

# Systematic evaluation of machine learning methods for identifying human–pathogen protein–protein interactions

Huaming Chen, Fuyi Li, Lei Wang, Yaochu Jin, Chi-Hung Chi, Lukasz Kurgan, Jiangning Song and Jun Shen

Corresponding authors: Jun Shen, Faculty of Engineering and Information Science, School of Computing and Information Technology, University of Wollongong, NSW 2522, Australia. Tel.: +61-2-4221-3873. E-mail: jshen@uow.edu.au; Lukasz Kurgan, Department of Computer Science, Virginia Commonwealth University, 401 West Main Street, Room E4225, Richmond, VA 23284, USA. Tel.: +1-804-827-3986. E-mail: lkurgan@vcu.edu; Jiangning Song, Biomedicine Discovery Institute, Department of Biochemistry and Molecular Biology, and Monash Centre of Data Science, Monash University, Melbourne, VIC 3800, Australia. Tel.: +61-3-9902-9304. E-mail: jiangning.song@monash.edu

## Abstract

In recent years, high-throughput experimental techniques have significantly enhanced the accuracy and coverage of protein–protein interaction identification, including human–pathogen protein–protein interactions (HP-PPIs). Despite this progress, experimental methods are, in general, expensive in terms of both time and labour costs, especially considering that there are enormous amounts of potential protein-interacting partners. Developing computational methods to predict interactions between human and bacteria pathogen has thus become critical and meaningful, in both facilitating the detection of interactions and mining incomplete interaction maps. **In this paper, we present a systematic evaluation of machine learning-based computational methods for human–bacterium protein–protein interactions (HB-PPIs).** We first reviewed a vast number of publicly available databases of HP-PPIs and then critically evaluate the availability of these databases. Benefitting from its well-structured nature, we subsequently preprocess the data and identified six bacterium

**Huaming Chen** received his BEng and MEng degrees from Lanzhou University, China, in 2012 and 2015, respectively. He is currently a PhD candidate in the University of Wollongong. His research interests are bioinformatics, machine learning and neural network.

**Fuyi Li** received his BEng and MEng degrees from Northwest A&F University, China. He is currently a PhD candidate in the Department of Biochemistry and Molecular Biology and Biomedicine Discovery Institute, Monash University, Australia. His research interests are bioinformatics, computational biology, machine learning and data mining.

**Lei Wang** received his BEng and MEng degrees from Southeast University, China, in 1996 and 1999, respectively, and his PhD degree from Nanyang Technological University, Singapore, in 2004. He is currently an Associate Professor with the Faculty of Informatics, University of Wollongong.

**Yaochu Jin** received his BSc, MSc and PhD degrees from Zhejiang University, Hangzhou, China, in 1988, 1991 and 1996, respectively, and the Dr.-Ing. degree from Ruhr University Bochum, Bochum, Germany, in 2001. He is a Professor in the Department of Computer Science, University of Surrey, United Kingdom. His research interests include computational intelligence, computational systems biology and nature-inspired problem-solving.

**Chi-Hung Chi** received the PhD degree from Purdue University, USA. He is currently a Senior Principal Research Scientist of Data 61 in CSIRO (Commonwealth Scientific and Industrial Research Organization), Australia. His research areas include cybersecurity, behaviour modelling, knowledge graph, data engineering and analytics, cloud and service computing, social computing, Internet-of-Things and distributed computing.

**Lukasz Kurgan** is a Robert J. Mattauch Endowed Professor of Computer Science at the Virginia Commonwealth University. He is a Fellow of American Institute for Medical and Biological Engineering (AIMBE) and has published close to 150 peer-reviewed journal articles that focus on structural and functional characterization of proteins and small RNAs. More details about his research group are available at <http://biomine.cs.vcu.edu/>.

**Jiangning Song** is an associate professor and group leader in the Biomedicine Discovery Institute, Monash University, Australia. He is also affiliated with the Monash Centre for Data Science, Faculty of Information Technology, Monash University. His research interests include bioinformatics, computational biology, machine learning, data mining, and pattern recognition.

**Jun Shen** received the PhD degree from Southeast University, China, in 2001. He is currently an Associate Professor in the University of Wollongong, Wollongong, NSW, Australia. His expertise is on cloud computing and big data.

Submitted: 31 January 2020; Received (in revised form): 31 March 2020

© The Author(s) 2020. Published by Oxford University Press. All rights reserved. For Permissions, please email: journals.permissions@oup.com

pathogens that could be used to study bacterium subjects in which a human was the host. Additionally, we thoroughly reviewed the literature on 'host-pathogen interactions' whereby existing models were summarized that we used to jointly study the impact of different feature representation algorithms and evaluate the performance of existing machine learning computational models. Owing to the abundance of sequence information and the limited scale of other protein-related information, we adopted the primary protocol from the literature and dedicated our analysis to a comprehensive assessment of sequence information and machine learning models. A systematic evaluation of machine learning models and a wide range of feature representation algorithms based on sequence information are presented as a comparison survey towards the prediction performance evaluation of HB-PPIs.

**Key words:** bioinformatics; human-pathogen interactions; protein-protein interactions; systematic evaluation; sequential analysis; machine learning

## Introduction

Infectious diseases are predominantly caused by many pathogenic species, such as fungi, viruses, bacteria and so on. These infectious species actively interact with their hosts in a variety of ways, which place host-pathogen interactions (HPIs) in a complicated, but also critical, role in the study of infectious disease mechanisms. In most cases, the host-pathogen system is studied from different perspectives to further our understanding of infectious mechanisms [1]. A major approach is studying the interactions of inter-species proteins, in which one protein is from the host and the other is from the pathogen.

While protein interactions occur extraordinarily between human and bacterium pathogens, one of the earliest studies illustrated the importance of human-bacterium interactions (HBI) in relation to the symptoms caused by anthrax [2]. In this study, *Bacillus anthracis* was conclusively demonstrated as the primary cause of anthrax. Additional studies of *B. anthracis* were conducted, aimed at fully understanding the mechanisms of a complete protein interaction network between *B. anthracis* (the bacterium pathogen) and *Homo sapiens* (the host) [3, 4]. These studies encouraged researchers to study a broad range of infectious diseases by exploring human-bacterium protein-protein interactions (HB-PPIs).

However, the investigation of HBIs consumes lots of time, money and resources in determining the complete interaction network and understanding their mechanisms. Currently, investigations of the interactions between host and pathogens are still very limited. Even though large-scale biomedical technologies, such as yeast two-hybrid assay and the affinity purifications-mass spectrometry (AP-MS) method, have allowed us to detect interactions (positive or negative) in a faster and more accurate way, the amount of possible HB-PPIs is large. Other small-scale technologies, like nuclear magnetic resonance (NMR), are often labour-intensive and time-consuming. Thus, it is critical to formulate a computational model for the prediction of HB-PPIs.

Several reviews studied current computational approaches [5, 6] as well as researches on applying machine learning-based models to predict host-pathogen protein-protein interactions (HP-PPIs) [7–11]. In particular, how to deploy machine learning models as a generic approach in predicting novel HBIs based on sequence information is considered as an important category of research, which involves many challenges and opportunities. However, there is currently no comprehensive evaluation study that has focused on machine learning models as the primary computational method and further comparatively evaluated their corresponding performances across a wide range of HBI systems.

In this paper, we implemented an evaluation protocol based on literature reviews by first collecting HBI data from a wide range of host-pathogen databases. The systematic evaluation was subsequently achieved from two aspects. The first considered the application of feature representation algorithms to the protein data, while the other was related to different machine learning-based models.

The remainder of this paper is organized as follows. Section 2 summarizes the literature review from four different perspectives, including the review of host-pathogen interaction studies, the review of available host-pathogen interaction databases, the review of computational methods for host-pathogen protein-protein interaction predictions and the review of sequential representation algorithms and the machine learning-based methods for prediction. In Section 3, the materials collected for evaluation and the details of curated datasets are presented. Section 4 discusses the evaluation results in detail, and we conclude the paper in Section 5.

## Literature review

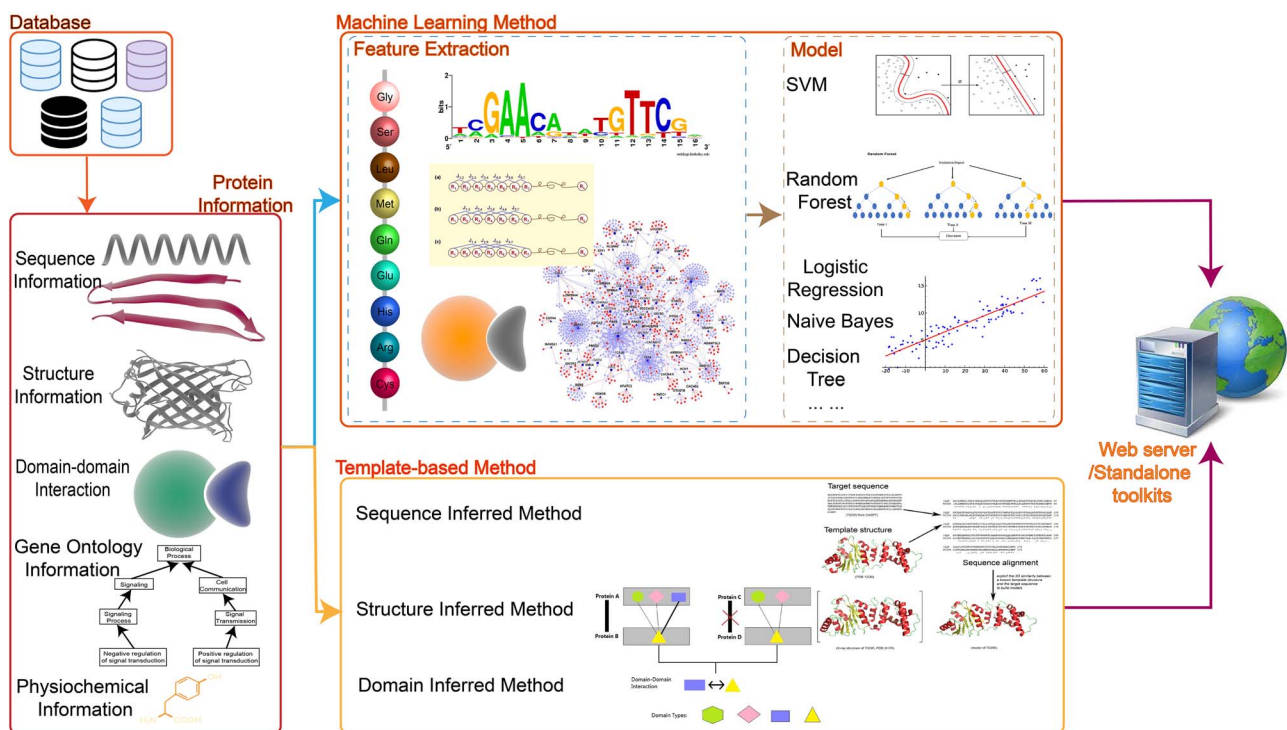
Although there has been a long history of research on PPI prediction, so far there are only a small number of publications that have focused on host-pathogen interaction reviews [5, 6, 12, 13]. A broad search has resulted in four major review papers, and Table 1 summarizes the reviews. The studies by [6, 12] have a wide coverage of HPIs, which include predictions as well as analyses, while the reviews by [5, 13] focused on the computational prediction of HPIs. These reviews aimed at describing the progress of HPIs, without anchors of naming pathogens, and they collectively reported on potential computational methods such as homology-based approaches, structure-based approaches, domain and motif interaction-based approaches and machine learning-based approaches. Furthermore, no systematic evaluation with details was implemented or reported in these reviews. Recently, [14] conducted a sequence-based predictors review; however, they focused on the prediction of protein-binding residues via single-sequence methods.

Adapted from these reviews, we subsequently collected all published predictors that focused on HB-PPIs and HP-PPIs, which are summarized in Table 2. The frameworks of the two different types of computational models for predicting HP-PPIs, including machine learning-based models and template-based models, are shown in Figure 1.

A template-based model utilizes different types of protein information to build the prediction model, including sequence information, structure information and domain information [15–17]. Template-based models use different protein

**Table 1.** Overview of the reviews of HP-PPIs

Review	Reviewed methods for HP-PPI	Reviewed database	Released dataset	Evaluation results of methods		
				Methods	Independent dataset	Measurement
[5] (2015)	Machine learning- and data mining-based approaches, homology-based approaches, structure-based approaches, domain and motif-based approaches	None	N/A	N/A	N/A	N/A
[6] (2015)	N/A	Web-based databases (HCVPro, PATRIC, HPIDB, PHISTO and so on)	N/A	N/A	N/A	N/A
[12] (2016)	Homology-based prediction, structure-based prediction, domain/motif interaction-based prediction, machine learning-based predictions of host-pathogen interactions	PATRIC, PIG, VirHostNet, HPIDB, VirusMINT and so on	N/A	N/A	N/A	N/A
[13] (2013)	Homology-based approaches, structure-based approach, domain and motif interaction-based approach, machine learning-based approach	VirusMINT, PHI-base, MINT, VirHostNet, BioGRID, IntAct, APID, PATRIC and so on	N/A	N/A	N/A	N/A

**Figure 1.** Flowcharts of template-based methods and machine learning-based methods. For each type of method, the key steps are summarized and visualized.

information to detect high score homology which might yield similar functions. However, template-based models may fail to predict whether the remote homology is interacting with known proteins. Another type of computational model is based on machine learning models. The protein information is first vectorized as the input to learn their inherent relationships automatically, which are thus used to build the model and

predict the interactions. Specifically, for PPIs, the relevant protein information can be sequence information, gene ontology information, domain information, gene expression information and interaction network information.

