

1 **Deep learning integration of molecular and interactome data for protein-compound**  
2 **interaction prediction**

3

4 Narumi Watanabe, Yuuto Ohnuki, Yasubumi Sakakibara

5

6 Department of Biosciences and Informatics, Keio University, Yokohama, Kanagawa 223-

7 8522, Japan

8

9 **Corresponding Author**

10 Yasubumi Sakakibara

11 3-14-1 Hiyoshi, Kohoku-ku, Yokohama, 223-8522, Japan

12 Phone/FAX: +81-45-566-1791

13 E-mail: yasu@bio.keio.ac.jp.

14

15 **Abstract**

16 **Motivation:** Virtual screening, which can computationally predict the presence or absence of  
17 protein-compound interactions, has attracted attention as a large-scale, low-cost, and short-  
18 term search method for seed compounds. Existing machine learning methods for predicting  
19 protein-compound interactions are largely divided into those based on molecular structure  
20 data and those based on network data. The former utilize information on proteins and  
21 compounds, such as amino acid sequences and chemical structures, while the latter utilize  
22 interaction network data, such as data on protein-protein interactions and compound-  
23 compound interactions. However, few attempts have been made to combine both types of  
24 data in molecular information and interaction networks.

25 **Results:** We developed a deep learning-based method that integrates protein features,

26 compound features, and multiple types of interactome data to predict protein-compound  
27 interactions. We designed three benchmark datasets with different difficulties and evaluated  
28 the performance on them. The performance evaluations show that our deep learning  
29 framework for integrating molecular structure data and interactome data outperforms state-of-  
30 the-art machine learning methods for protein-compound interaction prediction tasks. The  
31 performance improvement is proven to be statistically significant by the Wilcoxon signed-  
32 rank test. This reveals that the multi-interactome captures different perspectives than amino  
33 acid sequence homology and chemical structure similarity, and both type of data have a  
34 synergistic effect in improving prediction accuracy. Furthermore, experiments on three  
35 benchmark datasets show that our method is more robust than existing methods in accurately  
36 predicting interactions between proteins and compounds that are unseen in the training  
37 samples.

38

### 39 **Keywords**

40 Protein-compound interaction, Deep learning, Heterogeneous interaction network, Integration  
41

### 42 **Introduction**

43 Most compounds that currently act as drugs bind to target proteins that can cause disease, and  
44 these compounds can control their functions. Therefore, it is necessary to search for  
45 compounds that can interact with the target protein when developing new drugs, and this  
46 process must be performed efficiently. However, determining the interaction of a large  
47 number of protein-compound pairs via experiments is expensive in terms of time and cost.  
48 Virtual screening that can computationally classify the presence or absence of protein-  
49 compound interactions has attracted attention as a large-scale, low-cost, short-term search  
50 method for hit compounds. In particular, the method of using machine learning for virtual

51 screening is considered to be applicable to a wide variety of proteins and compounds.

52         Machine learning-based methods for predicting protein-compound interactions are  
53 largely divided into those based on molecular structure data and those based on network data.  
54 The former use protein and compound data represented in amino acid sequences and  
55 chemical structure formulas, and they can be applied to proteins when a docking simulation  
56 cannot be performed because the three-dimensional structure is unknown. In our previous  
57 study [1-3], using positive interactions between drug compounds and their target proteins  
58 downloaded from DrugBank (a database that contains information on existing drug  
59 compounds) [4] and negative interactions consisting of randomly combined compounds and  
60 proteins, we performed binary classification using a support vector machine (SVM). A  
61 prediction accuracy of 85.1% was achieved. Based on this result, we developed COPICAT, a  
62 comprehensive prediction system for protein-compound interactions, which enabled us to  
63 search for lead compounds from a huge compound database, PubChem [5], consisting of tens  
64 of millions of compounds.

65         Deep learning, a method developed in the field of machine learning, has been used in  
66 a variety of fields in recent years because it has achieved high prediction accuracy in fields  
67 such as image recognition, speech recognition, and compound activity prediction [6]. Deep  
68 learning-based protein-compound interaction prediction methods have been developed based  
69 on molecular structure data [7-10]. However, these existing deep learning-based methods  
70 utilize only information based on amino acid sequences and chemical structures, so the  
71 functional properties of proteins and compounds have not yet been incorporated into  
72 prediction.

73         The other type of machine learning approach for protein-compound interaction  
74 prediction is based on network data. An interaction network is commonly used to  
75 comprehensively represent interactions between molecules. For example, the protein-protein

76 interaction network represents the relationships among physically interacting proteins. In the  
77 protein-protein interaction network, a node is a protein, and an edge is drawn between a pair  
78 of proteins that interact with each other.

79         Some previous studies incorporated data from multiple interaction networks to predict  
80 molecular interactions. For instance, multi-modal graphs were proposed to handle three types  
81 of interactions: protein-protein, protein-drug, and polypharmacy side effects. A deep learning  
82 method, Decagon [11], for multi-modal graphs was proposed to predict polypharmacy side  
83 effects. DTINet [12] and NeoDTI [13] were designed and developed as graph-based deep  
84 learning frameworks to integrate heterogeneous networks for drug-target interaction  
85 predictions and drug repositioning. In particular, NeoDTI exhibited a substantial performance  
86 improvement over other state-of-the-art prediction methods based on multiple interaction  
87 network data.

88         In addition to predicting protein-compound interactions, several studies have  
89 predicted other types of molecular interactions. Protein-protein interactions induce many  
90 biological processes within a cell, and experiential and computational methods have been  
91 developed to identify various protein-protein interactions. High-throughput experimental  
92 methods such as yeast two-hybrid screening were developed to discover and validate protein-  
93 protein interactions on a large scale. Computational methods for protein-protein interaction  
94 predictions employ various machine learning methods, such as SVM with feature extraction  
95 engineering [14]. The recurrent convolutional neural network (CNN), which is a deep  
96 learning method, was applied to sequence-based prediction for protein-protein interactions  
97 [15]. Compounds that can interact with each other are often represented as compound-  
98 compound interactions (also known as chemical-chemical interactions); interactive  
99 compounds tend to share similar functions. Compound-compound interactions, called drug-  
100 drug interactions, can be used to predict side effects based on the assumption that interacting

101 compounds are more likely to have similar toxicity [16]. A computational approach to  
102 compound-compound interaction predictions has been studied with various machine learning  
103 methods, including end-to-end learning with a CNN based on the SMILES representation  
104 [17].

105         The purpose of this study is to improve prediction accuracy by integrating molecular  
106 structure data and heterogeneous interactome data into a deep learning method for predicting  
107 protein-compound interactions. In addition to the molecular information (amino acid  
108 sequence and chemical structure) itself, protein-protein interaction network data with similar  
109 reaction pathways or physical direct binding and compound network data linking compounds  
110 with similar molecular activities are incorporated into the deep learning model as multiple-  
111 interactome data. To the best of our knowledge, there are no deep learning-based solutions  
112 for predicting protein-compound interactions that integrate multiple heterogeneous  
113 interactome data along with the direct input of amino acid sequences and chemical structures.

114         This study proposes a method for predicting protein-compound (drug-target)  
115 interactions by combining protein features, compound features, and network context for both  
116 proteins and compounds. The network context comes in the form of protein-protein  
117 interactions from the STRING database [18], and the compound-compound interactions come  
118 from the STITCH database [19]. The protein-protein interaction network and compound-  
119 compound interaction network are processed using node2vec [20] to generate feature vectors  
120 for each protein node and each compound node in the interaction networks in an  
121 unsupervised manner. Each network-based representation is then combined with additional  
122 features extracted from a CNN applied to the amino acid sequence of a protein and from the  
123 extended-connectivity fingerprint (ECFP) of a compound. The final combined protein  
124 representations and compound representations are used to make a protein-compound  
125 interaction prediction with an additional fully connected layer. The overall learning

126 architecture is illustrated in Figure 1.

