527 Tognolli, M., van Roey, K., Cesareni, G., Hermjakob, H.: The MIntAct project—IntAct as a common curation
528 platform for 11 molecular interaction databases. *Nucleic Acids Research* **42**(D1), 358–363 (2013).
529 doi:10.1093/nar/gkt1115. <https://academic.oup.com/nar/article-pdf/42/D1/D358/3585170/gkt1115.pdf>
- 530 4. Orchard, S., Hermjakob, H.: Shared resources, shared costs—leveraging biocuration resources. *Database* **2015**
531 (2015). doi:10.1093/database/bav009. bav009.
532 <https://academic.oup.com/database/article-pdf/doi/10.1093/database/bav009/7298544/bav009.pdf>
- 533 5. Consortium, T.U.: UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Research* **47**(D1), 506–515
534 (2018). doi:10.1093/nar/gky1049.
535 <https://academic.oup.com/nar/article-pdf/47/D1/D506/27437297/gky1049.pdf>
- 536 6. Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using
537 deep ensembles. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett,
538 R. (eds.) *Advances in Neural Information Processing Systems*, vol. 30, pp. 6402–6413. Curran Associates, Inc.,
539 ??? (2017). <https://proceedings.neurips.cc/paper/2017/file/9ef2ed4b7fd2c810847ffa5fa85bce38-Paper.pdf>

- 540 7. Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. In: Proceedings of
541 the 34th International Conference on Machine Learning - Volume 70. ICML'17, pp. 1321–1330. JMLR.org, ???
542 (2017)
- 543 8. Craven, M., Kumlien, J.: Constructing biological knowledge bases by extracting information from text sources.
544 In: Proceedings of the Seventh International Conference on Intelligent Systems for Molecular Biology, pp.
545 77–86. AAAI Press, ??? (1999)
- 546 9. Mintz, M., Bills, S., Snow, R., Jurafsky, D.: Distant supervision for relation extraction without labeled data. In:
547 Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint
548 Conference on Natural Language Processing of the AFNLP, pp. 1003–1011. Association for Computational
549 Linguistics, Suntec, Singapore (2009). <https://www.aclweb.org/anthology/P09-1113>
- 550 10. Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H., Kang, J.: Biobert: a pre-trained biomedical language
551 representation model for biomedical text mining. *Bioinformatics* **36**(4), 1234–1240 (2019).
552 doi:10.1093/bioinformatics/btz682
- 553 11. Bunescu, R., Ge, R., Kate, R.J., Marcotte, E.M., Mooney, R.J., Ramani, A.K., Wong, Y.W.: Comparative
554 experiments on learning information extractors for proteins and their interactions. *Artificial Intelligence in*
555 *Medicine* **33**(2), 139–155 (2005)
- 556 12. Pyysalo, S., Ginter, F., Heimonen, J., Björne, J., Boberg, J., Järvinen, J., Salakoski, T.: Bioinfer: a corpus for
557 information extraction in the biomedical domain. *BMC Bioinformatics* **8**(1), 50 (2007)
- 558 13. Hsieh, Y.-L., Chang, Y.-C., Chang, N.-W., Hsu, W.-L.: Identifying protein-protein interactions in biomedical
559 literature using recurrent neural networks with long short-term memory. In: Proceedings of the Eighth
560 International Joint Conference on Natural Language Processing (volume 2: Short Papers), pp. 240–245
- 561 14. Peng, Y., Lu, Z.: Deep learning for extracting protein-protein interactions from biomedical literature. In:
562 BioNLP 2017, pp. 29–38. Association for Computational Linguistics.
563 <https://www.aclweb.org/anthology/W17-2304><http://dx.doi.org/10.18653/v1/W17-2304>
- 564 15. Franceschini, A., Szklarczyk, D., Frankild, S., Kuhn, M., Simonovic, M., Roth, A., Lin, J., Minguez, P., Bork,
565 P., von Mering, C., Jensen, L.J.: STRING v9.1: protein-protein interaction networks, with increased coverage
566 and integration. *Nucleic Acids Research* **41**(D1), 808–815 (2012). doi:10.1093/nar/gks1094.
567 <https://academic.oup.com/nar/article-pdf/41/D1/D808/3617210/gks1094.pdf>
- 568 16. Tudor, C.O., Ross, K.E., Li, G., Vijay-Shanker, K., Wu, C.H., Arighi, C.N.: Construction of phosphorylation
569 interaction networks by text mining of full-length articles using the eFIP system. *Database* **2015** (2015).
570 doi:10.1093/database/bav020. bav020.
571 <https://academic.oup.com/database/article-pdf/doi/10.1093/database/bav020/16975967/bav020.pdf>
- 572 17. Torii, M., Arighi, C.N., Li, G., Wang, Q., Wu, C.H., Vijay-Shanker, K.: Rlims-p 2.0: A generalizable rule-based
573 information extraction system for literature mining of protein phosphorylation information **12**(1), 17–29 (2015).
574 doi:10.1109/TCBB.2014.2372765
- 575 18. Szklarczyk, D., Gable, A.L., Lyon, D., Junge, A., Wyder, S., Huerta-Cepas, J., Simonovic, M., Doncheva, N.T.,
576 Morris, J.H., Bork, P., Jensen, L.J., Mering, C.: STRING v11: protein-protein association networks with
577 increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids*
578 *Research* **47**(D1), 607–613 (2018). doi:10.1093/nar/gky1131.
579 <https://academic.oup.com/nar/article-pdf/47/D1/D607/27437323/gky1131.pdf>
- 580 19. Chen, Q., Panyam, N.C., Elangovan, A., Verspoor, K.: BioCreative VI Precision Medicine Track system
581 performance is constrained by entity recognition and variations in corpus characteristics. *Database* **2018** (2018).
582 doi:10.1093/database/bay122. bay122.