As indicated in Table 2, numerous feature representation algorithms for sequence information are incorporated with different machine learning models for predicting HP-PPIs [7, 8,

Table 2. Computational approaches for predicting HP-PPIs (sorted by published year)

Predictor	Pathogen	Data source	Training data		Independent data	Evaluation scheme	Protein information	Sequence representation algorithms	Algorithm/model	Stand-alone software/platform	Web server
			Positive pairs	Negative pairs							
[18]	Parasite	BIND, DIP, IntAct, Reactome	39,207 human, 18,412, 2643 <i>Plasmodium falciparum</i> intra-species interactions	N/A	N/A	Gene ontology	N/A	Bayesian statics	N/A	N/A	
[15]	Parasite	DIP	N/A	N/A	N/A	Sequence	PSSM	Remote homology detection	N/A	N/A	
[19]	Parasite	MINT, IntAct, Reactome, HPRD	1112	1136	N/A	Sequence	CTM variation	Random forest	N/A	N/A	
[7]	Virus	HPRD, MINT, BIND, DIP, IntAct, Reactome	1028	1:25, 1:50, 1:100 ratio of positive pairs	N/A	Domain, sequence, interaction network	k-mers	Support vector machine (linear kernel)	N/A	N/A	
[8]	Virus	I-MAP	500	500	N/A	Sequence	CTM variation	Support vector machine (RBF)	N/A	N/A	
[9]	Bacterium ( <i>B. anthracis</i> , <i>F. tularensis</i> , <i>Y. pestis</i> , <i>S. typhi</i> )	PHISTO	600 <i>B. anthracis</i> , 491 <i>F. tularensis</i> , 839 <i>Y. pestis</i> , 62 <i>S. typhi</i>	1:100 ratio of positive pairs	N/A	Sequence, gene ontology, gene expression, interaction network	k-mers	Multitask learning	<a href="http://www.cs.cmu.edu/~mkshirsa/i_smb2013_paper_r320.html">http://www.cs.cmu.edu/~mkshirsa/i_smb2013_paper_r320.html</a>	N/A	
[20]	Virus	RefSeq	3638	3638	Holding subcatalog PPI dataset	Independent test	Gene ontology	N/A	Transfer learning	N/A	N/A
[10]	Virus	IntAct	657	2910	N/A	Sequence, interaction network, tissue information, post-translational modifications	AAC, PAAC, PSSM	Ensemble learning	N/A	N/A	
[16]	Bacterium ( <i>Mycobacterium tuberculosis</i> )	N/A	-	N/A	N/A	Sequence, motif	N/A	Homologous method	N/A	<a href="http://cadd.pharmacy.nkai.edu.cn/tbdb">http://cadd.pharmacy.nkai.edu.cn/tbdb</a>	
[21]	Bacterium ( <i>B. anthracis</i> )	PARTIC	554	N/A	N/A	Sequence, graph properties	CTM variation, quadruples of consecutive amino acids	Four layers neural network	<a href="ftp://ftp.sanbi.aic.za.za/machine_learning/">ftp://ftp.sanbi.aic.za.za/machine_learning/</a>	N/A	



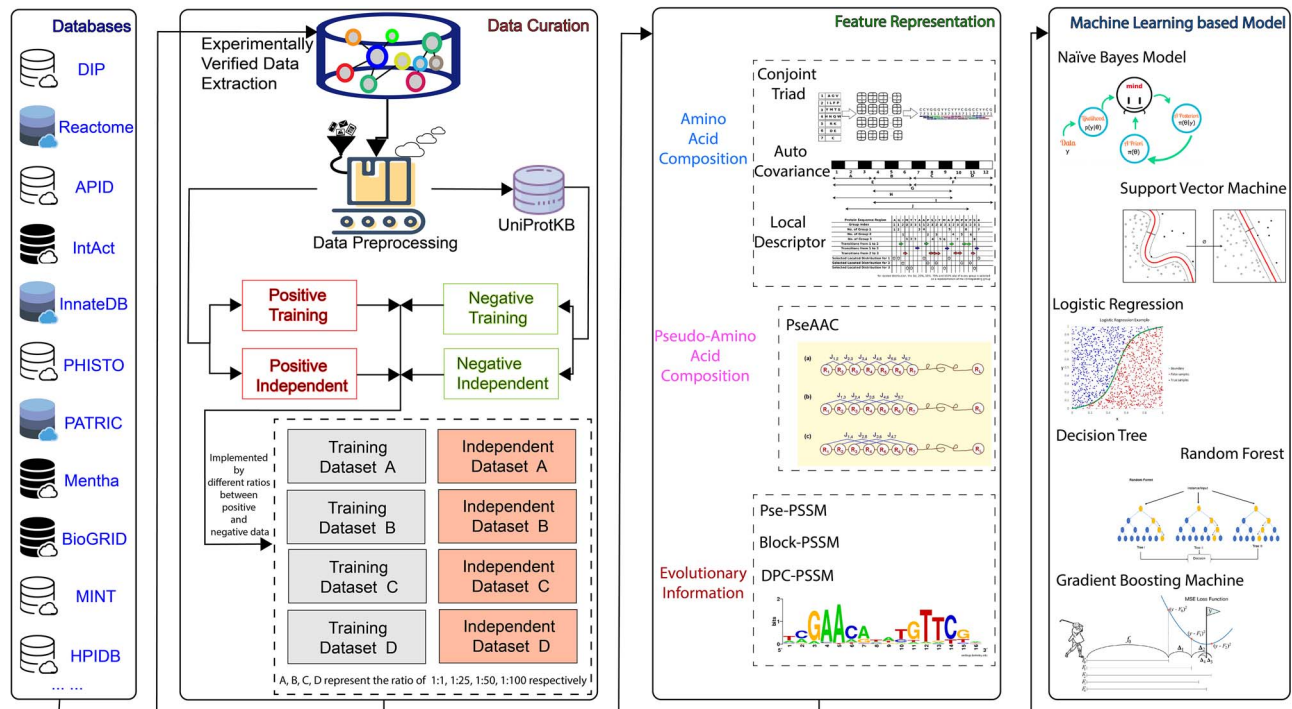


Figure 2. Overall framework of the performance-evaluation tests used in this study.

10, 15, 16, 18–21]. In this regard, we first grouped the sequential feature representation algorithms into three different types: amino acid composition, pseudo-amino acid composition and evolutionary information. It should be noted that not only the reported algorithms in Table 2 but also the related sequential representation algorithms from other protein sequence-specific topics, such as protein structure, protein-folding topics, are included in this section. The models from [8, 19], which are shown in Table 2, were selected as the representative models regardless of the pathogen species.

### HPI databases

There has been continuous effort spent on developing online HPI databases and repositories by many researchers. These developments mostly benefited from the National Institute of Allergy and Infectious Diseases (NIAID), which initialized a strategic plan to focus on biodefense research. Several ‘priority pathogens’ were defined. Several initial developments, including pathogen-interaction gateway (PIG [22]), BioHealthBase [23] and the Pathosystems Resource Integration Center (PATRIC [24]), were wholly or partially funded by the NIAID.

The first web-based database with massive annotated records for pathogen research was the Ecological Database of the World’s Insect Pathogens (EDWIP) [25]. EDWIP uses a one-to-one interaction relationship, which records the infection between a single host species and a single pathogenic species. This strategy resulted in 9400 records between 4454 host species and 2285 pathogen species when it was first released in 2003. PIG was designed as a collection of a number of public resources, which focussed on experimentally verified and manually curated HP-PPIs. This centralized database served as an easy-to-use database which transfers search results to the relevant database, such as the UniProt [26] database. Another important host-pathogen interaction database is the pathogen-host interaction search tool (PHISTO) [27]. This tool aims to provide researchers with a complete coverage of HPI data via monthly updates.

Proteomics Standards Initiative Common QUery InterfaGe (PSICQUIC) [28] service was installed to allow access to and extraction of HPI data the other web-based databases.

Although EDWIP is no longer available online, it is of particular interest for pathologists and ecologists to collect and analyse the HPI data. Concerning the HPI databases, one of the most critical factors in building a reliable benchmark dataset is the data sources. Typically, there are several different sources. One primary approach is to collect data from the literature and through manual verification by domain experts, such as the Database of Interacting Proteins (DIP) [29] and Reactome [30]. A second approach is to collect the data submitted by users. Finally, the data can also be novel derived or predicted data by computational models, such as the Pathogen–Host Interaction Data Integration and Analysis System (PHIDIAS) [31] and the Penicillium–Crop Protein–Protein Interactions database (PCPPI) [32]. Following the development of DIP and EDWIP, the HPI databases have become more interactive for the users. From Table 2, we can see that the most commonly used databases for HPI study, including DIP [29], IntAct [33], Mentha [34], the PHISTO [27], etc. can serve different purposes. According to the evaluation and data analysis of the databases, Mentha [34] has covered most of the PPI information for *H. sapiens*, while PHISTO focuses on studying human as the host species.

We, herein, have reviewed numerous publicly available databases, whose results were returned by searching specific keywords in the NCBI PubMed search engine. We manually examined the abstracts of the first 400 results ranked by ‘best relevance’ out of more than 4000 returned items based on the keywords ‘pathogen’ and ‘database’. As such, in this paper, a selection of 11 databases is reviewed and evaluated based on their contents. Details are provided in the following sections.

### Sequential representation algorithms

To encode proteins as feature vectors, several different features have been included in this study to predict PPIs between

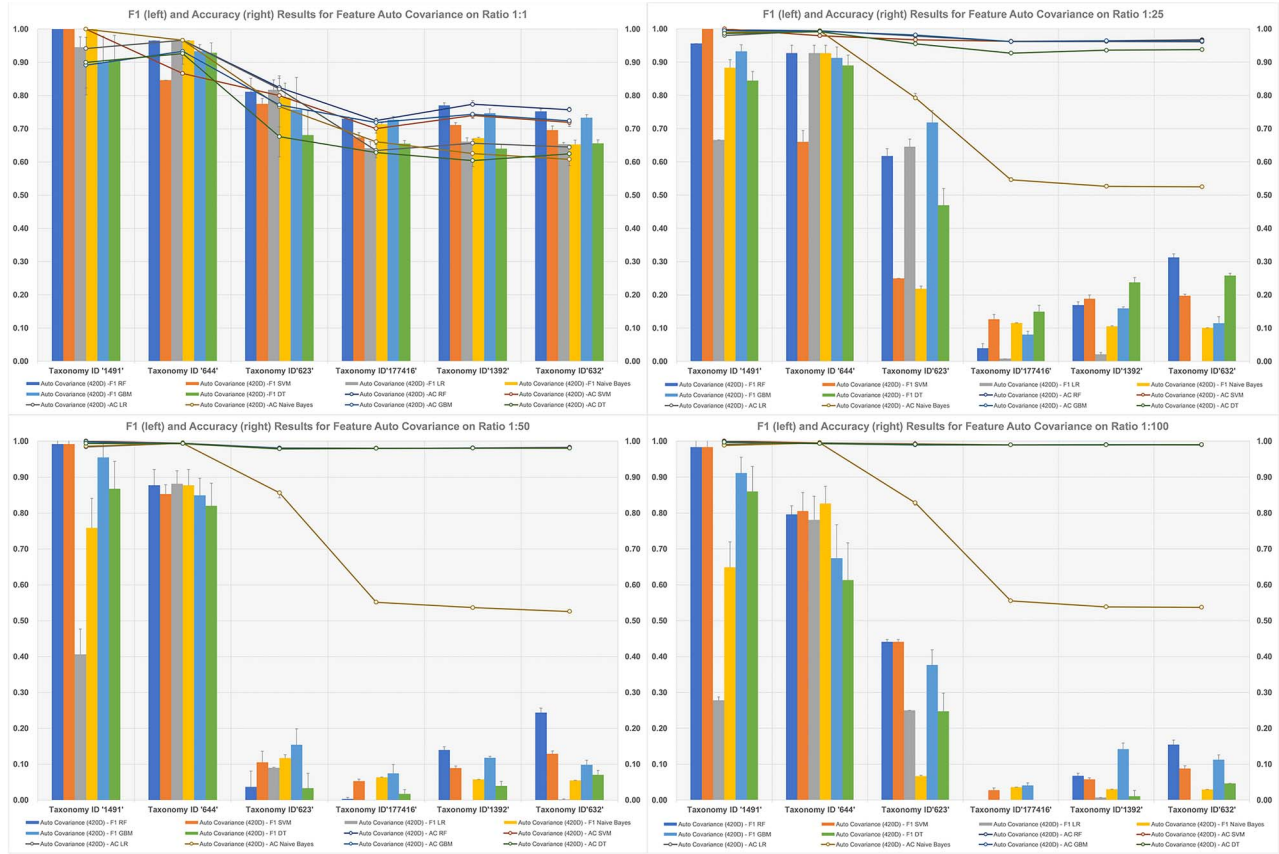


Figure 3. Performance evaluation results for AC in terms of F1 (bar charts) and ACC (line charts) based on different class ratios (1:1, 1:25, 1:50 and 1:100).