127 We designed three benchmark datasets with different difficulties and evaluated the  
128 performance on them. In the performance evaluations, we demonstrate that integrating the  
129 molecular structure data and multiple heterogeneous interactome data has a synergistic effect  
130 in improving the accuracy of protein-compound interaction prediction. Furthermore,  
131 performance comparisons with state-of-the-art deep learning methods based on molecular  
132 information [10] and those based on interaction network data [13] as well as the traditional  
133 machine learning methods SVM and random forest show that our model exhibits significant  
134 performance improvements in the most important evaluation measures: AUROC, AUPRC, F-  
135 measure and accuracy, while the other methods show low values of these measures. The  
136 improvement is proven to be statistically significant by the Wilcoxon signed-rank test.  
137 Finally, we analyse whether protein-protein interactions capture a different perspective than  
138 amino acid sequence homology and whether compound-compound interactions capture a  
139 different perspective than chemical structure similarity.

140

## 141 **Methods**

142

### 143 *1D-CNN for Encoding Protein Data*

144 First, the protein data were applied to a one-dimensional convolutional neural network (1D-  
145 CNN). For the protein input, a one-hot vector was used for the distributed representation of  
146 an amino acid sequence of 20 dimensions at a height and width of 8,923 dimensions with the  
147 maximum length of amino acid sequences.

148 An amino acid sequence is a linear structure (1-D grid). In this study, a filter (kernel)  
149 with a one-dimensional convolution operation was applied to the linear structure. Here, a  
150 “one-dimensional” convolutional operation for linear structures was interpreted as scanning

151 the input structure in only one direction along the linear structure with a filter of the same  
152 height (dimension) as that of the distributed representation of the input.

153

#### 154 ***One-Dimensional (1D) Convolutional Layer***

155 We denote  $A = [\mathbf{a}_1^{(1)}, \mathbf{a}_2^{(1)}, \dots, \mathbf{a}_q^{(1)}]$  as an input vector sequence that corresponds to the one-  
156 hot vector representation of an amino acid sequence (as illustrated in Figure 1). For a filter  
157 function in the  $l$ -th hidden layer of the CNN, the input is the set of feature maps in the  $(l-1)$ -th  
158 hidden layer  $\mathbf{x}_{i:i+r-1,j}^{(l-1)} = c_{i,j}^{(l-1)} \in \mathbb{R}^{m \times n}$ , where  $r$  is the size of the filter,  $m$  is the size of the  
159 feature map, and  $n$  is the number of feature maps. The output of the  $k$ -th filter is a feature  
160 map of the  $l$ -th layer  $c_i^{(l,k)} \in \mathbb{R}^m$ , which is defined as follows:

$$c_i^{(l,k)} = f(\mathbf{W}^{(l,k)} c_{i,j}^{(l-1)} + \mathbf{b}^{(l,k)}),$$

161 where  $f$  is an activation function (leaky-ReLU),  $\mathbf{W}^{(l,k)} \in \mathbb{R}^{m \times n \times d}$  is the weight matrix of the  
162  $k$ -th filter in the  $l$ -th convolutional layer, and  $\mathbf{b}^{(l,k)}$  is the bias vector. The average-pooling  
163 mechanism is applied to every convolution output. To obtain the final output  $\mathbf{y} =$   
164  $\{y^{(t,1)}, y^{(t,2)}, \dots, y^{(t,s)}\}$ , global max-pooling is used as follows:

$$y^{(t,k)} = \max_i(c_i^{(t,k)}),$$

165 where  $t$  represents the last layer of the CNN and  $s$  represents the number of filters in the last  
166 layer.

167

#### 168 ***Extended-Connectivity Fingerprint (ECFP) for Compound Data***

169 The extended-connectivity fingerprint (ECFP, also known as the circular fingerprint or  
170 Morgan fingerprint) [21] is the most commonly used feature representation for representing a  
171 property of the chemical structure of a compound. This algorithm first searches the partial  
172 structures around each atom recurrently, then assigns an integer identifier to each partial

173 structure and expresses this as a binary vector by using a hash function. Potentially, an  
174 infinite number of structures exist in the chemical space; consequently, the ECFP requires  
175 vectors with a large number of bits (usually 1,024 - 2,048 bits). In this study, we employed an  
176 ECFP with 1024 bits as the feature representation for the chemical structure of a compound.  
177

## 178 *Feature Representation Learning for Protein-protein and Compound-Compound*

### 179 *Interactions*

180 A protein-protein interaction network that connects physically interacting proteins and a  
181 compound-compound interaction network that connects compounds with similar molecular  
182 activities were input as multiple-interactome data. First, each network was represented as a  
183 graph. A node is a protein in the protein-protein network and a compound in the compound-  
184 compound network. An edge is drawn between a pair of proteins (compounds) that interact  
185 with each other. By applying this graph to “node2vec” [20], the feature vector of each node  
186 was obtained in an unsupervised manner; node2vec is a deep learning method that learns the  
187 feature representation of nodes in a graph and obtains a feature vector for each node.  
188 Node2vec is a graph embedding algorithm that can be applied to any type of graph, and it can  
189 learn a feature vector such that nodes that are nearby on the graph are also close in the  
190 embedded feature space. In other words, the inner product of the feature vectors of the nearby  
191 nodes is high. It is known that the accuracy of the node classification task and the link  
192 prediction task using the obtained feature representations of nodes is higher than that of the  
193 existing methods.

194 The node2vec algorithm was applied to the protein-protein interaction network and  
195 the compound-compound interaction network. Using a protein and a compound as vertices,  
196 the interaction networks were converted into graphs with edge weights based on the  
197 reliability of the experimental data and the similarity in molecular activity. Node2vec



198 (version 0.2.2) from the Python library, which implemented the node2vec algorithm, was  
199 applied to the converted graph. The node2vec parameters used the default values (embedding  
200 dimensions: 128; number of nodes searched in one random walk: walk\_length=80; number of  
201 random walks per node: num\_walk=10; control of probability of revisiting a walk node: p=1;  
202 control of the search speed and range: r=1; whether to reflect the graph weight:  
203 weight\_key=weight).

204 Let a protein-protein interaction network be expressed by a weighted graph  
205  $G_{protein} = (V_{protein}, E_{protein})$  and a compound-compound interaction network by a  
206 weighted graph  $G_{compound} = (V_{compound}, E_{compound})$ . By applying node2vec to these  
207 graphs, the feature representations can be obtained and are denoted as  $N_{protein} =$   
208  $node2vec(G_{protein}) \in \mathbb{R}^d$  and  $N_{compound} = node2vec(G_{compound}) \in \mathbb{R}^d$  for a dimension  
209 of  $d$ .

210

### 211 ***Deep Learning Model Structure for Integrating Molecular Information and the*** 212 ***Interaction Network***

213 The feature vectors obtained from the 1D-CNN for the amino acid sequence and node2vec  
214 for the protein-protein interaction network were concatenated and fed to the final output  
215 layer. Similarly, the feature vectors from the ECFP for the chemical structure and node2vec  
216 for the compound-compound interaction network were concatenated and fed to the final  
217 output layer.

218 We designed an output layer consisting of an element-wise product calculation  
219 followed by a fully connected layer, which is an extension of cosine similarity. The  
220 architecture is illustrated in Figure 2. First, the feature vectors for the proteins and  
221 compounds were mapped onto the same latent space with a fixed dimension  $d$  by applying  
222 fully connected layers. The similarity between the vector for proteins and the vector for

223 compounds on the latent space was calculated by the element-wise product calculation  
224 method followed by a fully connected layer. When a pair of proteins and compounds was  
225 input, if the similarity was higher than some predefined threshold (where the default was 0.5),  
226 it was predicted that there was an interaction between the input pair. If the similarity was  
227 lower, it was predicted that there was no interaction. This model is denoted as the “integrated  
228 model”.