- 583 <https://academic.oup.com/database/article-pdf/doi/10.1093/database/bay122/27329364/bay122.pdf>
- 584 20. Huang, H., Arighi, C.N., Ross, K.E., Ren, J., Li, G., Chen, S.-C., Wang, Q., Cowart, J., Vijay-Shanker, K., Wu,
585 C.H.: iPTMnet: an integrated resource for protein post-translational modification network discovery. *Nucleic*
586 *Acids Research* **46**(D1), 542–550 (2017). doi:10.1093/nar/gkx1104.
587 <https://academic.oup.com/nar/article-pdf/46/D1/D542/23162331/gkx1104.pdf>
- 588 21. Hornbeck, P.V., Zhang, B., Murray, B., Kornhauser, J.M., Latham, V., Skrzypek, E.: PhosphoSitePlus, 2014:
589 mutations, PTMs and recalibrations. *Nucleic Acids Research* **43**(D1), 512–520 (2014).
590 doi:10.1093/nar/gku1267. <https://academic.oup.com/nar/article-pdf/43/D1/D512/17437800/gku1267.pdf>
- 591 22. Elangovan, A., He, J., Verspoor, K.: Memorization vs. generalization : Quantifying data leakage in NLP
592 performance evaluation. In: Proceedings of the 16th Conference of the European Chapter of the Association for
593 Computational Linguistics: Main Volume, pp. 1325–1335. Association for Computational Linguistics, Online
594 (2021). <https://www.aclweb.org/anthology/2021.eacl-main.113>
- 595 23. Futoma, J., Simons, M., Panch, T., Doshi-Velez, F., Celi, L.A.: The myth of generalisability in clinical research
596 and machine learning in health care. *The Lancet Digital Health* **2**(9), 489–492 (2020).
597 doi:10.1016/S2589-7500(20)30186-2
- 598 24. Li, G., Wu, C., Vijay-Shanker, K.: Noise reduction methods for distantly supervised biomedical relation
599 extraction. In: BioNLP 2017, pp. 184–193. Association for Computational Linguistics, Vancouver, Canada,
600 (2017). doi:10.18653/v1/W17-2323. <https://www.aclweb.org/anthology/W17-2323>
- 601 25. Zhang, H., Guan, R., Zhou, F., Liang, Y., Zhan, Z.-H., Huang, L., Feng, X.: Deep residual convolutional neural
602 network for protein-protein interaction extraction. *IEEE Access* **7**, 89354–89365 (2019).
603 doi:10.1109/ACCESS.2019.2927253
- 604 26. Orchard, S., Kerrien, S., Abbani, S., Aranda, B., Bhate, J., Bidwell, S., Bridge, A., Briganti, L., Brinkman,
605 F.S., Cesareni, G.: Protein interaction data curation: the international molecular exchange (imex) consortium.
606 *Nature Methods* **9**(4), 345 (2012)
- 607 27. Wei, C.-H., Kao, H.-Y., Lu, Z.: Gnormplus: an integrative approach for tagging genes, gene families, and
608 protein domains. *BioMed Research International* **2015** (2015)
- 609 28. Brown, G.R., Hem, V., Katz, K.S., Ovetsky, M., Wallin, C., Ermolaeva, O., Tolstoy, I., Tatusova, T., Pruitt,
610 K.D., Maglott, D.R., Murphy, T.D.: Gene: a gene-centered information resource at NCBI. *Nucleic Acids*
611 *Research* **43**(D1), 36–42 (2014). doi:10.1093/nar/gku1055.
612 <https://academic.oup.com/nar/article-pdf/43/D1/D36/7314925/gku1055.pdf>
- 613 29. Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for
614 language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the
615 Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers),
616 pp. 4171–4186. Association for Computational Linguistics.
617 <https://www.aclweb.org/anthology/N19-1423><http://dx.doi.org/10.18653/v1/N19-1423>
- 618 30. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent
619 neural networks from overfitting. *Journal of Machine Learning Research* **15**(56), 1929–1958 (2014)
- 620 31. Yadav, S., Ekbal, A., Saha, S., Kumar, A., Bhattacharyya, P.: Feature assisted stacked attentive shortest
621 dependency path based bi-lstm model for protein-protein interaction. *Knowledge-Based Systems* (2018)
- 622 32. McCoy, R.T., Min, J., Linzen, T.: BERTs of a feather do not generalize together: Large variability in
623 generalization across models with similar test set performance. In: Proceedings of the Third BlackboxNLP
624 Workshop on Analyzing and Interpreting Neural Networks for NLP, pp. 217–227. Association for Computational
625 Linguistics, Online (2020). doi:10.18653/v1/2020.blackboxnlp-1.21.