*H. sapiens* and bacterium pathogens, which are (1) protein amino acid composition information [35–37], (2) protein pseudo-amino acid composition information [38–40] and (3) protein evolutionary information features [41, 42]. We discuss the related feature-encoding algorithms below.

#### Amino acid composition

**Conjoint triad method.** It was proposed by [35] to classify 20 amino acids into 7 groups according to the dipole scale and volume scale of each amino acid, which describe their respective electrostatic and hydrophobic properties. Since this proposal, several variations of encoding algorithms for sequence representation have been devised based on this classification scheme. Among these, one popular approach is to consider the relationship of the properties of one amino acid and its vicinal amino acids as a descriptor [35], which is named the conjoint triad method (CTM). The conjoint triad information of several adjacent amino acids makes it easy to represent every single protein sequence as a class-based feature with the same length, which is also called its  $k$ -mer feature. Each amino acid type is indicated as a number ranging from 1 to 7 according to its group. The frequency of three conjoint triad data (3-mer) of a sequence is calculated. In total, there will be a combination set including  $\{(1, 1, 1), (1, 2, 1), \dots, (1, 7, 1), \dots, (1, 7, 7), \dots, (7, 7, 7)\}$ . As a result, 3-mer features will encode a sequence to a vector of 343 dimensions. For other 2-mer, 4-mer and 5-mer features, the features number would be 49, 2401 and 16 807, respectively.

**Auto covariance.** The auto covariance (AC) relationship among the amino acids based on the order of the sequence information was utilized in another feature representation algorithm

by [36]. It is a popular transformation algorithm used to adopt numerical vectors to uniform matrices by analysing sequences in the auto cross covariance (ACC) information. Between two different vectors, there are two covariance relationships: cross covariance (CC) and ACC. Only ACC variables are calculated [36]. The basic idea is to derive the physicochemical properties of the amino acid, which include its hydrophobicity (H), volume of side chains (VSCs), polarity (P1), polarizability (P2), solvent-accessible surface area (SASA) and the net charge index of the side chains (NCISC).

In the AC method, each single protein sequence is first translated into a numerical value corresponding to seven different physicochemical properties. Because the ranges of these seven physicochemical properties vary a lot from each other, a first step to normalize the numerical values is required. These values were hence normalized to a distribution whose mean is zero and the standard deviation is one. The normalization equation is shown in Eq. 1:

$$\bar{p}_{ij} = \frac{p_{ij} - \text{mean}_j}{\text{sd}_j} \quad (i = 1, 2, 3, \dots, 20; j = 1, 2, 3, 4, 5, 6, 7) \quad (1)$$

where  $p_{ij}$  represents the  $j$ th property value of the  $i$ th amino acid,  $\text{mean}_j$  is the mean value of the  $j$ th property over the 20 amino acids and  $\text{sd}_j$  is the standard deviation of the  $j$ th property over the 20 amino acids. Via this operation, every protein sequence is translated into an  $N * M$  matrix with zero mean and a standard deviation of unity in each column. With a proper range of these numerical values for each single protein sequence, AC can be used to represent them in a uniform matrix. Based on

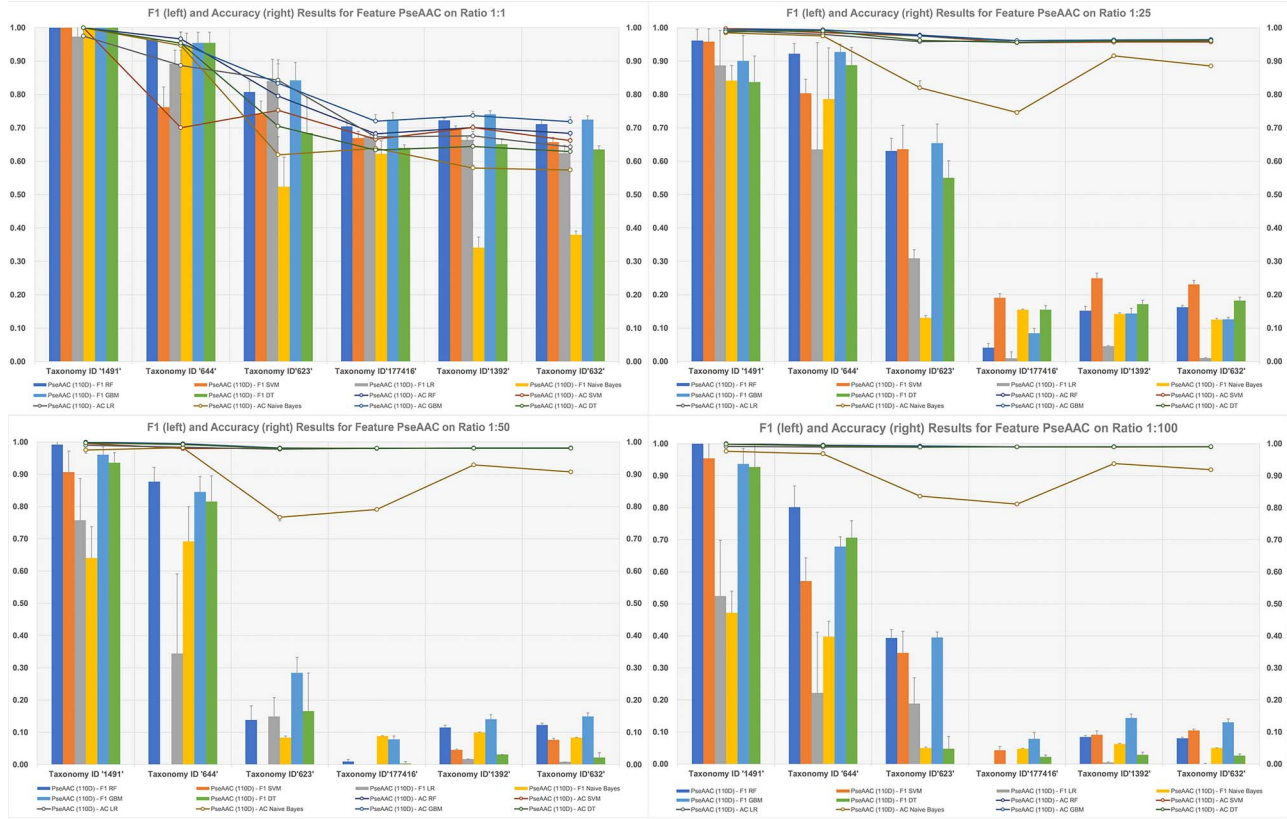


Figure 4. Performance evaluation results for PseAAC in terms of F1 (bar charts) and Acc (line charts) based on different class ratios (1:1, 1:25, 1:50 and 1:100).

Eq. 2, a matrix of  $lg \times 7$  was calculated, where lag is the distance threshold between two amino acids, and  $0 < lg \leq lag$ .

$$AC(lag, j) = \frac{1}{N - lag} \sum_{i=1}^{N-lag} \left( p_{ij} - \frac{1}{N} \sum_{i=1}^N p_{ij} \right) * \left( p_{i+lag, j} - \frac{1}{N} \sum_{i=1}^N p_{ij} \right) \quad (2)$$

For  $z$  properties chosen from the seven physicochemical properties, the length of AC is  $lag * z$ .  $p_{ij}$ , which corresponds to the value from  $\{p_{ij}\}$ . Here,  $N$  is the length of the protein sequence. After ACC transformation, a representation of the PPI is a concatenation of these two AC transform calculation results. **Local descriptor.** Another sequence-based feature representation method is a local descriptor [37]. The most important feature of an HP-PPI is that the interaction often occurs in some specific intermittent fragments. To better extract this continuous or discrete knowledge from sequence information, [37] proposed using region descriptors to first divide a protein sequence into 10 regions via 6 different methods: quarter regions, half regions (E, F), central 50% region (G), first 75% region (H), last 75% region (I) and the central 75% region (J). With these 10 regions, a local descriptor is utilized to transform the region sequence into three related descriptors [37]: composition (C), transition (T) and distribution (D). C is the composition ratio of each group of amino acid within a separate region, T represents the percentage of which amino acid group is followed by another amino acid group, and D describes the specific location information obtained by selecting the first, 25, 50, 75% and last of each amino acid group. When using a local descriptor, the extracted feature vector contains 7 C features, 21 T features and 35 D features. When multiplied by 10 different local regions, the local descriptor method generates

630 features for a single protein sequence. For an HB-PPI pair, this local descriptor contains 1260 features.

There are also some other schemes that can be used to extract different types of features of a protein sequence, for example, the Moran autocorrelation score [43] and the amino acid triplet [8]. As protein sequence information is directly linked to PPI, a further novel representation of PPIs, especially for HB-PPIs, might include any other information related to the specific host species and pathogenic species, which may be a better alternative for predicting HP-PPIs [9].

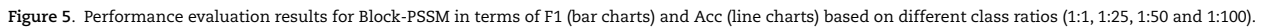
#### Pseudo-amino acid composition

Directly converting a protein sequence to a vectorized feature according to the amino acid composition (AAC) might result in sequence order information loss. The pseudo-amino acid composition (PseAAC) method was proposed as a novel protein sequence representation of a discrete model, which has remarkable prediction performance as an important feature representation algorithm [38, 44–47].

Various modes of PseAAC have been introduced in the literature. The key is to combine the sequence order correlation information from the protein sequence. In the work of [38], the original version of PseAAC was introduced, as shown in Eq. 3:

$$\begin{aligned} \theta_1 &= \frac{1}{T-1} \sum_{i=1}^{T-1} \Theta(S_i, S_{i+1}) \\ \theta_2 &= \frac{1}{T-2} \sum_{i=1}^{T-2} \Theta(S_i, S_{i+2}) \quad \lambda < T \\ \theta_\lambda &= \frac{1}{T-\lambda} \sum_{i=1}^{T-\lambda} \Theta(S_i, S_{i+\lambda}) \end{aligned} \quad (3)$$





block substitution matrix BLOSUM [55]). The PSSM is calculated according to Eq. 5:

$$P_{m,n} = \sum_{k=1}^{20} w(m,k) * \theta(n,k) \quad (5)$$

where  $w(m, k)$  is the probability that the  $k$ th amino acid appears at position  $m$  and  $\theta(n, k)$  is the value of the position of  $(n, k)$  in the similarity matrix.

In this study, PSI-BLAST was employed to create PSSMs with three iterations, where the *e*-value was set to 0.001. Accordingly, the various lengths of the protein sequences resulted in matrices with different dimensions, which introduced different encoding features based on the PSSM profiles. The following parts present several PSSM-based feature representation algorithms.

Pse-PSSM. The pseudo position-specific score matrix (Pse-PSSM) was first introduced for the task of predicting whether or not an uncharacterized protein was a membrane protein [44]. Pse-PSSM extends the idea of corrupting the PSSM descriptor vertically as a mean value, as shown in Eq. 7, where the value of the PSSM is first processed by a standardization procedure horizontally by rows in Eq. 6. The concept of the PseAAC is to generate correlation information between different amino acid locations.

$$p'_{m,n} = \frac{p_{m,n} - \frac{1}{20} \sum_{k=1}^{20} p_{m,k}}{\sqrt{\frac{1}{20} \sum_{k=1}^{20} \left( p_{m,k} - \frac{1}{20} \sum_{k=1}^{20} p_{m,k} \right)^2}} \quad (6)$$

$$\bar{p}_n = \frac{1}{T} \sum_{m=1}^T p'_{m,n} \quad (n = 1, 2, \dots, 20) \quad (7)$$