229 More precisely, let  $\mathbf{a}_{protein}$  denote the feature vector output by the 1D-CNN for an  
230 amino acid sequence, and let  $\mathbf{b}_{compound}$  denote the feature vector of the ECFP for the  
231 chemical structure of a compound. Let  $\mathbf{N}_{protein}$  and  $\mathbf{N}_{compound}$  denote the feature  
232 representations obtained from node2vec for the protein-protein interaction network and the  
233 compound-compound interaction network. Two feature vectors  $\mathbf{a}_{protein}$  and  $\mathbf{N}_{protein}$  were  
234 concatenated as one vector  $\mathbf{v}_{protein}$  for the protein multi-modal feature. Two feature vectors  
235  $\mathbf{b}_{compound}$  and  $\mathbf{N}_{compound}$  were concatenated as one vector  $\mathbf{v}_{compound}$  for the compound  
236 multi-modal feature. The concatenated feature vectors  $\mathbf{v}_{protein}$  and  $\mathbf{v}_{compound}$  were mapped  
237 onto the same latent space with a fixed dimension  $d$  by applying the fully connected layers  $f$   
238 and  $g$ . From this, the similarity between the two vectors for the latent space was calculated.

$$\begin{aligned}\mathbf{v}_{protein} &= \text{concat}(\mathbf{a}_{protein}, \mathbf{N}_{protein}), \\ \mathbf{v}_{compound} &= \text{concat}(\mathbf{b}_{compound}, \mathbf{N}_{compound}), \\ (x_1, x_2, \dots, x_d) &= f(\mathbf{v}_{protein}), \\ (y_1, y_2, \dots, y_d) &= g(\mathbf{v}_{compound}), \\ \text{output}_{integrated} &= h(x_1 \cdot y_1, x_2 \cdot y_2, \dots, x_d \cdot y_d).\end{aligned}$$

239 As described above, to handle data from different modalities such as proteins and  
240 compounds, we adopted a method of embedding data of different modalities into a common  
241 latent space. Defining the similarity in the obtained latent space enables the measurement of

242 the similarity between the data for different modalities. Visual semantic embedding (VSE) is  
243 a typical example of a method that handles data from different modalities and can associate  
244 images with text data in acquiring these multi-modal representations [22]. VSE was  
245 developed to generate captions from images (image captioning). The image feature and the  
246 sentence feature are linearly transformed and embedded into a common latent space.

247

### 248 ***Single-Modality Models***

249 To see the effect of integrating multi-modal features, two baseline models were constructed  
250 for the performance comparison. One was based on molecular structure data and used only  
251 amino acid sequence and chemical structure information, and the other was based on  
252 interaction network data and used only protein-protein interaction and compound-compound  
253 interaction information. The single-modality model based on molecular structure data,  
254 denoted the “single-modality model (molecular)”, is defined as follows:

$$(x_1, x_2, \dots, x_n) = f(\mathbf{a}_{protein}),$$

$$(y_1, y_2, \dots, y_n) = g(\mathbf{b}_{compound}),$$

$$\text{output}_{molecule} = h(x_1 \cdot y_1, x_2 \cdot y_2, \dots, x_n \cdot y_n),$$

255 and the single-modality model based on interaction network data, denoted the “single-  
256 modality model (network)”, is defined as follows:

$$(x_1, x_2, \dots, x_n) = f(\mathbf{N}_{protein}),$$

$$(y_1, y_2, \dots, y_n) = g(\mathbf{N}_{compound}),$$

$$\text{output}_{network} = h(x_1 \cdot y_1, x_2 \cdot y_2, \dots, x_n \cdot y_n).$$

257

### 258 ***Loss Function***

259 For the similarity output  $x$  of the model, the output value was restricted to the range 0 to 1 by  
260 the sigmoid function, and cross entropy was applied as the loss function  $L(\theta)$  to calculate the

261 training error.

262

### 263 ***Hyperparameter Optimization***

264 The hyperparameters, the number and size of the filters in the convolutional layers in the 1D-  
265 CNN, and the number of units in the fully connected output layers were optimized by the  
266 Bayesian optimization tool Optuna [23], which is an automatic hyperparameter optimization  
267 software framework particularly designed for machine learning. For the hyperparameter  
268 optimization, the validation dataset was obtained by dividing the training samples into a set  
269 for training and a set for validation.

270

### 271 ***Regularization***

272 Regularization is important for avoiding overfitting and improving the prediction accuracy in  
273 deep learning for complex model architectures with a large number of parameters.  
274 Regularization is especially important in our deep learning model, which integrates multiple  
275 datasets of different modalities; hence, we employed several regularization methods.

276 We employed batch normalization [24], which allowed us to use much higher  
277 learning rates and be less careful about initialization, after each convolutional layer. We also  
278 inserted dropout [25] after the fully connected layers. Furthermore, we added an L2  
279 regularization term to the training-loss function  $L(\theta)$ . When incorporating weight decay, the  
280 objective function to be optimized is as follows:

$$L(\theta) + \lambda \frac{1}{2} \sum_w \|w\|^2,$$

281 where  $w$  refers to the parameters of the entire model, and the second term of the above  
282 equation indicates taking the sum of the squared values of all the parameters and dividing by  
283 2.  $\lambda$  is a parameter that controls the strength of regularization. Adding this term to the

284 objective function has the effect of preventing the absolute value of the network weight from  
285 becoming too large, which helps prevent overfitting.

286

### 287 *Comparison with State-of-the-Art Existing Methods*

288 The prediction performance of the proposed models was compared with that of state-of-the-  
289 art deep learning methods based on molecular structure data and interaction network data.

290 The first method was based on a graph CNN for protein-compound prediction [10]. It  
291 employed a graph CNN for encoding chemical structures and a CNN for  $n$ -grams of amino  
292 acid sequences. The second method was NeoDTI [13], which demonstrated superior  
293 performance over other previous methods based on multiple-interaction-network data. We  
294 also compared our method with the traditional machine learning methods SVM and random  
295 forest [26] as the baseline prediction methods. These traditional methods require structured  
296 data as input. For the protein information, the 3-mer (3-residue) frequency in the amino acid  
297 sequence was used as the feature vector for 8,000 dimensions. For the compound  
298 information, an ECFP with a length of 1,024 and a radius of 2 was used. The radial basis  
299 function (RBF) was used as the kernel function of SVM, and all other parameters of SVM  
300 and random forest used the default values. In implementing these machine learning methods,  
301 scikit-learn (version 0.19.1) and chainer (version 5.0.0) were used.

302

### 303 *Datasets*

304 The protein-compound interaction data and compound-compound networks were retrieved  
305 from the database STITCH [19], and the protein-protein networks were retrieved from the  
306 database STRING [18].

307

### 308 *Protein-Compound Interaction Data*

309 Protein-compound interaction data can be obtained from the STITCH database [19]. STITCH  
310 contains data on the interaction of 430,000 compounds with 9.6 million proteins from 2,031  
311 species. The STITCH data sources consist of (1) structure-based prediction results, such as  
312 the genome context and co-expression; (2) high-throughput experimental data; (3) automatic  
313 text mining; and (4) information from existing databases. When a protein-compound dataset  
314 is downloaded from STITCH, a score based on the reliability is created for each of the above  
315 four items for each protein-compound pair. For the protein-compound interaction data used in  
316 this study (as a “positive” example), the threshold value for the reliability score of item (2)  
317 was set to 700, and the data with a reliability score of 700 or higher were extracted from  
318 STITCH so that interologs were eliminated and the data were composed of only  
319 experimentally reliable interactions; the data that did not meet this threshold were removed.  
320 For the STITCH data, interactions with a confidence score of 700 or more were determined  
321 based on the criterion that they were at least highly reliable [27]. Of the combinations of  
322 proteins and compounds, only pairs not stored in the STITCH database were taken as  
323 “negative” examples. In general, protein-compound pairs that are not stored in STITCH have  
324 very low confidence, with a score of 150 or less for their interaction [28], so these are  
325 considered to be non-interacting negative examples. The ratio of the positive and negative  
326 examples was 1 to 2.