- 626 <https://aclanthology.org/2020.blackboxnlp-1.21>
- 627 33. Hüllermeier, E., Waegeman, W.: Aleatoric and epistemic uncertainty in machine learning: An introduction to
628 concepts and methods. *Machine Learning* **110**(3), 457–506 (2021)
- 629 34. Kendall, A., Gal, Y.: What uncertainties do we need in bayesian deep learning for computer vision? In: Guyon,
630 I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural
631 Information Processing Systems*, vol. 30. Curran Associates, Inc., ??? (2017).
632 <https://proceedings.neurips.cc/paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf>
- 633 35. Mottin, L., Pasche, E., Gobeill, J., Rech de Laval, V., Gleizes, A., Michel, P.-A., Bairoch, A., Gaudet, P., Ruch,
634 P.: Triage by ranking to support the curation of protein interactions. *Database* **2017** (2017).
635 [doi:10.1093/database/bax040](https://doi.org/10.1093/database/bax040). [bax040](https://doi.org/10.1093/database/bax040).
636 <https://academic.oup.com/database/article-pdf/doi/10.1093/database/bax040/19234641/bax040.pdf>
- 637 36. Raisamo, R., Rakkolainen, I., Majoranta, P., Salminen, K., Rantala, J., Farooq, A.: Human augmentation: Past,
638 present and future. *International Journal of Human-Computer Studies* **131**, 131–143 (2019).
639 [doi:10.1016/j.ijhcs.2019.05.008](https://doi.org/10.1016/j.ijhcs.2019.05.008). 50 years of the *International Journal of Human-Computer Studies*. Reflections
640 on the past, present and future of human-centred technologies

641 **Figures**

642 **Tables**

643 **Additional Files**

644 Additional file 1 — Sample additional file title

645 Additional file descriptions text (including details of how to view the file, if it is in a non-standard format or the file
646 extension). This might refer to a multi-page table or a figure.

Table A14 Human verification, without taking into account predictions in multiple abstracts, of randomly sampled subset of 30 PTM-PPI per PTM type. The samples were chosen after confidence calibration.