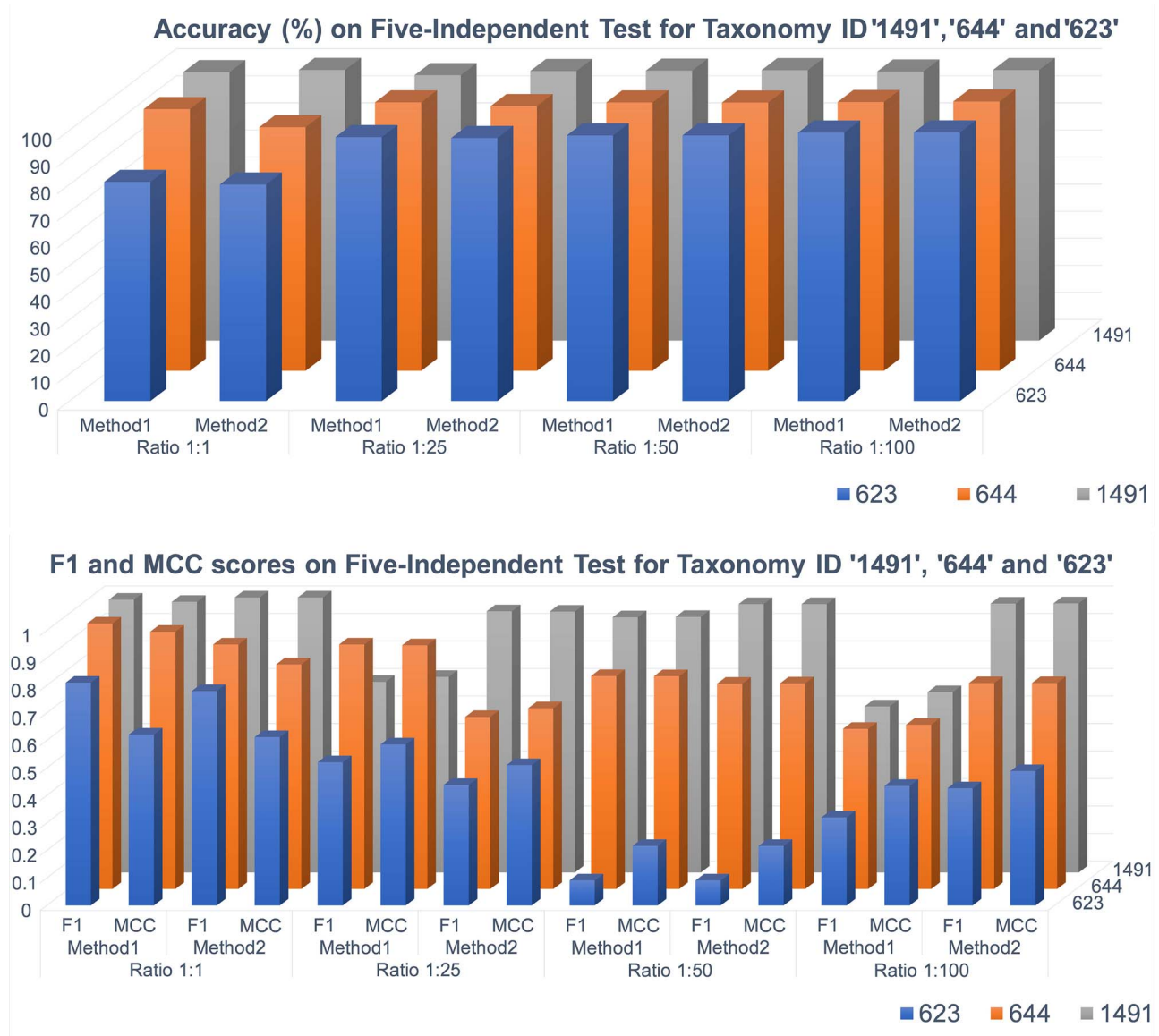


Figure 6. The performance results of taxonomy IDs '1491', '644' and '623' in terms of Acc, F1 and MCC.

Thus, the original PSSM profile is converted to a 20-dimension vector,  $\bar{p} = \{\bar{p}_n, n = 1, 2, \dots, 20\}$ . This derived feature focuses on representing the average score of each amino acid type according to the reference database, which loses the sequence order information of the protein. Thus, [44] proposed considering supplementary information from the PseAAC, which slices the PSSM profile according to Eq. 8:

$$p_{se_n} = \frac{1}{T-c} \sum_{m=1}^{T-c} [p'_{m,n} - p'_{(m+c),n}]^2 \quad (n = 1, 2, \dots, 20; c < T) \quad (8)$$

This process generates a 40-dimension vector  $P_{se} = \{\bar{P}_1, \bar{P}_2, \dots, \bar{P}_{20}, \bar{P}_{se1}, \bar{P}_{se2}, \dots, \bar{P}_{se20}\}$  where  $0 < c < \min(T)$ . For a given set of protein sequences, the upper bound of  $c$  should be smaller than the shortest length of the protein sequences.

**Block-PSSM.** By considering the PSSM profile in a dimension format of  $T \times 20$ , [56] proposed dividing the whole sequence into 20 equal blocks, where each represents 5% of the total sequence.

Each block generates a 20-dimension vector, which is finally combined as a  $20 \times 20 = 400$ -dimension vector.

The  $i$ th block is calculated according to Eq. 9:

$$pblock_{i,j} = \frac{1}{B_i} \sum_{i=1}^{B_i} p_{i,j} \quad i = 1, 2, \dots, 20; j = 1, 2, \dots, 20 \quad (9)$$

where  $i$  represents the block number. Since each 5% of a sequence is considered as a block,  $i$  ranges from 1 to 20 and  $j$  is the number of amino acid types. In short,  $pblock_{i,j}$  is extracted as a  $1 \times 20$  vector, thus  $pblock = pblock_1, pblock_2, \dots, pblock_{20}$  is calculated as the Block-PSSM feature in the form of  $1 \times 400$  vector feature.

**DPC-PSSM.** Another variation of the PSSM-based feature was proposed by [57]. The original PSSM profile is scaled to the range 0–1 by following a sigmoid function, as shown in Eq. 10:

$$p''_{m,n} = \frac{1}{1 + e^{-p_{m,n}}} \quad (10)$$

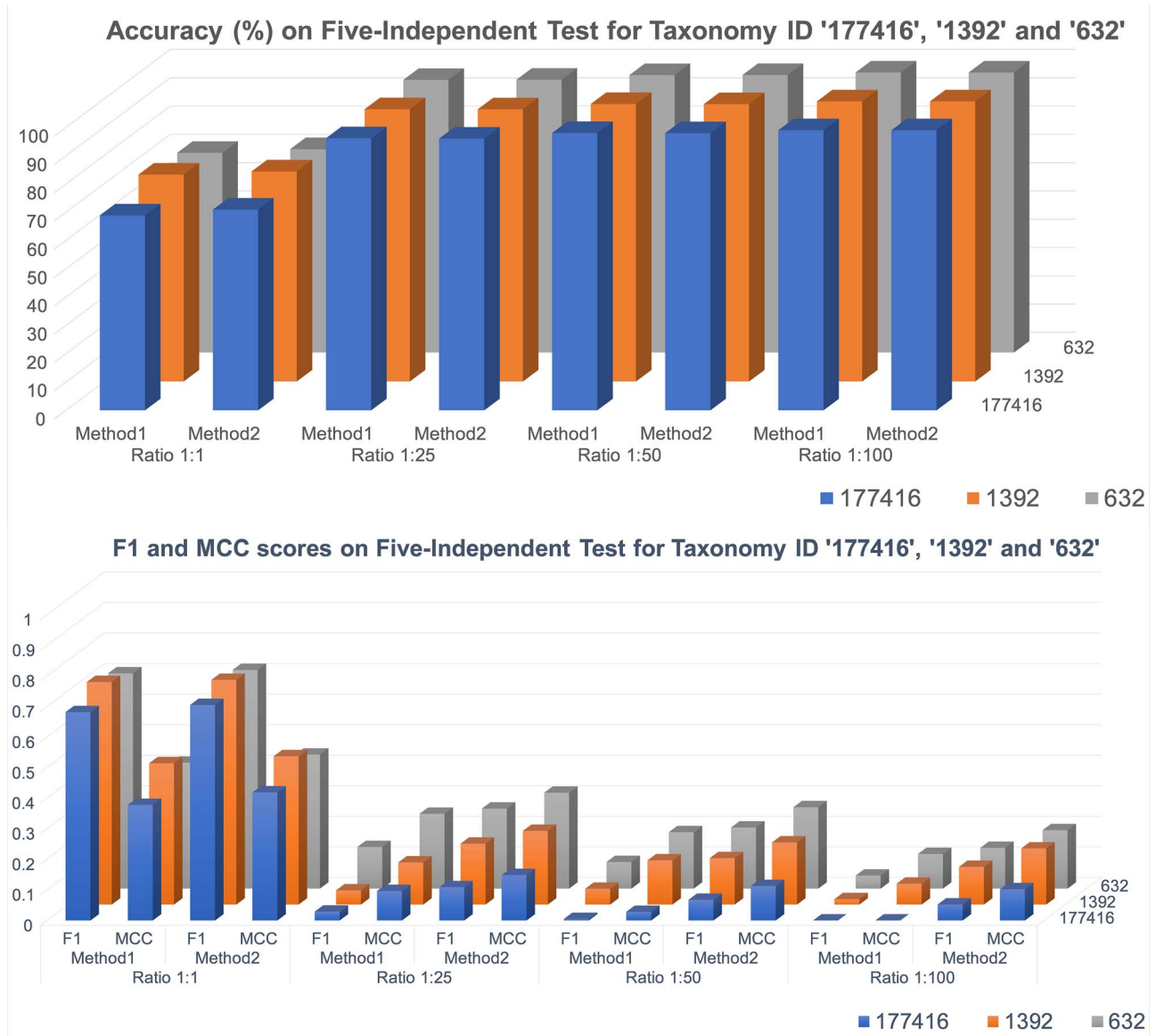


Figure 7. The performance results of taxonomy IDs '177419', '1392' and '632' in terms of Acc, F1 and MCC.

where  $p''_{m,n}$  is also used in the transition probability composition PSSM. The AAC-PSSM method is used to extract the corresponding AAC information from  $p = \{p''_{m,n}, m = 1, \dots, T; n = 1, \dots, 20\}$ . The vector in Eq. 11 represents the average mutation score of the amino acid types in the protein during the evolution process, namely, the AAC-PSSM. This calculation generates a 20-dimension feature vector.

As a supplementary approach, the traditional dipeptide composition (DPC) of the protein sequence was extended by [57], which was named DPC-PSSM. The calculation of DPC-PSSM is based on the covariance between two adjacent amino acid residues, denoted in Eq. 12. This process produces a 400-dimension feature vector.

$$P_{aacn} = \frac{1}{T} \sum_{m=1}^T p''_{m,n} \quad m = 1, \dots, T; n = 1, \dots, 20 \quad (11)$$

$$P_{dpc_{ij}} = \frac{1}{T-1} \sum_{k=1}^{T-1} p''_{k,i} * p''_{(k+1),j} \quad i, j = 1, \dots, 20 \quad (12)$$

### Machine learning-based methods for prediction

Applying computational approaches to predict bioinformatics tasks is considered an important supplementary method for identifying specific targets and high-fidelity interactions in experiments. Recently, we have witnessed numerous applications focusing on domains containing an abundance of unknown data, which require hypothesis verification [5, 13, 58, 59].

In Table 2, the predictors of [8, 19], which are based on machine learning methods and protein sequence information, were selected for our study. These machine learning models include support vector machine (SVM) and random forest (RF). In this section, we will first briefly review most of the potential machine learning models that can be utilized for HB-PPI prediction, which include logistic regression (LR), the Naïve Bayes (NB) model and gradient boosting machine (GBM). These models have demonstrated their capability in other applications for protein structure prediction; however, this is the first time they have been presented in an overall performance evaluation

in relation to different feature representation algorithms for HB-PPIs.

**Support vector machine (SVM).** SVM is one of the most widely used models in the literature, which was originally developed by [60]. The introduced structural risk minimization theory ensures the performance of SVM to be widely and successfully applied to many classification and regression tasks in computational biology. SVM contains a radial basis function (RBF) kernel, which is given the task of classifying HP-PPI pairs [8, 11]. Given a dataset of HB-PPIs denoted as  $\{x_i, y_i\}, i = 1, 2, \dots, N$ , where  $x_i \in \mathbb{R}^n$  and  $y_i \in \{+1, -1\}$ ,  $y_i$  is calculated as shown in Eq. 17:

$$y(x) = \text{sign} \left[ \sum_{i=1}^N y_i \alpha_i * K(x, x_i) + b \right] \quad (17)$$

where  $K(x, x_i) = \exp(-\gamma \|x_i - x\|^2)$  stands for the RBF kernel and  $\alpha_i$  contains the parameters from a convex quadratic programming problem.

**Decision tree (DT).** The DT was designed as a non-parametric supervised model [61]. It uses a tree-like graph to predict an incoming instance based on learned decision rules from given data samples and represented features. DTs are simple to understand and interpret, and they are capable of handling both numerical and categorical data.

**Random forest.** Derived from the DT model, RF adopts a random learning method to construct a combination of DTs [62]. RF has superior performance compared with other machine learning algorithms for classification tasks [63, 64], regression tasks and so on. Technically, RF is an ensemble learning model based on the tree bagging method, which builds a bunch of random DTs to avoid the latent problem caused by potentially biased data. In this study, we implemented RF using the scikit-learn toolkit [65] in Python.

**Logistic regression.** LR is an important machine learning model, which targets modelling  $y_i$  between 0 and 1 given unseen data  $x_i$  [66, 67]. Accordingly, the LR returns results via Eq. 18:

$$\begin{aligned} P(y_i = 1 | x_i) &= h_\theta(x_i) = 1 / (1 + \exp(-\theta^T * x_i)) \\ P(y_i = 0 | x_i) &= 1 - P(y_i = 1 | x_i) = 1 - h_\theta(x_i) \end{aligned} \quad (18)$$

where  $\theta$  is the combination of the model parameters and the optimization of  $\theta$  is solved with either the cross-entropy function  $J_1$  or the mean square error loss function  $J_2$ , as shown in Eq. 19:

$$\begin{aligned} J_1(\theta) &= -\sum_i (y_i \log(h_\theta(x_i)) + (1 - y_i) \log(1 - h_\theta(x_i))) \\ J_2(\theta) &= \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2 \end{aligned} \quad (19)$$

**Naïve Bayes model.** Based on Bayes' theorem [68, 69], the Naïve Bayes model consists of a probabilistic classifier and considers features as independent variables when the class label is given. Given  $X = (x_1, x_2, \dots, x_n)$ , where  $x_i$  is the  $i$ th feature, the probability of being in category  $y_k$  is calculated via Eq. 20:

$$p(y_k | X) = \frac{p(y_k)}{p(X)} \prod_{i=1}^n p(x_i | y_k) \quad (20)$$

In this study, we selected the Gaussian Naïve Bayes (GNB) model to deal with the continuous data produced by the various feature representation algorithms. The distribution of the

data was assumed to be a Gaussian distribution, which follows Eq. 21:

$$p(x_i | y_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \quad (21)$$

where  $\mu_k$  is the mean of  $X$  and  $\sigma_k^2$  is the corresponding variance.