327

### 328 ***Protein-Protein Interaction Data***

329 The protein-protein interaction information was obtained from the STRING database [18],  
330 which contains data for protein-protein interactions covering 24.6 million proteins from 5,090  
331 species. The STRING data sources consist of (1) experimental data; (2) pathway databases;  
332 (3) automatic text mining; (4) co-expression information; (5) neighbouring gene information;  
333 (6) gene fusion information; and (7) co-occurrence-based information. In particular, item (1)

334 is interaction data obtained from actual experiments, which include biochemical, biophysical,  
335 and genetic experiments. These are extracted from databases organized by the BioGRID  
336 database [29] and the IMEx consortium [30]. When the protein-protein interaction data from  
337 STRING were downloaded, a score based on the reliability was created for each of the above  
338 seven items for each protein-protein pair. Regarding the protein-protein interaction network,  
339 the threshold value for the reliability score of item (1) was set to 150. Data that did not satisfy  
340 this criterion were removed.

341

#### 342 ***Compound-Compound Interaction Data***

343 The compound-compound interaction data were also obtained from the STITCH database.  
344 The compound-compound interaction data in STITCH are based on (1) the chemical  
345 reactions obtained from the pathway databases; (2) structural similarity; (3) association with  
346 previous literature; and (4) correspondence between the compounds based on molecular  
347 activity similarity. For the similarity of the molecular activities in item (4), the activity data  
348 obtained by screening the model cell line NCI60 were used. When the compound-compound  
349 interaction data were downloaded from STITCH, a score based on the reliability was created  
350 for each of the above four items for each compound pair. For the compound-compound  
351 interaction data used in this study, the threshold value for the reliability score in item (4) was  
352 set to 150. Data that did not satisfy this criterion were removed.

353

#### 354 ***Construction of the Baseline, Unseen Compound-Test, and Hard Datasets for Evaluation***

355 From the STITCH and STRING databases, a total of 22,881 protein-compound interactions,  
356 175,452 protein-protein interactions and 69,231 compound-compound interactions were  
357 downloaded. Using the downloaded dataset in which the protein-protein interaction,  
358 compound-compound interaction and protein-compound interaction data were all available,

359 the three types of datasets below were constructed to perform five-fold cross validation.

360 In typical  $k$ -fold cross validation, all positive and negative examples are randomly  
361 split into  $k$  folds. One of them is used as a test sample, and the remaining  $k-1$  are used as  
362 training samples; then, the  $k$  results obtained are averaged. We call the cross-validation  
363 dataset the *baseline dataset*.

364 In this study, as more difficult and more practical tasks, we constructed two more  
365 cross-validation datasets, called the *unseen compound-test dataset* and the *hard dataset*. In  
366 the unseen compound-test dataset, we split the data into  $k$  folds so that none of the folds  
367 contain the same compounds as the others. In the unseen compound-test dataset, the  
368 compounds in the test sample do not appear in the training sample. In other words, the  
369 interaction of new (unseen) candidate compounds with the target proteins must be accurately  
370 predicted. In the hard dataset, we split the data into  $k$  folds so that none of the folds contain  
371 the same proteins and compounds as the others. In the hard dataset, neither the proteins nor  
372 the compounds in the test sample appear in the training sample. In other words, interactions  
373 in which neither the proteins nor the compounds are found in the training sample must be  
374 accurately predicted.

375

## 376 **Results**

377 The following measures were used for the performance evaluation criteria: AUROC (area  
378 under the receiver operating characteristic curve), AUPRC (area under the precision-recall  
379 curve), F-measure, and accuracy.

$$\text{F-measure} = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}},$$

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN},$$

380 where  $TP$  is the number of true positives,  $TN$  is the number of true negatives,  $FP$  is the



381 number of false positives,  $FN$  is the number of false negatives, Recall is defined by

382  $TP/(TP+FN)$ , and Precision is defined by  $TP/(TP+FP)$ .

383

### 384 *Effectiveness of Integrating Molecular Structure Data and Interaction Network Data*

385 The performance of our three models was evaluated to determine the effectiveness of

386 integrating the molecular structure data and the interaction network data. The results on the

387 three datasets are shown in Tables 1-3. In the tables, the mean and standard deviation (SD)

388 for the five folds are shown. Furthermore, the symbol “\*” indicates that there was a

389 significant difference in the Wilcoxon signed-rank test, with p-value  $p < 0.05$ , in comparison

390 with the integrated model.

391

392 **Table 1.** Performance comparison of three proposed models with existing methods on the

393 baseline dataset.

	AUROC	AUPRC	F-measure	Accuracy
Integrated model (molecular+network)	<b>0.972±0.004</b>	<b>0.954±0.005</b>	<b>0.900±0.006</b>	<b>0.933±0.004</b>
Single-modality model (molecular)	0.956±0.004*	0.927±0.006*	0.868±0.009*	0.911±0.006*
Single-modality model (network)	0.947±0.008*	0.920±0.010*	0.853±0.015*	0.904±0.009*
Graph CNN [10]	0.917±0.006*	0.850±0.006*	0.794±0.014*	0.864±0.008*
NeoDTI [13]	0.956±0.005*	0.905±0.016*	0.872±0.006*	0.917±0.004*
SVM	0.805±0.009*	0.651±0.012*	0.743±0.012*	0.837±0.006*
Random forest	0.873±0.009*	0.767±0.015*	0.837±0.012*	0.895±0.007*

394

395 **Table 2.** Performance comparison on the unseen compound-test dataset.

	AUROC	AUPRC	F-measure	Accuracy
Integrated model (molecular+network)	<b>0.890±0.039</b>	<b>0.842±0.050</b>	<b>0.727±0.085</b>	<b>0.843±0.038</b>
Single-modality model (molecular)	0.869±0.027	0.786±0.023*	0.657±0.053	0.802±0.017
Single-modality model (network)	0.831±0.053	0.759±0.055*	0.661±0.073*	0.809±0.030*
Graph CNN [10]	0.804±0.037*	0.679±0.031*	0.637±0.027	0.773±0.009*
NeoDTI [13]	0.823±0.067	0.773±0.064*	0.621±0.062*	0.805±0.024*
SVM	0.765±0.020*	0.603±0.029*	0.689±0.029	0.810±0.016
Random forest	0.770±0.023*	0.635±0.026*	0.697±0.036	0.828±0.014

396

397 **Table 3.** Performance comparison on the hard dataset.

	AUROC	AUPRC	F-measure	Accuracy
Integrated model (molecular+network)	<b>0.882±0.035</b>	<b>0.834±0.041</b>	<b>0.714±0.064</b>	<b>0.836±0.030</b>
Single-modality model (molecular)	0.851±0.023	0.770±0.023*	0.662±0.038*	0.806±0.020*
Single-modality model (network)	0.780±0.051*	0.706±0.040*	0.601±0.057*	0.784±0.023*
Graph CNN [10]	0.707±0.038*	0.563±0.083*	0.427±0.132*	0.719±0.043*
NeoDTI [13]	0.790±0.039*	0.715±0.046*	0.297±0.084*	0.719±0.018*
SVM	0.652±0.019*	0.500±0.023*	0.481±0.044*	0.755±0.012*
Random forest	0.605±0.033*	0.452±0.046*	0.364±0.075*	0.728±0.026*

398

399 Compared with the two single-modality models, the integrated model significantly

400 improved the prediction accuracy in all evaluation measures. For example, in terms of

401 AUPRC, which is a more informative evaluation index in a dataset that is imbalanced  
402 between positive and negative samples, the integrated model showed significant  
403 improvements of 3.0%, 7.1% and 8.3% over the single-modality model (molecular) and  
404 3.7%, 10.9% and 18.1% over the single-modality model (network) in the baseline dataset, the  
405 unseen compound-test dataset and the hard dataset, respectively. This demonstrates that  
406 integrating multiple heterogeneous interactome data with molecular structure data brought a  
407 synergistic effect in improving the accuracy of protein-compound interaction prediction.

408

#### 409 *Performance Comparison with Other Existing Methods*

410 The prediction performance of our three models was compared with that of state-of-the-art  
411 deep learning methods and traditional machine learning methods based on molecular  
412 structure data and interaction network data. The results on the three datasets are shown in  
413 Tables 1-3.