Pubmed Id	Participant1 Uniprot	Participant1 Name	Participant2 Uniprot	Participant2 Name	Prediction	Human result
22732500	P48431	Sox2	Q01860	Oct3/4	ubiquitination	Incorrect - No trigger word
21492879	P48023	Fas ligand	Q81XH7	Th1	ubiquitination	Incorrect - No trigger word
19428371	O75943	RAD24	Q99638	RAD9	ubiquitination	Incorrect - No trigger word
17168739	P02786	transferrin receptor	P09544	IRP	ubiquitination	Incorrect - No trigger word
7789533	P19338	C23	Q9Y2Q3	glutathione S-transferase	phosphorylation	Incorrect - No trigger word
19729590	P23443	S6 kinase	P35568	insulin receptor substrate 1	phosphorylation	Incorrect - relationship not described
23478265	O00571	DEAD box protein 3	Q9UHD2	TANK-binding kinase 1	phosphorylation	Incorrect - relationship not described
16365045	O75581	low density lipoprotein receptor-related prote...	P49841	GSK3beta	phosphorylation	Incorrect - relationship not described
19306938	O60500	nephrin	Q9Y5K6	CD2-associated protein	phosphorylation	Incorrect - relationship not described
10720129	P08700	IL3	P36888	FLT3	phosphorylation	Incorrect - relationship not described
27456486	P20701	LFA-1	Q05397	Protein tyrosine kinase 2	phosphorylation	Incorrect - relationship not described
2780290	P19338	nucleolin	P50613	Protein kinase	phosphorylation	Incorrect - NER
17108171	P04629	TrkA	P08138	p75 neurotrophin receptor	phosphorylation	Incorrect - relationship not described
10884023	P16220	Ca2+/cAMP responsive element binding protein	P18509	pituitary adenylate cyclase-activating polypep...	phosphorylation	Incorrect - relationship not described
27118568	O14974	myosin phosphatase-targeting subunit 1	P61586	Rho A	phosphorylation	Incorrect - relationship not described
17353368	Q13094	SH2 domain-containing leukocyte protein of 76 kD	Q92918	hematopoietic progenitor kinase 1	phosphorylation	Correct
20122754	P05067	APP	P45983	JNK	phosphorylation	Correct
23023514	P42574	caspase 3	Q16555	Collapsin response mediator protein-2	phosphorylation	Incorrect - relationship not described
7852364	P78356	PtdIns(4)P 5-kinase	Q9Y2I7	Fab1p	phosphorylation	Incorrect - relationship not described
26464283	P00533	EGFR	P35354	cyclooxygenase-2	phosphorylation	Correct
17211494	Q06124	PTPN11	Q13480	Grb2 associated binder 1	phosphorylation	Incorrect - No trigger word
29514920	P30291	Wee1	Q01850	Cdr2	phosphorylation	Correct
17045592	P49841	GSK-3beta	Q00535	cdk5	phosphorylation	Incorrect - relationship not described
22633971	P35548	Msx2	Q13950	Runx2	phosphorylation	Incorrect - relationship not described
9852063	P42681	tyrosine kinase	P63252	Kir2.1	phosphorylation	Correct
19879273	P31749	Akt	P35222	beta-catenin	phosphorylation	Correct
3014310	P04271	S-100	P50613	protein kinase	phosphorylation	Incorrect - relationship not described
9883577	Q15049	myosin light chain	Q32MK0	MLC kinase	phosphorylation	Incorrect - NER
29162743	P16220	CREB	P40763	signal transducer and activator of transcripti...	phosphorylation	Incorrect - NER
11567986	P62993	Grb2	Q92835	SH2-containing inositol 5'-phosphatase	phosphorylation	Incorrect - relationship not described
30158515	Q00535	CDK5	Q15078	p35	phosphorylation	Incorrect - relationship not described
1707345	O00459	p85 beta	P27986	p85 alpha	phosphorylation	Incorrect - relationship not described
25684187	P01588	Erythropoietin	P17302	connexin43	phosphorylation	Incorrect - relationship not described
25446109	P04792	HSP27	P45983	JNK	phosphorylation	Incorrect - relationship not described
21080372	P08670	eukaryotic translation initiation factor 4A-I ...	Q99873	PRMT1	methylation	Correct
26503212	Q8NCA5	FAM98A	Q99873	PRMT1	methylation	Correct
24412544	Q96L73	NSD1	Q9BQQ3	p65	methylation	Incorrect - relationship not described
25749972	O75164	KDM4A	P42345	mTOR	methylation	Not - sure
25554733	P04637	p53	Q9NQR1	SETD8	methylation	Correct