**Gradient boosting machine.** GBM was first developed as a greedy optimization model [70] for both regression and classification tasks. Among the variants of GBM, gradient tree boosting is a frequently used model integrated with DTs. Given  $X = (x_1, x_2, \dots, x_n)$ , in which  $x_i$  is related to label  $y_i$ , gradient tree boosting builds an ensemble of trees sequentially by distilling the gradient descent algorithm into the process of new tree construction. A new tree is constructed under the discrepancy between the target function  $f(x)$  and current model, in which  $f(x_i) = y_i$ . The discrepancy between the target function  $f(x)$  and the current model is also called the residual of GBM.

## Materials

### Human-bacterium interaction resources

In this section, we first collected and reviewed 11 public databases, as summarized in Table 3: the Database of Interacting Proteins (DIP) [29], Reactome [30], the Agile Protein Interaction DataAnalyzer (APID) [71], IntAct [33], the Molecular Interaction Database (MINT) [72], the InnateDB [73], the PHISTO [27], the Pathosystems Resource Integration Center (PATRIC) [24], Mentha [34], the Host Pathogen Interaction Database (HPIDB) [74, 75] and the Biological General Repository for Interaction Datasets (BioGRID) [76].

As humans are one of the primary host species among infectious diseases, the HPI resources are considered as the preliminary investigation subjects from all these databases. The column 'HPI number' indicates the corresponding recorded interaction number from the databases, which contain both inter-species interactions and intra-species interactions. These 11 databases were selected because their data sources mainly come from the literature, which have been subjected to expert manual verification, and public archival databases, which also contain high confidence of the presented data.

Taking database PATRIC [24] as an example, the data source was built upon several public archival databases, such as MINT [72], IntAct [33], BioGRID [76] and DIP [29]. The cross-archived databases have extended the availability of HPI resources; however, some duplicates inevitably occur during the combination of these 11 databases. Thus, we followed the traditional data collection and cleansing methods from the literature [7–9].

### Curation

In this section, we briefly describe the major statistics for 'golden dataset' curation, which will be thoroughly surveyed in the following sections.

**Positive interactions.** Six different types of bacteria were selected, and the related data were preprocessed from the available databases. We identified the bacteria by mapping the taxonomy IDs according to the NCBI Taxonomy database. In Table 4, the corresponding information, including taxonomy ID, organism name, total pair number from the database and the number after cleansing, are presented. These 11 databases were accessed and downloaded in September 2018.

**Table 3.** HPI resources

Database	Data source	Data type	HPI number
DIP	Literature and domain expert manual verification	Protein–protein interactions	76 882
Reactome	Literature and domain expert manual verification	Comprehensive data portal including pathway and analysis	1 016 953
APID	Public archival databases	Protein–protein interactions	133 994
IntAct	Public archival databases and literature	Molecular interaction database	857 826
MINT	Literature	Protein–protein interactions	123 892
InnateDB	Literature	Mammalian innate immunity networks, pathways and genes	24 077
PHISTO	Public archival databases	Host–pathogen and human intra-species protein–protein interactions	90 453
PATRIC	Public archival databases	Comprehensive data portal for bacterium pathogens	618 737
Mentha	Public archival databases	Protein–protein interactions	1 272 096
HPIDB	Public archival databases and literature	Host–pathogen interactions	62 783
BioGRID	Literature	Comprehensive data portal for protein, genetic and chemical interactions	1 568 115

**Table 4.** The positive interaction statistics of selected bacterium species

Taxonomy ID	Bacterium pathogen	Total number from databases	After cleansing
1491	<i>Clostridium botulinum</i>	61	57
644	<i>Aeromonas hydrophila</i>	73	73
623	<i>Shigella paradysenteriae</i>	118	105
177 416	<i>Francisella tularensis</i> subsp. <i>tularensis</i> (strain SCHU S4/Schu 4)	1319	1207
1392	<i>B. anthracis</i> bacterium	3275	2810
632	<i>Yersinia pseudotuberculosis</i> subsp. <i>pestis</i> (Lehmann and Neumann 1896) Bercovier et al. 1981	4114	3528

In Table 4, the statistics refer to the results of the representative proteins. Meanwhile, any proteins with fewer than 50 amino acids were removed since these proteins may be non-functional fragments. The protein sequence information was primarily from the SwissProt/UniProtKB database [77].

**Negative interactions.** How to select feasible negative PPIs remains an active topic for the prediction of PPIs. Currently, there is not a standard protocol defining both the negative pairing strategy and the ratio to positive interactions. In most cases, building a negative interaction dataset by randomly selecting protein pairs from a set of unknown interacting relationships between protein pairs is utilized. This heuristic approach works well in practice as the interaction ratio (i.e. the number of positive interactions in a large, random set of protein pairs) is expected to be very low, which in the work of [7] was defined as 25, 50 and 100 times as many negative examples as positive examples. In the study by [9], the ratio was set to 1/100. The assumption in this approach is that the probability that the selected negatives contain true positives is negligible.

Thus, we followed the traditional approaches from the literature [7, 9, 10, 21]. A random pairing for a negative PPI was first undertaken between different proteins sets, which in this study was between the chosen bacterium pathogens (listed in Table 4) and *H. sapiens* proteins (taxonomy ID: 9606). Then, we randomly selected a subset from this random pairing set to be the negative

dataset. The negative interactions were selected with different ratios: 1:1, 1:25, 1:50 and 1:100.

**Protein information.** When building machine learning models to predict PPIs, HB-PPIs are needed to utilize the diverse protein information, which can be divided into three groups: structure-based, domain-based and sequence-based protein information.

Numerous studies have utilized and examined different information when predicting specific HP-PPIs [7, 78, 79]. Particularly, domain–domain and structure–structure interaction methods are the two main approaches used to complement existing high-confidence interactions [7, 79]. Also, structural similarity, which refers to a result of homology-based modelling, is an important alternative for detecting proteins with a homogeneous structure based on experimentally verified HP-PPIs [78].

Although structure-based and domain-based information have some benefits for exploring HPIs [80, 81], they can limit the scope of studied HP-PPIs to specific genres and species, such as HIV-1, HCV, Ebola viruses and so on [7, 10, 58, 82–84]. One dominant reason is the limited amount of available experimentally determined structures and domain information, particularly for bacteria. Imputation remains a core technology to compensate for the dearth of protein information and helps to address the challenge of interaction predictions [79]. Imputation for missing data also impacts the prediction performance since it brings putative information, which might not be accurate. Thus, utilizing structure-based and domain-based information



**Table 5.** Overview of the protein information for the dataset preparation process

Taxonomy ID	Whole proteome information		Positive information		Positive pairs	Total no. of HB-PPI	Negative pairs			
	Human	Bacterium	Human	Bacterium			1:1	1:25	1:50	1:100
1491	18 181	524	9	7	57	9.5 M	57	1425	2850	5700
644	18 181	511	66	4	73	9.3 M	73	1825	3650	7300
623	18 181	1724	75	60	105	31.3 M	105	2625	5250	10 500
177 416	18 181	550	889	306	1207	10.0 M	1207	30 175	60 350	120 700
1392	18 181	1501	1537	844	2810	27.3 M	2810	70 250	140 500	281 000
632	18 181	1893	1866	1092	3528	34.4 M	3528	88 200	176 400	352 800

Note: Only the proteins processed used in the positive interactions are counted in this table; M is short for 'million'. For each human-bacterium PPI dataset, the number of pathogen proteins, the size of the dataset and other such statistics are shown.

limits the availability and scalability to a wide range of studies of HB-PPIs.

Alternatively, there has been a research trend of predicting PPIs from sequence-based protein information [35, 85]. Sequence-based protein information is one of the most abundant sources of protein data, which has stimulated ongoing research to improve the prediction performance of novel feature representation and machine learning models [8, 14, 21, 86, 87]. Sequence-based methods enable the models to be applied to large datasets and various species and genres.

**Independent datasets.** To help our readers understand each dataset's information, in Table 5, all the protein numbers related to the different subsets were included. This information, which was related to the reviewed sequence information from the UniProtKB database [26], was last updated on 30 October 2018. In all, we collected 18 181 *H. sapiens* protein sequence information, and the corresponding protein numbers for each taxonomy ID are reported in Table 5.

Evaluation of the models requires careful preparation of the independent datasets. Generally, cross validation shows better performance than the independent testing model for an unseen dataset. To give a general performance evaluation, we followed [8] when we built the independent datasets. The difference was that we further built 5-fold independent datasets, which helped us to better measure the means and variations of the machine learning models. The independent datasets were not used during the training, and various measurements were included to evaluate the performance of different models based on the independent datasets. Thus, we first randomly selected one-fifth of the PPIs from both positive and negative interactions to be the independent dataset. The remaining PPIs of positive and negative interactions were then combined as the training set. We assembled the negative interactions with a random sampling method, where random sampling of the negative interactions was conducted five times, which allowed us to evaluate the different models with statistic means and variations to reduce the bias caused by negative interactions. The involved protein numbers for *H. sapiens* and the corresponding bacterium taxonomy IDs are reported in detail in the appendix. We have reported the number of utilized proteins for each species for different ratio settings. We anticipate that this experimental setting and details will help to provide more information to build novel machine learning methods in future work.

The framework of our evaluation study is presented in Figure 2. In Figure 2, a clear process procedure from databases to training and independent datasets, followed by the feature

representation algorithms and machine learning model evaluations, are mapped in a coherent line. The best model selection and prediction are given as the main outcome of this framework.

## Evaluation results

### Evaluation metrics

A set of six popular performance evaluation metrics, including precision (Pre), accuracy (Acc), sensitivity (Sn), specificity (Sp), F1-score and Matthew's correlation coefficient (MCC) were applied to evaluate the overall prediction performance of the models [46, 88–94]. The measurements are defined as follows:

$$\begin{aligned}
 \text{Pre} &= \frac{TP}{TP+FP} \\
 \text{Acc} &= \frac{TP+TN}{TP+FP+TN+FN} \\
 \text{Sn} &= \frac{TP}{TP+FN} \\
 \text{Sp} &= \frac{TN}{TN+FP} \\
 \text{F1} &= \frac{2 * \text{Pre}}{\text{Pre} + \text{Rec}} \\
 \text{MCC} &= \frac{(TP * TN) - (FN * FP)}{\sqrt{(TP+FN) * (TN+FP) * (TP+FP) * (TN+FN)}}
 \end{aligned} \tag{22}$$

where TP, FP, TN and FN represent the numbers of true positives, false positives, true negatives and false negatives, respectively. Also, the receiver operating characteristic (ROC) curve and the area under the curve (AUC) were included to quantify the model performances.