414 The integrated model yielded superior prediction performance compared with the  
415 other existing methods. In the baseline dataset, the integrated model achieved significant  
416 improvements compared with the graph CNN-based method [10], NeoDTI [13] and the  
417 traditional machine learning methods SVM and random forest (Table 1). In fact, the  
418 Wilcoxon signed-rank test [31] verification showed that the performance difference was  
419 statistically significant, with a p-value  $p < 0.05$ , and hence proved the superiority of the  
420 integrated model.

421 In the unseen compound-test dataset and the hard dataset, a more remarkable  
422 difference in the performance of the integrated model was confirmed. We compared the  
423 integrated model with the graph CNN-based method and NeoDTI in terms of AUROC,  
424 AUPRC and F-measure. The integrated model greatly outperformed the others, with  
425 significant improvements (10.7% in terms of AUROC, 24.0% in terms of AUPRC and 14.1%

426 in terms of F-measure on the unseen compound-test dataset, and 24.8% in terms of AUROC,  
427 48.1% in terms of AUPRC and 67.2% in terms of F-measure on the hard dataset) over the  
428 graph CNN-based method. In comparison with NeoDTI, significant improvements were also  
429 confirmed: 8.1% in terms of AUROC, 8.9% in terms of AUPRC and 17.1% in terms of F-  
430 measure on the unseen compound-test dataset, and 11.6% in terms of AUROC, 16.6% in  
431 terms of AUPRC and 140.4% in terms of F-measure on the hard dataset. Based on the above  
432 results, the integrated model can predict protein-compound interactions with stable accuracy,  
433 regardless of the difficulty of the dataset and the types of proteins and compounds that make  
434 up the test data, compared to other existing methods. This is due to the integrated model  
435 using features based on sequence information and compound structure information and  
436 features obtained from the interaction network as well as the effect of using the element-wise  
437 product of the protein feature vector and the compound feature vector in the output layer.

438         The single-modality model also yielded superior prediction performance compared  
439 with the existing methods using the same-modality input data. The graph CNN-based  
440 prediction method [10] obtains a compound feature vector by converting the chemical  
441 structure into a graph and applying it to the graph CNN, and it obtains a protein feature vector  
442 by splitting the amino acid sequence into  $n$ -grams and applying it to the CNN. Therefore, the  
443 graph CNN-based method can be defined as having the same molecular structure data-based  
444 prediction model as the single-modality model (molecular). In the baseline dataset, the  
445 unseen compound-test dataset and the hard dataset, the single-modality model (molecular)  
446 outperformed the graph CNN-based prediction method. For example, in the hard dataset, the  
447 single-modality model (molecular) achieved an improvement of 20.4% in terms of AUROC,  
448 36.8% in terms of AUPRC and 55.0% in terms of F-measure on the hard dataset over the  
449 graph CNN-based method (Table 3). From this result, in protein-compound interaction  
450 prediction, it is sufficient to use the ECFP as a feature representation for the compound

451 structure, compared with the deep learning method in which the compound structure is  
452 converted into a graph structure and a graph CNN is applied.

453 NeoDTI takes protein-protein interaction and compound-compound interaction  
454 information as input and predicts whether an edge is drawn between the compound and  
455 protein nodes by learning to reconstruct the network. Therefore, NeoDTI can be defined as an  
456 interaction network-based prediction model, which is the same as the single-modality model  
457 (network). The difference is that the single-modality model (network) first uses unsupervised  
458 deep learning (node2vec) to automatically learn feature representations for nodes in the given  
459 heterogeneous interaction networks and then applies supervised learning to predict protein-  
460 compound interactions based on the learned features, while NetoDTI simultaneously learns  
461 the feature representations of nodes and protein-compound interactions in a supervised  
462 manner. In the three datasets, the prediction performance of the single-modality model  
463 (network) was comparable to that of NetoDTI.

464

## 465 **Discussion**

466 To interpret the accuracy improvement obtained by integrating multiple interactome data  
467 with molecular structure data, which was shown in the previous section, we analysed whether  
468 the protein-protein interaction captured a different perspective than amino acid sequence  
469 homology and whether the compound-compound interaction captured a different perspective  
470 than chemical structure similarity. More concretely, we investigated the relationship between  
471 the amino acid sequence homology and the similarity of proteins in the protein-protein  
472 interaction network as well as the relationship between the chemical structure similarity and  
473 the similarity in the compound-compound interaction network.

474 For every pair of proteins in the dataset used in the experiments, the amino acid  
475 sequence similarity was calculated using DIAMOND, and the cosine similarity between two

476 vectors of the pair output by node2vec using the protein-protein interaction network was  
477 calculated. All of the protein pairs were plotted with the amino acid sequence similarity on  
478 the x-axis and the cosine similarity in the protein-protein interaction network on the y-axis.  
479 The scatter plot is shown in Figure 3 (top). Similarly, for every pair of compounds, the  
480 Jaccard coefficient of the ECFPs of the two compounds and the cosine similarity between the  
481 two vectors output by node2vec using a compound-compound interaction network were  
482 calculated. All of the compound pairs were plotted with the Jaccard coefficient on the x-axis  
483 and the cosine similarity in the compound-compound interaction network on the y-axis, as  
484 shown in Figure 3 (bottom). In both scatter plots, no clear correlation was observed. In fact,  
485 the correlation coefficients for each scatter plot were 0.186 and 0.199, respectively. In other  
486 words, it was confirmed that the amino acid sequence similarity and the similarity in the  
487 protein-protein interaction network were not proportional. Similarly, it was confirmed that  
488 the chemical structure similarity and the similarity in the compound-compound interaction  
489 network were not proportional. Therefore, we concluded that the protein-protein interaction  
490 network captured a different perspective than the amino acid sequence homology and  
491 compensated for it. The compound-compound interactions captured a different perspective  
492 than the chemical structure similarity and compensated for it.

493 For example, the protein “5-hydroxytryptamine (serotonin) receptor 6, G protein-  
494 coupled (HTR6)” and the compound “Mesulergine” in the test sample in the “hard dataset”  
495 have a positive interaction [32], and our model succeeded in correctly predicting it. However,  
496 the single-modality model (molecular) and graph CNN-based method failed to predict the  
497 positive interaction; that is, both predicted that the pair would not interact. The most similar  
498 protein-compound pair in the training sample to the pair HTR6 and Mesulergine was the  
499 protein “adrenoceptor alpha 2A (ADRA2A)” and the compound “Pergolide” [33]. The  
500 protein ADRA2A and the compound Pergolide have a positive interaction in the training

501 sample. The sequence similarity score between HTR6 and ADRA2A is rather low at 100.5,  
502 but the similarity of the two proteins in the protein-protein interaction network is relatively  
503 high at 0.805. A part of the protein-protein interaction network around HTR6 and ADRA2A  
504 is displayed in Figure 4 (left). Similarly, the Jaccard coefficient of the ECFPs between  
505 Mesulergine and Pergolide is relatively low 0.273 (in general, compound pairs with a Jaccard  
506 coefficient of ECFPs below 0.25 are considered not to have chemically similar structures  
507 [34]), but the cosine similarity of the two compounds in the compound-compound interaction  
508 network is high at 0.735. A part of the compound-compound interaction network around  
509 Mesulergine and Pergolide is displayed in Figure 4 (right).

510

## 511 **Conclusions**

512 This study aimed to improve the performance of predicting protein-compound interactions by  
513 integrating molecular structure data and interactome data. This was achieved by integrating  
514 multiple heterogeneous interactome data into predictions of protein-compound interactions.  
515 An end-to-end learning method was developed that combined a 1D-CNN for amino acid  
516 sequences, an ECFP representation for compounds, and feature representation learning with  
517 node2vec for protein-protein and compound-compound interaction networks. The proposed  
518 integrated model exhibited significant performance differences with respect to the accuracy  
519 measures in comparison to the current state-of-the-art deep learning methods. The  
520 performance improvement was verified by the Wilcoxon signed-rank test as being  
521 statistically significant. The results indicated that the proposed model was able to more  
522 accurately predict the protein-compound interactions even in the hard dataset, where neither  
523 the proteins nor the compounds in the test sample appear in the training sample.