20615470	Q15047	KMT1E	Q8WTV6	KMT7	methylation	Incorrect - relationship not described
24151879	P04637	p53	Q96KQ7	G9a	methylation	Correct
23469257	O15047	Set1	O75934	Dam1	methylation	Correct
18472002	P04637	p53	P61978	hnRNP K	methylation	Correct
19386523	Q01196	RUNX1	Q06455	ETO	methylation	Incorrect - DNA Methylation
20614940	P04637	p53	Q96KQ7	G9a	methylation	Correct
27411844	P04637	p53	Q9NQR1	SETD8	methylation	Correct
20231378	P11274	B cell antigen receptor	P11912	Igalpha	methylation	Incorrect - NER
24129573	P01562	IFN-	Q01628	Interferon-induced transmembrane protein 3	methylation	Incorrect - relationship not described
26391684	P49757	Numb	Q9NQR1	SET8	methylation	Correct
25350748	Q13283	G3BP1	Q99873	PRMT1	methylation	Correct
26126536	Q99873	PRMT1	Q9HAU4	Smurf2	methylation	Correct
27563394	P35270	SPR	Q8TEK3	disruptor of telomeric silencing 1-like protein	methylation	Incorrect - DNA Methylation
25748791	Q8WWM7	ataxin-2-like	Q99700	ataxin-2	methylation	Incorrect - relationship not described
2967285	Q15746	myosin light chain kinase	Q7Z406	myosin	dephosphorylation	Incorrect - Opposite type
9417127	P01375	tumor necrosis factor-alpha	P29350	SHP-1	dephosphorylation	Incorrect - relationship not described
6245872	P10145	NaF	P50613	protein kinase	dephosphorylation	Incorrect - NER
11264356	Q15181	protein phosphatase 1	Q9BY44	eIF-2alpha	dephosphorylation	Correct
9079814	P35670	WC1	P60568	interleukin (IL)-2	dephosphorylation	Incorrect - relationship not described
8083198	P18031	PTP-1B	P50391	PP-1	dephosphorylation	Incorrect - relationship not described
12853978	P00533	epidermal growth factor receptor	P42574	caspase-3	dephosphorylation	Incorrect - relationship not described
12138178	P42224	Stat1	P42226	Stat6	dephosphorylation	Incorrect - relationship not described
10391142	P01308	insulin	P06213	insulin receptor	dephosphorylation	Incorrect - relationship not described
9119896	P49840	glycogen synthase kinase-3alpha (kinase FA/GSK...	Q9Y2R2	protein tyrosine phosphatase	dephosphorylation	Incorrect - NER
16620785	P31749	Akt	P49840	glycogen synthase kinase-3alpha/beta	dephosphorylation	Correct
7683660	P01308	Insulin	P19338	nucleolin	dephosphorylation	Incorrect - relationship not described
10405762	O14494	PAP2	O43688	hPAP2c	dephosphorylation	Incorrect - relationship not described
10405762	O14494	PAP2	O14495	hPAP2b	dephosphorylation	Incorrect - relationship not described
10496881	P06493	p34cdc2	P30307	cdc25-C	dephosphorylation	Correct
11773439	P17706	TC-PTP	P42229	signal transducer and activator of transcripti...	dephosphorylation	Correct
15823043	P19634	Na(+)/H(+) exchanger isoform 1	Q15181	PP1	dephosphorylation	Correct
15269244	P06213	insulin receptor	P18031	protein-tyrosine phosphatase (PTP) 1B	dephosphorylation	Correct
11495355	P21128	P11	P60903	P10	dephosphorylation	Incorrect - relationship not described
11707519	P16220	cAMP response element-binding protein	P27361	extracellular-signal related kinases 1 and 2	dephosphorylation	Incorrect - No trigger word
1321126	P06213	insulin receptor	P10586	LAR	dephosphorylation	Correct
9989818	P04637	p53	P28749	p107/E2F	dephosphorylation	Incorrect - relationship not described
17046078	P12931	Src	P29350	SHP1	dephosphorylation	Incorrect - relationship not described
9380758	P13569	cystic fibrosis transmembrane conductance regu...	P35813	PP2Calpha	dephosphorylation	Correct
24691491	P20042	eukaryotic translation initiation factor 2	P63000	Rac1	dephosphorylation	Incorrect - relationship not described
11495355	P21128	P11	Q15257	PP2A	dephosphorylation	Incorrect - relationship not described
9989818	P28749	p107/E2F	Q15257	PP2A	dephosphorylation	Correct
11222372	Q15257	PP2A	Q92934	BAD	dephosphorylation	Correct
8550570	Q15181	PP-1	Q15257	phosphatase-2A	dephosphorylation	Correct
14731392	P25440	Brd2	Q92831	PCAF	acetylation	Not - sure