### Performance evaluation based on different class ratios

One primary evaluation of this study was the ratio impact of different predictors, which was the ratio between positive and negative protein interactions. We herein present the F1 score and Acc value from our measurements for feature 'ACC' for the evaluation discussion. Since the curated HP-PPI datasets involved different ratios between the positive and negative interactions data, Acc could be used to more precisely measure the performance of the model in a more accurate way at the ratio of 1:1. However, if the ratios become more skewed, such as 1:25 to 1:100, the F1 score would be a more suitable performance measurement. The mean value and deviation of each of the five independent tests were calculated in terms of different bacterial species and building ratio settings between the positive and negative pairs. In general, the ability to predict positive interactions as negative pairs decreases both the F1 and Acc results. Here, we found that the Acc was as high as 0.990099 when all the test data were predicted as negative interactions for a ratio of 1:100 between the positive and negative interactions. For ratios of 1:25, 1:50 and 1:100 between the positive and negative interactions, the

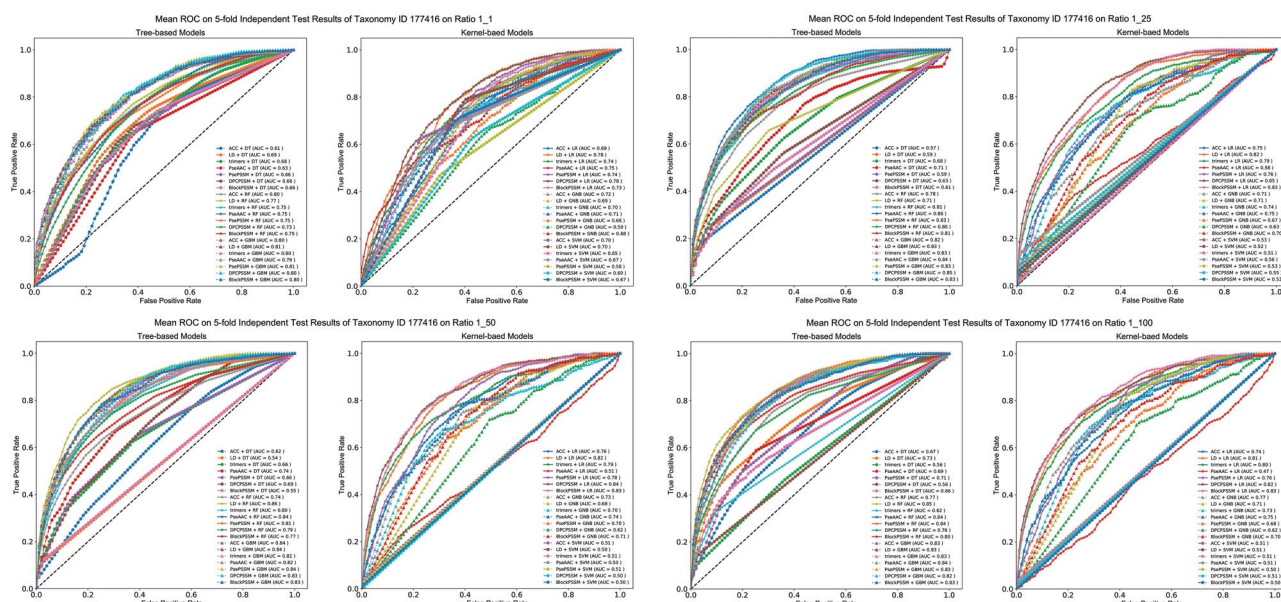


Figure 8. ROC curves for taxonomy ID '177416'.

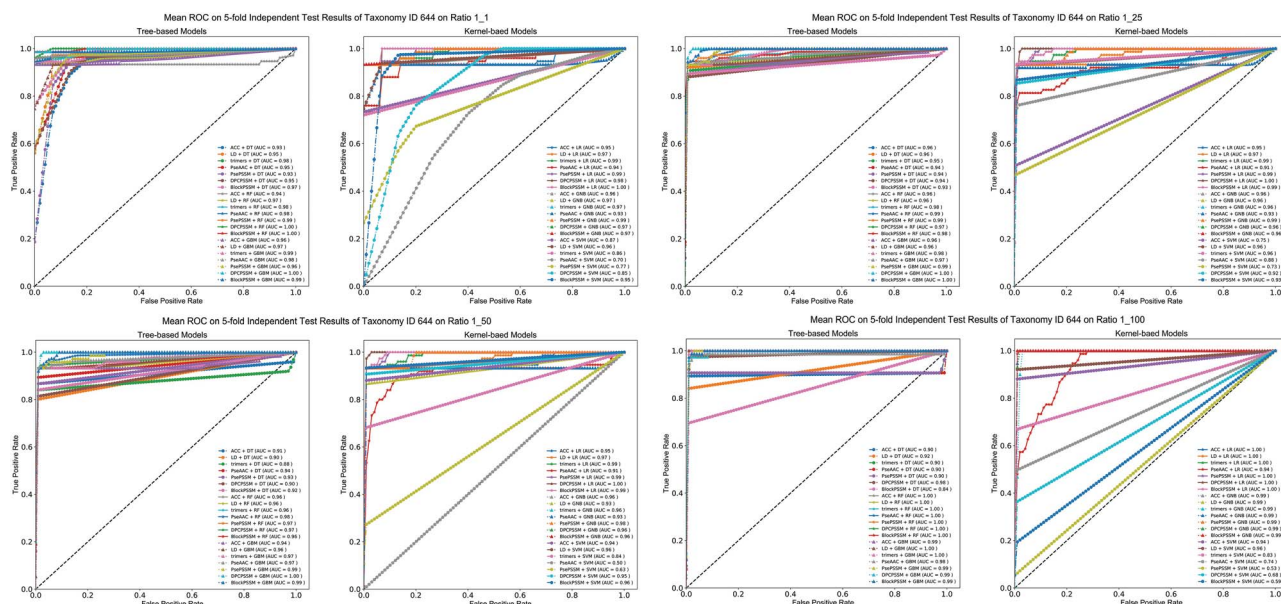


Figure 9. ROC curves for taxonomy ID '644'.

datasets were considered as imbalanced datasets. Therefore, the F1 score was more suitable for measuring the performance of imbalanced datasets.

From Figure 3, it is easy to see that the F1 scores present a trend of getting worse as the dataset becomes larger and more complex, which means more protein nodes and edges are involved in the dataset. For example, when the positive to negative ratio was 1:1, a  $1.0 \pm 0.0$  F1 score was found for the RF algorithm and the taxonomy ID is '1491'. However, the F1 score became  $0.96 \pm 0.0$  with RF for ID '644',  $0.817555 \pm 0.029558$  with LR for ID '623',  $0.730386 \pm 0.005192$  with RF for ID '177416',  $0.770171 \pm 0.007703$  with RF for ID '1392' and  $0.752226 \pm 0.006632$  with RF for ID '632'.

In Figures 4 and 5, feature 'PseAAC' from the PseAAC method and feature 'Block-PSSM' from the evolutionary information method are also included for different ratios. The performance comparison between these two different sequence-based features also indicates the impact of the ratio upon the F1 and Acc results.

From Figure 3, we can see that all the predictors have worse performance for all datasets when the ratio increases from 1:1 to 1:25, 1:50 and 1:100, especially when the dataset is with more than 100 000 samples. For example, for taxonomy ID '632', the F1 score was  $0.752226 \pm 0.006632$  for a 1:1 ratio; however, the F1 scores dropped to  $0.312530 \pm 0.010944$  for a 1:25 ratio,  $0.243679 \pm 0.012883$  for a ratio of 1:50 and  $0.154535 \pm 0.012569$

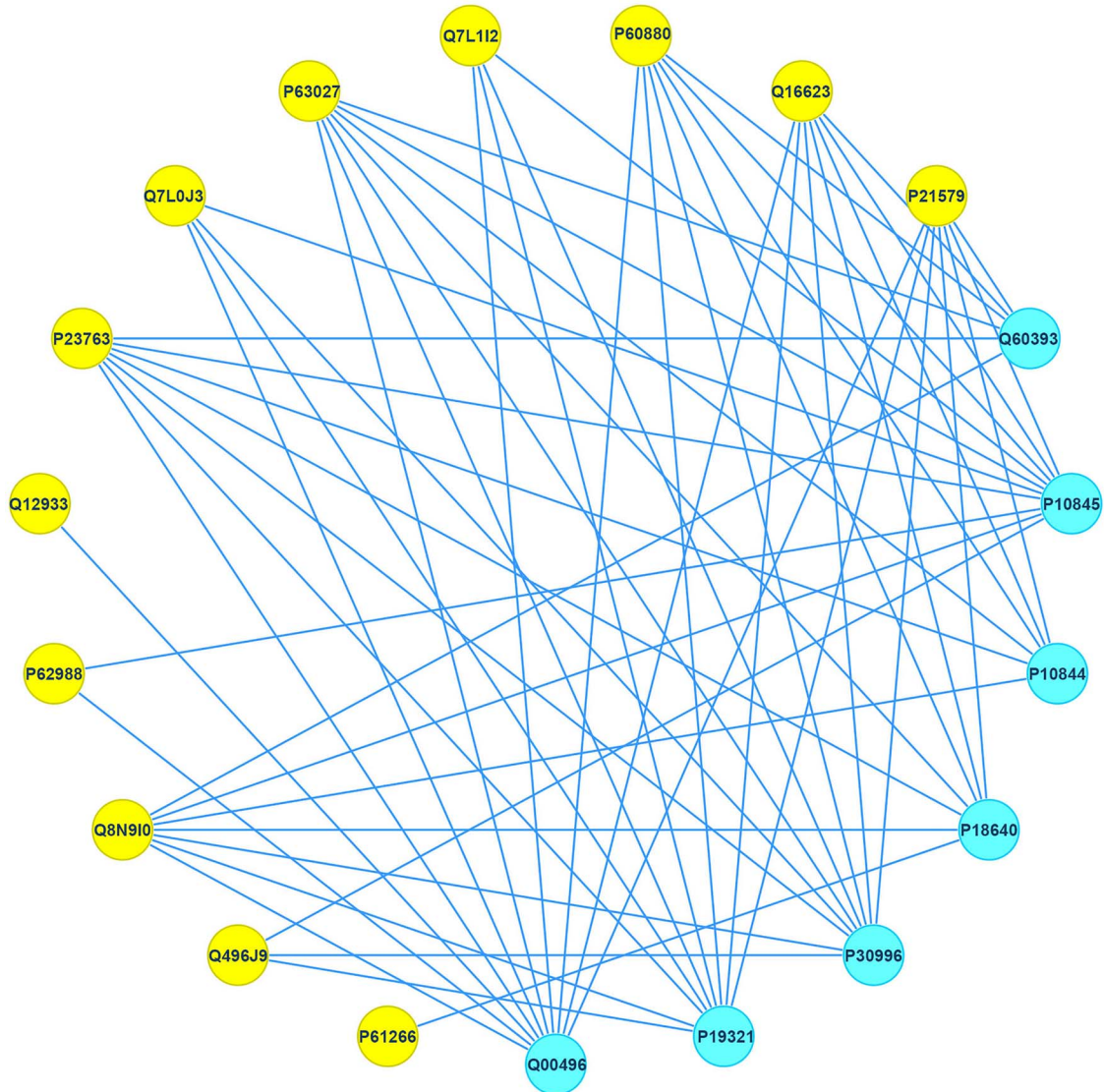


Figure 10. Protein interactions network for taxonomy ID 1491.

for the 1:100 ratio. These results were all achieved with the RF algorithm.

In Figures 6 and 7, the results of the existing available methods are included. Figure 6 contains the Acc, F1 and MCC scores for IDs '1491', '644' and '623', and Figure 7 contains the results for IDs '177419', '1392' and '632'. Both Figures 6 and 7 indicate the performance variation when the dataset changes from taxonomy ID '1491' to '644' and '623', which becomes worse for taxonomy IDs '177419', '1392' and '632'. Even though the existing methods in Figures 6 and 7 have incorporated several novel sequential feature representation algorithms, their performance has not improved.

### Overall performance

Figures 8 and 9 show the ROC curves for taxonomy IDs '177416' and '644', respectively.

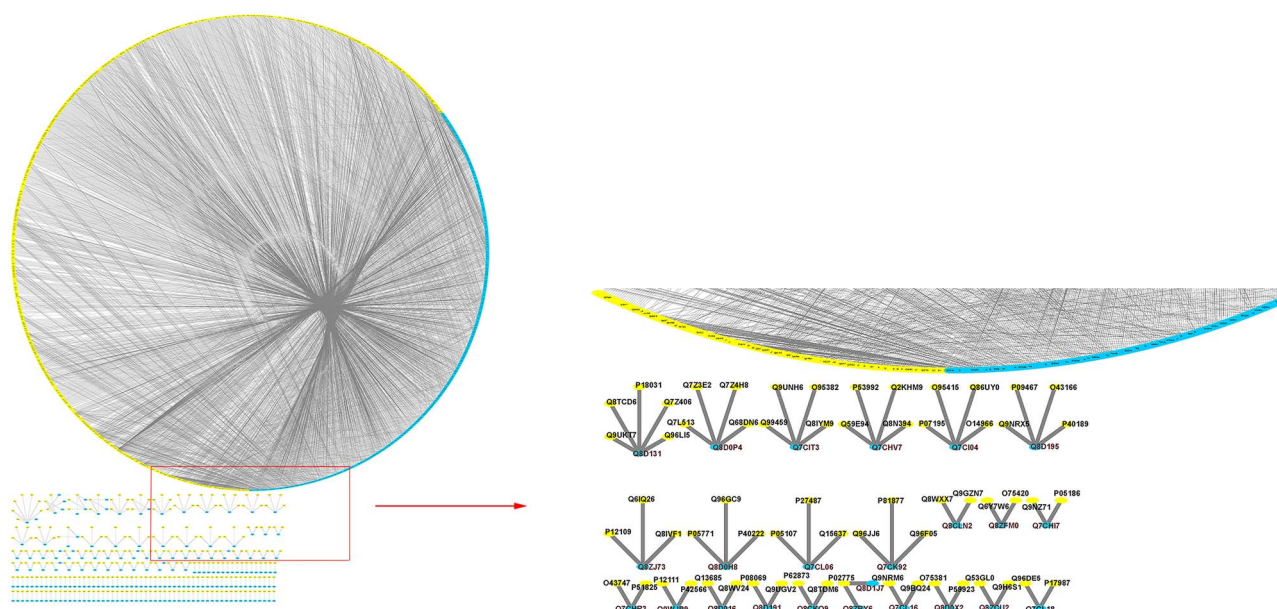
We have listed the six evaluated machine learning models as two groups. One group contains tree-based models, which

includes DT, RF and GBM. The other group consists of kernel-based models including SVM, LR and the Naïve Bayes model. The performances are presented as mean ROC curves from 5-fold independent test results for different ratios.