524 An important future task is to integrate the gene regulatory network as additional  
525 interactome data to further improve protein-compound interaction prediction. A large number

526 of gene expression profiles for various tissues and cell lines are available in public databases,  
527 and gene regulatory networks can be effectively inferred from the gene expression profiles.

528

529

### 530 **List of Abbreviations**

531 SVM: Support Vector Machines

532 CNN: Convolutional Neural Network

533 ECFP: Extended-Connectivity Fingerprint

534 VSE: Visual Semantic Embedding

535 AUROC: Area Under the Receiver Operating characteristic Curve

536 AUPRC: Area Under the Precision-Recall Curve

537 SD: Standard Deviation

538

### 539 **Declarations**

540

#### 541 *Availability of Data and Materials*

542 The source code for the implementation of this deep learning method, along with the dataset  
543 for the performance evaluation, is available at [https://github.com/Njk-901aru/multi\\_DTI.git](https://github.com/Njk-901aru/multi_DTI.git).

544

#### 545 *Competing Interests*

546 The authors declare that they have no competing interests.

547

#### 548 *Funding*

549 This work was supported by a Grant-in-Aid for Scientific Research on Innovative Areas

550 “Frontier Research on Chemical Communications” [no. 17H06410] from the Ministry of



551 Education, Culture, Sports, Science and Technology of Japan, and a Grant-in-Aid for  
552 Scientific Research (A) (KAKENHI) [No. 18H04127] from the JSPS.

553

#### 554 *Authors' Contributions*

555 NW; Implemented the software, analysed data, and co-wrote the paper. YO; analysed data  
556 and compared with the existing methods. YS; designed and supervised the research, analysed  
557 data, and co-wrote the paper. All authors read and approved the final manuscript.

558

#### 559 *Acknowledgements*

560 Not applicable.

561

#### 562 **References**

- 563 1. Nagamine N, Sakakibara Y (2007) Statistical prediction of protein–chemical  
564 interactions based on chemical structure and mass spectrometry data. *Bioinformatics*  
565 23:2004-2012.
- 566 2. Nagamine N, Shirakawa T, Minato Y, Torii K, Kobayashi H, Imoto M, Sakakibara Y  
567 (2009) Integrating statistical predictions and experimental verifications for enhancing  
568 protein-chemical interaction predictions in virtual screening. *PLoS Comput Biol*  
569 5:e1000397.
- 570 3. Sakakibara Y, Hachiya T, Uchida M, Nagamine N, Sugawara Y, Yokota M,  
571 Nakamura M, Popenoerf K, Komori T, Sato K (2012) COPICAT: a software system  
572 for predicting interactions between proteins and chemical compounds. *Bioinformatics*  
573 28:745-746.
- 574 4. Wishart DS, Knox C, Guo AC, Cheng D, Shrivastava S, Tzur D, Gautam B, Hassanali  
575 M (2008) DrugBank: a knowledgebase for drugs, drug actions and drug targets.

- 576 Nucleic Acids Res 36:D901-D906.
- 577 5. Kim S, Thiessen PA, Bolton EE, Chen J, Fu G, Gindulyte A, Han L, He J, He S,  
578 Shoemaker BA, Wang J, Yu B, Zhang J, Bryant SH (2016) PubChem substance and  
579 compound databases. *Nucleic Acids Res* 44:D1202-D1213.
- 580 6. LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521:436-444.
- 581 7. Tian K, Shao M, Wang Y, Guan J, Zhou S (2016) Boosting compound-protein  
582 interaction prediction by deep learning. *Methods* 110:64-72.
- 583 8. Öztürk H, Özgür A, Ozkirimli E (2018) DeepDTA: deep drug-target binding affinity  
584 prediction. *Bioinformatics (Oxford, England)* 34:i821-i829.
- 585 9. Lee I, Keum J, Nam H (2019) DeepConv-DTI: prediction of drug-target interactions  
586 via deep learning with convolution on protein sequences. *PLoS Comput Biol*  
587 15:e1007129.
- 588 10. Tsubaki M, Tomii K, Sese J (2019) Compound–protein interaction prediction with  
589 end-to-end learning of neural networks for graphs and sequences. *Bioinformatics*  
590 35:309-318.
- 591 11. Zitnik M, Agrawal M, Leskovec J (2018) Modeling polypharmacy side effects with  
592 graph convolutional networks. *Bioinformatics (Oxford, England)* 34:i457-i466.
- 593 12. Luo Y, Zhao X, Zhou J, Yang J, Zhang Y, Kuang W, Peng J, Chen L, Zeng J (2017)  
594 A network integration approach for drug-target interaction prediction and  
595 computational drug repositioning from heterogeneous information. *Nat Commun*  
596 8:573.
- 597 13. Wan F, Hong L, Xiao A, Jiang T, Zeng J (2019) NeoDTI: neural integration of  
598 neighbor information from a heterogeneous network for discovering new drug–target  
599 interactions. *Bioinformatics* 35:104-111.
- 600 14. Hosur R, Peng J, Vinayagam A, Stelzl U, Xu J, Perrimon N, Bienkowska J, Berger B

- 601 (2012) A computational framework for boosting confidence in high-throughput  
602 protein-protein interaction datasets. *Genome Biol* 13:R76.
- 603 15. Chen M, Ju CJT, Zhou G, Chen X, Zhang T, Chang KW, Zaniolo C, Wang W (2019)  
604 Multifaceted protein-protein interaction prediction based on Siamese residual RCNN.  
605 *Bioinformatics* 35:i305-i314.
- 606 16. Chen L, Lu J, Zhang J, Feng KR, Zheng MY, Cai YD (2013) Predicting chemical  
607 toxicity effects based on chemical-chemical interactions. *PLoS One* 8:e56517.
- 608 17. Kwon S, Yoon S (2019) End-to-end representation learning for chemical-chemical  
609 interaction prediction. *IEEE/ACM Trans Comput Biol Bioinform* 16:1436-1447.
- 610 18. Szklarczyk D, Morris JH, Cook H, Kuhn M, Wyder S, Simonovic M, Santos A,  
611 Doncheva NT, Roth A, Bork P, Jensen LJ, von Mering C (2017) The STRING  
612 database in 2017: quality-controlled protein-protein association networks, made  
613 broadly accessible. *Nucleic Acids Res* 45:D362-D368.
- 614 19. Kuhn M, von Mering C, Campillos M, Jensen LJ, Bork P (2016) STITCH: interaction  
615 networks of chemicals and proteins. *Nucleic Acids Res* 36:D684-D688.
- 616 20. Grover A, Leskovec J (2016) node2vec: scalable feature learning for networks. In:  
617 *Proceedings of KDD '16 (22nd ACM SIGKDD international conference on*  
618 *knowledge discovery and data mining)*. ACM, New York, NY, USA, p 855-864.
- 619 21. Rogers D, Hahn M (2010) Extended-connectivity fingerprints. *J Chem Inf Model*  
620 50:742-754.
- 621 22. Kiros R, Salakhutdinov R, Zemel RS (2014) Unifying visual-semantic embeddings  
622 with multimodal neural language models. *arXiv preprint arXiv:1411.2539*.
- 623 23. Akiba T, Sano S, Yanase T, Ohta T, Koyama M (2019) Optuna: a next-generation  
624 hyperparameter optimization framework. In: *Proceedings of 25th ACM SIGKDD*  
625 *international conference on knowledge discovery & data mining*. ACM, New York,

- 626 NY, USA, p 2623–2631.
- 627 24. Ioffe S, Szegedy C (2015) Batch normalization: accelerating deep network training by  
628 reducing internal covariate shift. arXiv preprint arXiv:150203167.
- 629 25. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout:  
630 a simple way to prevent neural networks from overfitting. *J Mach Learn Res* 15:1929-  
631 1958.
- 632 26. Bishop C (2006) *Pattern recognition and machine learning*. Springer-Verlag, Berlin,  
633 Heidelberg.
- 634 27. Liu R, Hameed MDMA, Kumar K, Yu X, Wallqvist A, Reifman J (2017) Data-driven  
635 prediction of adverse drug reactions induced by drug-drug interactions. *BMC*  
636 *Pharmacol Toxicol* 18:44.
- 637 28. Kuhn M, Szklarczyk D, Franceschini A, von Mering C, Jensen LJ, Bork P (2011)  
638 STITCH 3: zooming in on protein-chemical interactions. *Nucleic Acids Res* 40:D876-  
639 D880.
- 640 29. Chatr-Aryamontri A, Oughtred R, Boucher L, Rust J, Chang C, Kolas NK, O'Donnell  
641 L, Oster S, Theesfeld C, Sellam A, Stark C, Breitkreutz BJ, Dolinski K, Tyers M  
642 (2019) The BioGRID interaction database: 2017 update. *Nucleic Acids Res* 45:D369-  
643 D379.
- 644 30. Orchard S, Kerrien S, Abbani S, Aranda B, Bhate J, Bidwell S, Bridge A, Briganti L,  
645 Brinkman FSL, Cesareni G, Chatr-Aryamontri A, Chautard E, Chen C, Dumousseau  
646 M, Goll J, Hancock REW, Hannick LI, Jurisica I, Khadake J, Lynn DJ, Mahadevan U,  
647 Perfetto L, Raghunath A, Ricard-Blum S, Roechert B, Salwinski L, Stümpflen V,  
648 Tyers M, Uetz P, Xenarios I, Hermjakob H (2012) Protein interaction data curation:  
649 the International Molecular Exchange (IMEx) consortium. *Nat Methods* 9:345-350.
- 650 31. Wilcoxon F (1945) Individual comparisons by ranking methods. *Biom Bull* 1:80-83.

- 651 32. Krogsgaard-Larsen N, Jensen AA, Schroder TJ, Christoffersen CT, Kehler J (2014)  
652 Novel aza-analogous ergoline derived scaffolds as potent serotonin 5-HT<sub>6</sub> and  
653 dopamine D<sub>2</sub> receptor ligands. *J Med Chem* 57:5823-5828.
- 654 33. Millan MJ, Maiorini L, Cussac D, Audinot V, Boutin JA, Newman-Tancredi A (2002)  
655 Differential actions of antiparkinson agents at multiple classes of monoaminergic  
656 receptor. I. A multivariate analysis of the binding profiles of 14 drugs at 21 native and  
657 cloned human receptor subtypes. *J Pharmacol Exp Ther* 303:791-804.
- 658 34. Childs-Disney JL, Tran T, Vummidi BR, Velagapudi SP, Haniff HS, Matsumoto Y,  
659 Crynen G, Southern MR, Biswas A, Wang ZF, Tellinghuisen TL, Disney MD (2018)  
660 A massively parallel selection of small molecule-RNA motif binding partners informs  
661 design of an antiviral from sequence. *Chem* 4:2384-2404.

662

### 663 **Figure legends**

664 **Figure 1.** Deep learning architecture that integrates molecular structure data and interactome  
665 data to predict protein-compound interactions. It integrates graph-based and sequence-based  
666 representations for the target protein and compound. The amino acid sequence of the protein  
667 input was embedded into a one-hot vector of 20 dimensions in height. The ECFP  
668 representation of the compound input was embedded into a 1024-dimensional vector. The  
669 feature vectors were also extracted from the protein-protein and compound-compound  
670 interaction network using node2vec, a feature representation learning method for graphs.  
671 These feature vectors were combined as a protein vector and a compound vector. The  
672 interaction was predicted in the output unit.

673

674 **Figure 2.** The output layer architecture. The integrated model predicts the protein-compound  
675 interactions by embedding the protein and compound data from different modalities into a

676 common latent space. The feature vectors for the proteins and compounds are mapped onto  
677 the same latent space by applying a fully connected layer. Then, their similarity in the latent  
678 space is calculated with an element-wise product calculation followed by a fully connected  
679 layer.

680

681 **Figure 3.** (Top) Relationship between the amino acid sequence similarity and the similarity  
682 in protein-protein interaction network. (Bottom) Relationship between the chemical-structure  
683 similarity and the similarity in compound-compound interaction network. The amino acid  
684 sequence similarity was calculated using DIAMOND, and the chemical structure similarity  
685 was calculated as the Jaccard coefficient of the ECFPs of the two compounds. The correlation  
686 coefficients are 0.186 and 0.199, respectively.

687

688 **Figure 4.** (Left) Part of the protein-protein interaction network around ABL1 and YES1.  
689 (Right) Part of the compound-compound interaction network around Crizotinib and Ceritinib.