Because there are 1207 positive interaction pairs for taxonomy ID '177416', the dataset size is 121 907 for a ratio of 1:100, which is larger than that of taxonomy ID '644'. Somehow, the predictors' performance became worse for the larger dataset. Although the tree-based models still outperformed the kernel-based models for each dataset, the overall performance was not stable across the different host-bacterium systems.

In Table 6, the best results of all the predictors are listed accordingly for taxonomy ID '632'. For example, for the AC feature representation algorithm dataset, the best results for ratios of 1:1, 1:25, 1:50 and 1:100 were all achieved by RF model with accuracies of  $0.757082 \pm 0.008000$ ,  $0.967350 \pm 0.000365$ ,  $0.982521 \pm 0.000128$  and  $0.990674 \pm 0.000043$ , respectively. The tree-based models, including DT, RF and GBM, have demonstrated a strong generalization ability in terms of providing





**Figure 11.** Protein interaction network for taxonomy ID 632.

effective and efficient performance. The other models, such as kernel-based model, including SVM, Gaussian Naïve Bayes (GNB) model and the LR model, however, are less robust compared with the tree-based models. Meanwhile, the training time was in higher demand than for the tree-based models. Taking CTM as the feature representation algorithm, the time spent training GBM for the dataset of ratio 1:100 on taxonomy ID '632' was over 1500 s. However, the time spent training the SVM model was more than 23 000 seconds.

## Further discussion

Given different PPI networks, such as the HB-PPI between *H. sapiens* and *Clostridium botulinum* (ID: 1491), and the interaction between *H. sapiens* and *Yersinia pseudotuberculosis subsp. pestis* (ID: 632), the positive interactions networks have presented different complexities. As we can see, it still requires huge amounts of work towards the completeness of HB-PPIs network. They have indicated different pathways between the different species. [Figures 10](#) and [11](#) show the diagrams of two different interaction networks for taxonomy IDs 1491 and 632, respectively.

To accomplish a robust predictive performance of HB-PPIs, the relationship between the positive and negative protein interactions requires further consideration. There have been several methods dedicated to one-class classification tasks, such as semi-supervised learning [95–97], to leverage the power of singularly labelled data and unlabelled data. This may help to improve the performance of protein interaction prediction regardless of the ratio between the positive and negative protein interactions. Meanwhile, since sequential feature representation algorithms have been an active and challenging area, a better feature representation algorithm is needed to help build a sequence-based end-to-end machine learning model [98–100] for predicting HB-PPIs. In this regard, cutting-edge machine learning algorithms are expected to more effectively decipher the code of protein information, in particular deep learning algorithms such as graph neural networks [101], long short-term memory and convolutional neural network model [102]. It remains a challenge

in regard to how these advanced deep learning techniques can be better leveraged to efficiently distil the useful information and extract informative features from the HP-PPIs networks to further enhance the predictive performance. By benefitting from the advanced machine learning models, the study on PPI networks will eventually shed the lights on our understanding of the mechanisms of infectious disease. Since development of accessible portals for computational analysis and prediction has become a common practice, it is essential to construct web servers [103–108] to support and publish stand-alone tools to enhance the research communication and facilitate future discoveries of HP-PPIs. Given the increasing number of developed models for high-throughput prediction of HP-PPIs, the future work is anticipated to involve the development of user-friendly tools and web servers for HP-PPI evaluation and prediction.

## Conclusions

In this study, we evaluated HB-PPIs in a systematic manner, where the focus was on leveraging machine learning-based models as the primary computational method. We first presented a wide and deep review on currently available data sources and tools. As noted in the literature review (Section 2) of computational tools developed for prediction tasks of HP-PPIs, a careful data curation phase was implemented and a pipeline for HB-PPI studies was summarized, which included numerous sequential feature representation algorithms and machine learning models. Several other computational methods concerning HB-PPIs were also evaluated.

Given the study of HP-PPIs, we have tried to determine the impacts caused by different ratios of benchmark datasets, different feature representation algorithms and different machine learning models. The experimental results indicated that to better utilize machine learning models and harness the power of accumulated protein interaction data, a more robust and more powerful computational model is required to achieve better performance across different HB-PPI prediction tasks. To facilitate the usage and study of HB-PPIs, a complete evaluation report and



Table 6. Performance on taxonomy ID '632'

Feature	Accuracy		F1 score		MCC		1:100	1:50	1:25	1:50	1:100	1:500	1:1000
	1:1		1:1		1:1								
Auto covariance (420D)	0.757082	0.967350	0.982521	0.990674	0.752226	0.312530	0.243679	0.154535	0.514740	0.389241	0.335746	0.253434	0.212697
	±0.008000 (RF)	±0.000365 (RF)	±0.000128 (RF)	±0.000043 (RF)	±0.006632 (RF)	±0.010944 (RF)	±0.012883 (RF)	±0.012569 (RF)	±0.016240 (RF)	±0.010464 (RF)	±0.010207 (RF)	±0.010138 (RF)	±0.018621 (GBM)
Local descriptor (1260D)	0.720963	0.965377	0.981676	0.990444	0.727218	0.255139	0.177423	0.173899	0.442457	0.328314	0.256948	0.233817	0.212697
	±0.016687 (GBM)	±0.000487 (RF)	±0.000091 (RF)	±0.000060 (RF)	±0.013162 (GBM)	±0.009452 (RF)	±0.010255 (DT)	±0.010245 (GBM)	±0.032907 (GBM)	±0.013210 (RF)	±0.008037 (RF)	±0.012761 (GBM)	±0.019050 (SVM)
Conjoint triad method (686D)	0.700283	0.965039	0.981760	0.990523	0.700275	0.185780	0.180318	0.129115	0.400747	0.297864	0.249646	0.219050	0.209724
	±0.010306 (GBM)	±0.000311 (RF)	±0.000208 (SVM)	±0.000051 (SVM)	±0.006187 (GBM)	±0.010755 (RF)	±0.006771 (SVM)	±0.010062 (RF)	±0.020562 (GBM)	±0.012409 (RF)	±0.012510 (RF)	±0.009724 (SVM)	0.212697
PseAAC (110D)	0.718697	0.964374	0.981415	0.990391	0.724976	0.230770	0.148855	0.130497	0.437930	0.270015	0.241630	0.212697	0.212697
	±0.014061 (GBM)	±0.000203 (RF)	±0.000113 (RF)	±0.000145 (GBM)	±0.010361 (GBM)	±0.011551 (SVM)	±0.011666 (GBM)	±0.010625 (GBM)	±0.027559 (GBM)	±0.008225 (RF)	±0.013639 (GBM)	±0.018621 (GBM)	0.212697
Pse-PSSM (80D)	0.709632	0.966216	0.982049	0.990624	0.720259	0.256988	0.191165	0.143488	0.420486	0.348116	0.294669	0.242970	0.242970
	±0.005540 (GBM)	±0.000450 (RF)	±0.000120 (RF)	±0.000044 (RF)	±0.004842 (GBM)	±0.009757 (RF)	±0.012116 (RF)	±0.008093 (RF)	±0.010975 (GBM)	±0.013818 (RF)	±0.010245 (RF)	±0.007387 (RF)	0.242970
DPCPSSM (800D)	0.734278	0.966053	0.982060	0.990585	0.742534	0.259213	0.205636	0.154714	0.469595	0.344422	0.300172	0.240708	0.240708
	±0.009506 (GBM)	±0.000333 (RF)	±0.000208 (RF)	±0.000061 (RF)	±0.008470 (GBM)	±0.010695 (RF)	±0.012626 (RF)	±0.013181 (DT)	±0.018729 (GBM)	±0.010812 (RF)	±0.014612 (RF)	±0.009831 (RF)	0.240708
Block-PSSM (800D)	0.729037	0.965279	0.981698	0.990551	0.739192	0.207103	0.175258	0.157700	0.459515	0.321348	0.274589	0.227935	0.227935
	±0.008095 (GBM)	±0.000464 (RF)	±0.000293 (GBM)	±0.000078 (RF)	±0.007676 (GBM)	±0.009433 (RF)	±0.011638 (GBM)	±0.004793 (GBM)	±0.016168 (GBM)	±0.005442 (SVM)	±0.018411 (RF)	±0.004565 (GBM)	0.227935

databases analysis have released along with this review to the wider biomedical research community.

### Key Points

- A comprehensive review on currently available data sources and computational tools is presented.
- A comprehensive framework for HBI studies was summarized while both the datasets and computational methods were substantially reviewed and collected.
- A systematic evaluation of machine learning-based computational prediction was delivered. Although numerous existing studies have reported the performance of traditional machine learning methods separately, in this study, we evaluated a larger scope of machine learning models as well as feature representation algorithms. The evaluation was conducted by reporting multiple metrics and comparing between the different models.
- By developing the comprehensive pipeline for HBI studies, we have tried to answer the following questions: (a) How do machine learning-based models perform on the prediction task of HB-PPIs? (b) How do feature representation algorithms based on sequence information affect the model performance? (c) Do the ratios between the positive and negative interactions have an impact on the model performance?

### Acknowledgment

This work was supported by grants from the University Global Partnership Network (UGPN), the National Health and Medical Research Council of Australia (NHMRC) (1144652), the Australian Research Council (ARC) (LP110200333 and DP120104460), the National Institute of Allergy and Infectious Diseases of the National Institutes of Health (R01 AI111965), a Major Inter-Disciplinary Research (IDR) project awarded by Monash University, and the Collaborative Research Program of Institute for Chemical Research, Kyoto University (2019-32).

### Conflict of interest

None declared.

### References

1. Prashanthi K, Chandra N. Host-pathogen interactions. In: Dubitzky W, Wolkenhauer O, Cho K-H et al. (eds). *Encyclopedia of Systems Biology*. New York, NY: Springer New York, 2013, 904–8.
2. Mock M, Fouet A. Anthrax. *Annu Rev Microbiol* 2001; 55:647–71.
3. Maresso AW, Garufi G, Schneewind O. *Bacillus anthracis* secretes proteins that mediate heme acquisition from hemoglobin. *PLoS Pathog* 2008;4:e1000132.
4. Dyer MD, Nef C, Dufford M, et al. The human-bacterial pathogen protein interaction networks of *Bacillus anthracis*, *Francisella tularensis*, and *Yersinia pestis*. *PLoS One* 2010;5:e12089.

5. Nourani E, Khunjush F, Durmuş S. Computational approaches for prediction of pathogen-host protein-protein interactions. *Front Microbiol* 2015;6:1–10.
6. Durmus S, Çakir T, Özgür A, et al. A review on computational systems biology of pathogen-host interactions. *Front Microbiol* 2015;6:1–19.
7. Dyer MD, Murali TM, Sobral BW. Supervised learning and prediction of physical interactions between human and HIV proteins, infection. *Genet Evol* 2011;11:917–23.
8. Cui G, Fang C, Han K. Prediction of protein-protein interactions between viruses and human by an SVM model. *BMC Bioinformatics* 2012;13:S5–5.
9. Kshirsagar M, Carbonell J, Klein-Seetharaman J. Multitask learning for host-pathogen protein interactions. *Bioinformatics* 2013;29:217–26.
10. Emamjomeh A, Goliaei B, Zahiri J, et al. Predicting protein-protein interactions between human and hepatitis C virus via an ensemble learning method. *Mol BioSyst* 2014;10:3147–54.
11. Eid FE, Elhefnawi M, Heath LS. DeNovo: virus-host sequence-based protein-protein interaction prediction. *Bioinformatics* 2016;32:1144–50.
12. Sen R, Nayak L, De RK. A review on host-pathogen interactions: classification and prediction. *Eur J Clin Microbiol Infect Dis* 2016;35:1581–99.
13. Zhou H, Jin J, Wong L. Progress in computational studies of host-pathogen interactions. *J Bioinform Comput Biol* 2013;11: 1230001.
14. Zhang J, Kurgan L. Review and comparative assessment of sequence-based predictors of protein-binding residues. *Brief Bioinform* 2017;1–17.
15. Krishnadev O, Srinivasan N. A data integration approach to predict host-pathogen protein-protein interactions: application to recognize protein interactions between human and a malarial parasite. *In Silico Biol* 2008;8: 235–50.
16. Huo T, Liu W, Guo Y, et al. Prediction of host-pathogen protein interactions between mycobacterium tuberculosis and *Homo sapiens* using sequence motifs. *BMC Bioinformatics* 2015;16:1–9.
17. Hwang H, Dey F, Petrey D, et al. Structure-based prediction of ligand-protein interactions on a genome-wide scale. *Proc Natl Acad Sci* 2017;114:13685–90.
18. Dyer MD, Murali TM, Sobral BW. Computational prediction of host-pathogen protein-protein interactions. *Bioinformatics* 2007;23:i159–66.
19. Wuchty S. Computational prediction of host-parasite protein interactions between *P. falciparum* and *H. sapiens*. *PLoS One* 2011;6(e26960):26961–8.
20. Mei S. Probability weighted ensemble transfer learning for predicting interactions between HIV-1 and human proteins. *PLoS One* 2013;8:1–13.
21. Ahmed I, Witbooi P, Christoffels A. Prediction of human-*Bacillus anthracis* protein-protein interactions using multi-layer neural network. *Bioinformatics* 2018;34:4159–64.
22. Driscoll T, Dyer MD, Murali TM, et al. PIG - the pathogen interaction gateway. *Nucleic Acids Res* 2009;37:647–50.
23. Squires B, Macken C, Garcia-Sastre A, et al. BioHealthBase: informatics support in the elucidation of influenza virus host-pathogen interactions and virulence. *Nucleic Acids Res* 2008;36:497–503.
24. Wattam AR, Abraham D, Dalay O, et al. PATRIC, the bacterial bioinformatics database and analysis resource. *Nucleic Acids Res* 2014;42:581–91.