690

**Figure 1**

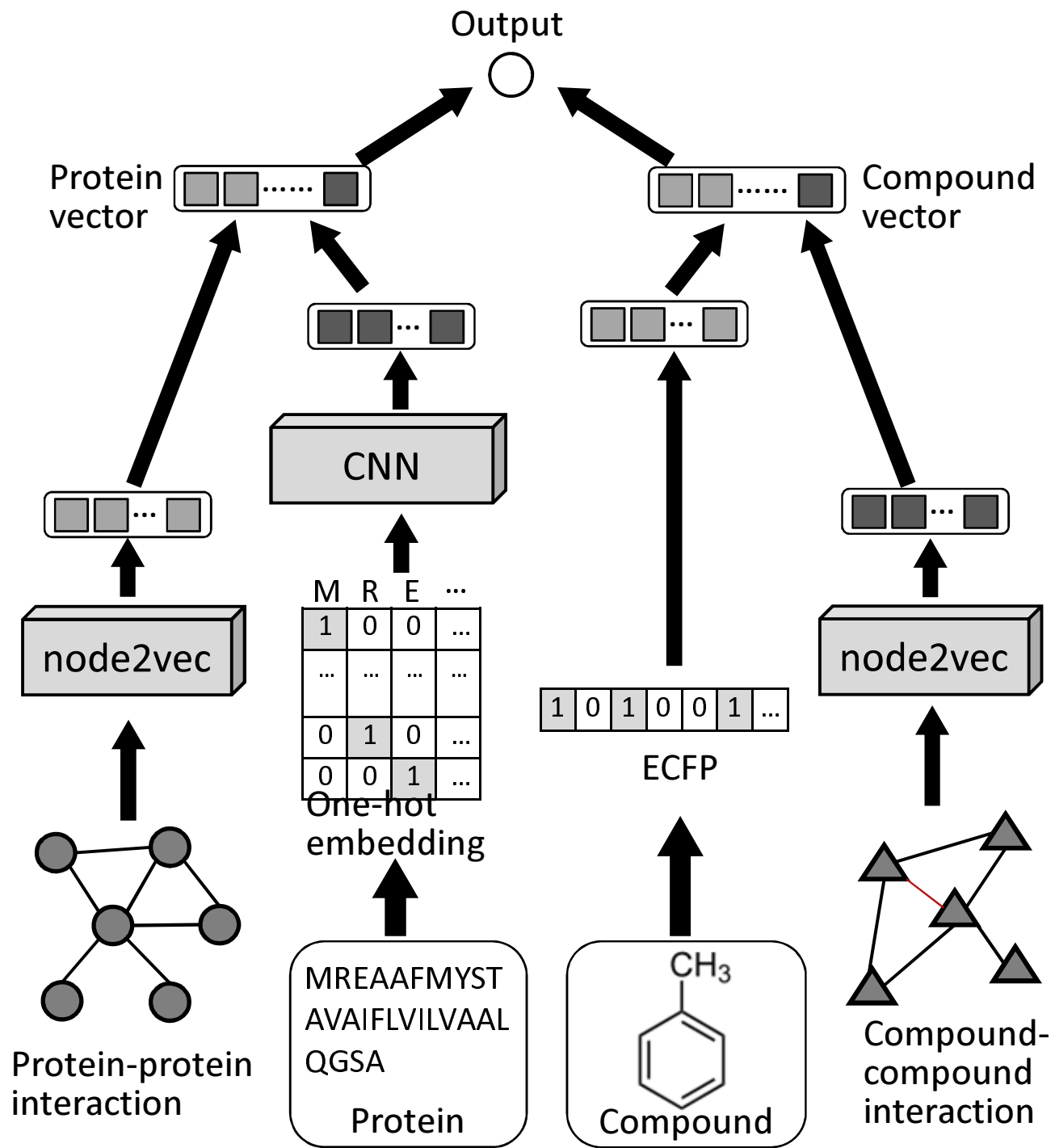
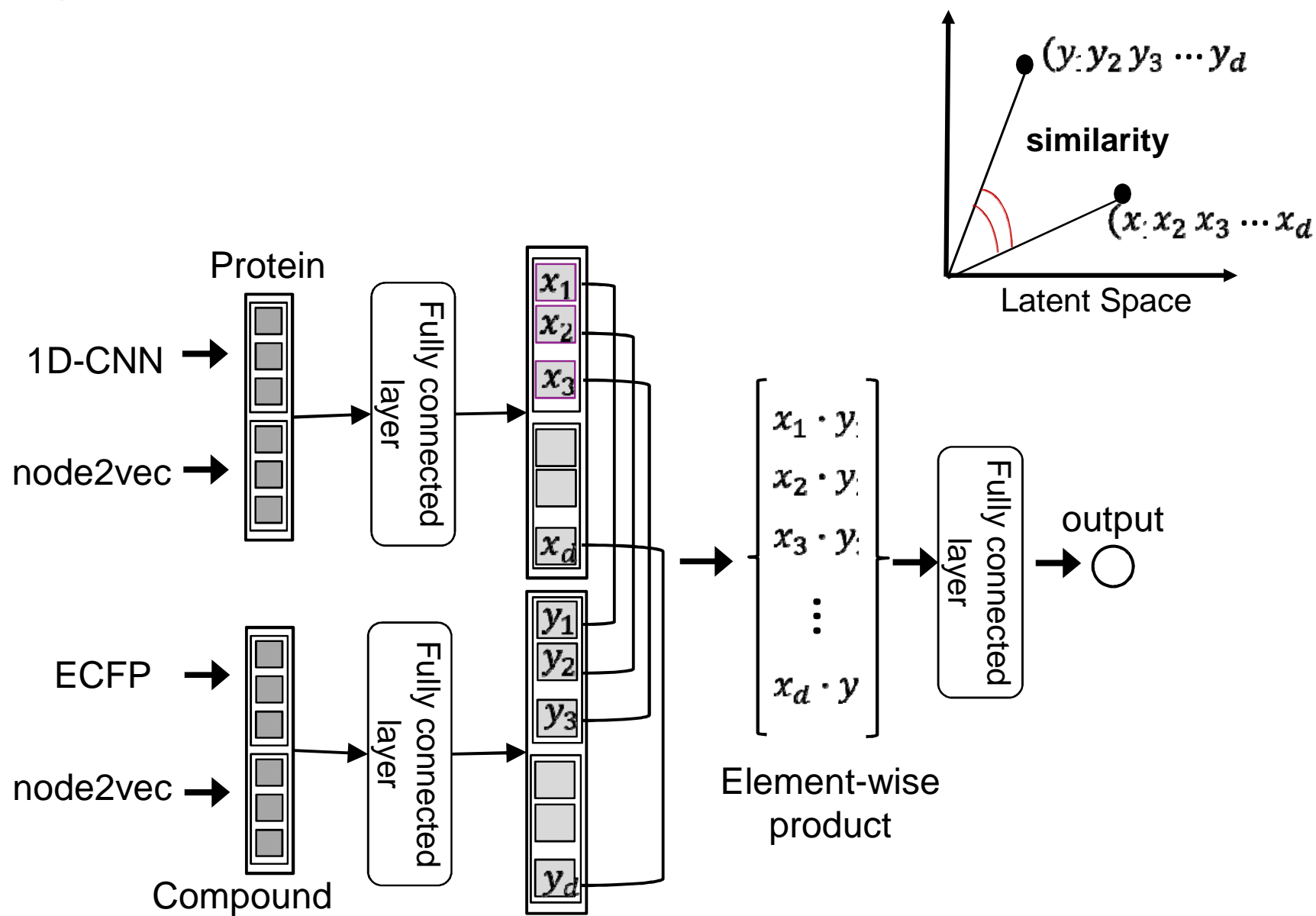
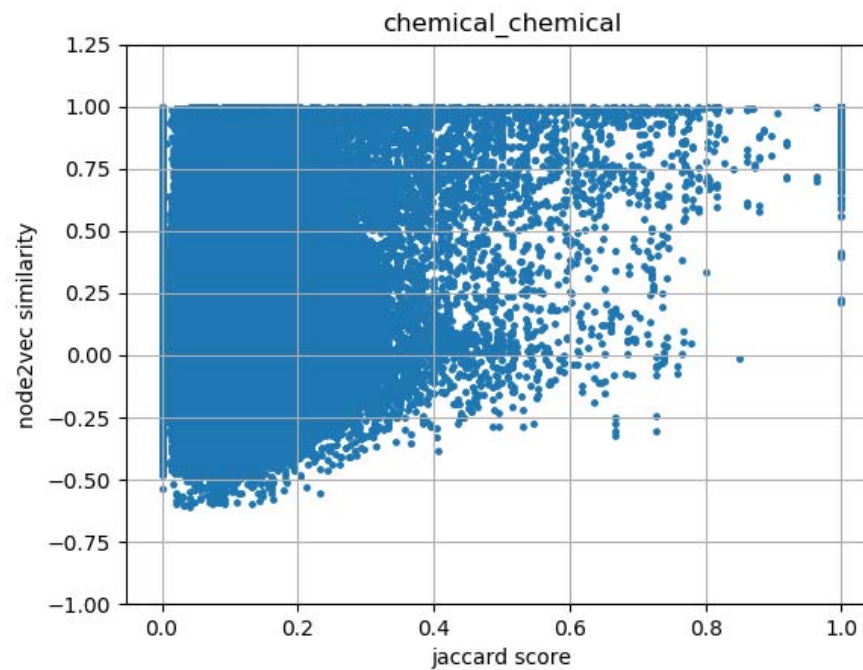
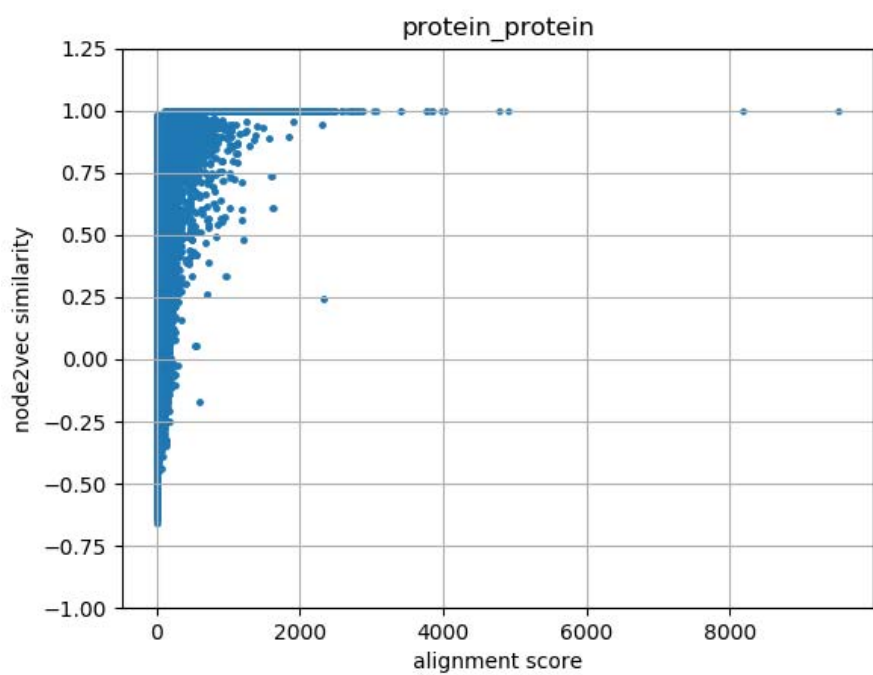


Figure 2





# Figure 3



**Figure 4**