25. Braxton SM, Onstad DW, Dockter DE, et al. Description and analysis of two internet-based databases of insect pathogens: EDWIP and VIDIL. *J Invertebr Pathol* 2003;**83**:185–95.
26. Consortium U. UniProt: the universal protein knowledgebase. *Nucleic Acids Res* 2017;**45**:D158–69.
27. Durmuş Tekir S, Çakır T, Ardiç E, et al. PHISTO: pathogen-host interaction search tool. *Bioinformatics* 2013;**29**:1357–8.
28. Chautard E, Dana JM, Rivas JDL, et al. PSICQUIC and PSIS-CORE: accessing and scoring molecular interactions. *Nat Methods* 2012;**8**:528–9.
29. Xenarios I, Salwinski L, Duan XJ, et al. DIP, the database of interacting proteins: a research tool for studying cellular networks of protein interactions. *Nucleic Acids Res* 2002;**30**:303–5.
30. Joshi-Tope G, Gillespie M, Vastrik I, et al. Reactome: a knowledgebase of biological pathways. *Nucleic Acids Res* 2005;**33**:428–32.
31. Xiang Z, Tian Y, He Y. PHIDIAS: a pathogen-host interaction data integration and analysis system. *Genome Biol* 2007;**8**:R150.
32. Yue J, Zhang D, Ban R, et al. PCPPI: a comprehensive database for the prediction of Penicillium-crop protein-protein interactions. *Database* 2017;**2017**:1–9.
33. Kerrien S, Aranda B, Breuza L, et al. The IntAct molecular interaction database in 2012. *Nucleic Acids Res* 2012;**40**:841–6.
34. Calderone A, Castagnoli L, Cesareni G. Mentha: a resource for browsing integrated protein-interaction networks. *Nat Methods* 2013;**10**:690–1.
35. Shen J, Zhang J, Luo X, et al. Predicting protein-protein interactions based only on sequences information. *Proc Natl Acad Sci* 2007;**104**:4337–41.
36. Guo Y, Yu L, Wen Z, et al. Using support vector machine combined with auto covariance to predict protein-protein interactions from protein sequences. *Nucleic Acids Res* 2008;**36**:3025–30.
37. Davies MN, Secker A, Freitas AA, et al. Optimizing amino acid groupings for GPCR classification. *Bioinformatics* 2008;**24**:1980–6.
38. Chou K-C. Prediction of protein cellular attributes using pseudo-amino acid composition. *Proteins: Struct Funct Genet* 2001;**43**:246–55.
39. Shen HB, Chou KC. PseAAC: a flexible web server for generating various kinds of protein pseudo amino acid composition. *Anal Biochem* 2008;**373**:386–8.
40. Chou K-C. Pseudo amino acid composition and its applications in bioinformatics, proteomics and system biology. *Curr Proteom* 2009;**6**:262–74.
41. Altschul SF, Madden TL, Schäffer AA, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997;**25**:3389–402.
42. Ahmad S, Sarai A. PSSM-based prediction of DNA binding sites in proteins. *BMC Bioinformatics* 2005;**6**:1–6.
43. Xia J-F, Han K, Huang D-S. Sequence-based prediction of protein-protein interactions by means of rotation forest and autocorrelation descriptor. *Protein Pept Lett* 2010;**17**:137–45.
44. Chou KC, Shen HB. MemType-2L: a web server for predicting membrane proteins and their types by incorporating evolution information through Pse-PSSM. *Biochem Biophys Res Commun* 2007;**360**:339–45.
45. Chou KC. Some remarks on protein attribute prediction and pseudo amino acid composition. *J Theor Biol* 2011;**273**:236–47.
46. Chen Z, Zhao P, Li F, et al. iLearn: an integrated platform and meta-learner for feature engineering, machine-learning analysis and modeling of DNA, RNA and protein sequence data. *Brief Bioinform* 2019;**10**.
47. Chen Z, Zhao P, Li F, et al. iFeature: a python package and web server for features extraction and selection from protein and peptide sequences. *Bioinformatics* 2018;**1**:1–4.
48. Zahiri J, Yaghoubi O, Mohammad-Noori M, et al. PPIevo: protein-protein interaction prediction from PSSM based evolutionary information. *Genomics* 2013;**102**:237–42.
49. Wang J, Yang B, Song J. Bastion6: a bioinformatics approach for accurate prediction of type VI secreted effectors. *Bioinformatics* 2018;**34**:2546–55.
50. Uddin MR, Sharma A, Farid DM, et al. EvoStruct-sub: an accurate gram-positive protein subcellular localization predictor using evolutionary and structural features. *J Theor Biol* 2018;**443**:138–46.
51. Göktepe YE, Kodaz H. Prediction of protein-protein interactions using an effective sequence based combined method. *Neurocomputing* 2018;**303**:68–74.
52. Zhang B, Li J, Lü Q. Prediction of 8-state protein secondary structures by a novel deep learning architecture. *BMC Bioinformatics* 2018;**19**:1–13.
53. Wang Y-B, You Z-H, Li L-P, et al. Improving prediction of self-interacting proteins using stacked sparse auto-encoder with PSSM profiles. *Int J Biol Sci* 2018;**14**:983–91.
54. Dayhoff MO. A model of evolutionary change in proteins. *Atlas Protein Seq Struct* 1972;**5**:89–99.
55. Henikoff S, Henikoff JG. Amino acid substitution matrices from protein blocks. *Proc Natl Acad Sci* 1992;**89**:10915–9.
56. Jeong JC, Lin X, Chen XW. On position-specific scoring matrix for protein function prediction. *IEEE/ACM Trans Comput Biol Bioinform* 2011;**8**:308–15.
57. Liu T, Zheng X, Wang J. Prediction of protein structural class for low-similarity sequences using support vector machine and PSI-BLAST profile. *Biochimie* 2010;**92**:1330–4.
58. Halder AK, Dutta P, Kundu M, et al. Review of computational methods for virus-host protein interaction prediction: a case study on novel Ebola-human interactions. *Brief Funct Genom* 2018;**17**:381–91.
59. Arnold R, Boonen K, Sun MGF, et al. Computational analysis of interactomes: current and future perspectives for bioinformatics approaches to model the host-pathogen interaction space. *Methods* 2012;**57**:508–18.
60. Cortes C, Vapnik V. Support-vector networks. *Mach Learn* 1995;**20**:273–97.
61. Safavian SR, Landgrebe D. A survey of decision tree classifier methodology. *IEEE Trans Syst Man Cybernet* 1991;**21**:660–74.
62. Breiman L. Random forests. *Mach Learn* 2001;**45**:5–32.
63. Li F, Li C, Revote J, et al. GlycoMine struct: a new bioinformatics tool for highly accurate mapping of the human N-linked and O-linked glycoproteomes by incorporating structural features. *Sci Rep* 2016;**6**:1–16.
64. Li F, Li C, Wang M, et al. GlycoMine: a machine learning-based approach for predicting N-, C- and O-linked glycosylation in the human proteome. *Bioinformatics* 2015;**31**:1411–9.

65. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: machine learning in python. *J Mach Learn Res* 2011;12:2825–30.
66. Song J, Li F, Leier A, et al. PROSPEROus: high-throughput prediction of substrate cleavage sites for 90 proteases with improved accuracy. *Bioinformatics* 2018;34:684–7.
67. Li F, Li C, Marquez-Lago TT, et al. Quokka: a comprehensive tool for rapid and accurate prediction of kinase family-specific phosphorylation sites in the human proteome. *Bioinformatics* 2018;34:4223–31.
68. Lewis DD. Naive (Bayes) at forty: The independence assumption in information retrieval. In: *European Conference on Machine Learning*, 1998, p. 4–15. Springer.
69. Zhang H. The optimality of naive Bayes. In: *The 17th International FLAIRS Conference*. Miami Beach, Florida, USA, 2004, p. 562–7.
70. Friedman JH. Greedy function approximation : a gradient boosting machine. *Ann Stat* 2001;29:1189–232.
71. Prieto C, De Las Rivas J. APID: agile protein interaction DataAnalyzer. *Nucleic Acids Res* 2006;34:298–302.
72. Licata L, Briganti L, Peluso D, et al. MINT, the molecular interaction database: 2012 update. *Nucleic Acids Res* 2012;40:857–61.
73. Breuer K, Foroushani AK, Laird MR, et al. InnateDB: systems biology of innate immunity and beyond - recent updates and continuing curation. *Nucleic Acids Res* 2013;41:1228–33.
74. Kumar R, Nanduri B. HPIDB - a unified resource for host-pathogen interactions. *BMC Bioinformatics* 2010;11:S16.
75. Ammari MG, Gresham CR, McCarthy FM, et al. HPIDB 2.0: a curated database for host-pathogen interactions. *Database : J Biol Database Curat* 2016;2016:1–9.
76. Chatr-Aryamontri A, Oughtred R, Boucher L, et al. The BioGRID interaction database: 2017 update. *Nucleic Acids Res* 2017;45:D369–79.
77. Boutet E, Lieberherr D, Tognolli M, et al. Uniprotkb/swiss-prot. *Plant Bioinform Springer* 2007;89–112.
78. Davis FP, Barkan DT, Eswar N, et al. Host pathogen protein interactions predicted by comparative modeling. *Protein Science : Publ Prot Soc* 2007;16:2585–96.
79. Mariano R, Wuchty S. Structure-based prediction of host-pathogen protein interactions. *Curr Opin Struct Biol* 2017;44:119–24.
80. Franzosa EA, Xia Y. Structural principles within the human-virus protein-protein interaction network. *Proc Natl Acad Sci* 2011;108:10538–43.
81. Franzosa EA, Garamszegi S, Xia Y. Toward a three-dimensional view of protein networks between species. *Front Microbiol* 2012;3:1–6.
82. Qi Y, Tastan O, Carbonell JG, et al. Semi-supervised multi-task learning for predicting interactions between HIV-1 and human proteins. *Bioinformatics* 2010;26:i645–52.
83. Tastan O, Qi Y, Carbonell JG, et al. Prediction of interactions between HIV-1 and human proteins by information integration. *Biocomputing 2009 World Scientific* 2009;516–27.
84. Tyagi N, Krishnadev O, Srinivasan N. Prediction of protein-protein interactions between *Helicobacter pylori* and a human host. *Mol BioSyst* 2009;5:1630–5.
85. Gomez SM, Noble WS, Rzhetsky A. Learning to predict protein-protein interactions from protein sequences. *Bioinformatics* 2003;19:1875–81.
86. Zhang L. Sequence-based prediction of protein-protein interactions using random tree and genetic algorithm. *Intell Comput Technol* 2012;334–41.
87. Yang S, Li H, He H, et al. Critical assessment and performance improvement of plant-pathogen protein-protein interaction prediction methods. *Brief Bioinform* 2017;1–11.
88. Mei S, Li F, Leier A, et al. A comprehensive review and performance evaluation of bioinformatics tools for HLA class I peptide-binding prediction. *Brief Bioinform* 2019;bbz051:051–17.
89. Li F, Zhang Y, Purcell AW, et al. Positive-unlabelled learning of glycosylation sites in the human proteome. *BMC Bioinformatics* 2019;1–17.
90. Zhang M, Li F, Marquez-Lago TT, et al. MULTiPLY: a novel multi-layer predictor for discovering general and specific types of promoters. *Bioinformatics* 2019;35:2957–65.
91. Chen Z, Liu X, Li F, et al. Large-scale comparative assessment of computational predictors for lysine post-translational modification sites. *Brief Bioinform* 2019;20:2267–90.
92. Li F, Wang Y, Li C, et al. Twenty years of bioinformatics research for protease-specific substrate and cleavage site prediction: a comprehensive revisit and benchmarking of existing methods. *Brief Bioinform* 2019;20:2150–66.
93. Song J, Wang Y, Li F, et al. iProt-sub: a comprehensive package for accurately mapping and predicting protease-specific substrates and cleavage sites. *Brief Bioinform* 2019;20:638–58.
94. Li F, Chen J, Leier A, et al. DeepCleave: a deep learning predictor for caspase and matrix metalloprotease substrates and cleavage sites. *Bioinformatics* 2020;36:1057–65.
95. Manevitz LM, Yousef M. One-class SVMs for document classification. *J Mach Learn Res* 2001;2:139–54.
96. Chidlovskii B, Hovelynck M. Multi-modality classification for one-class classification in social networks. United States patent US 8,386,574. 2013.
97. Ruff L, Vandermeulen R, Goernitz N et al. Deep one-class classification. In: *International Conference on Machine Learning*. 2018, p. 4393–402.
98. Min S, Lee B, Yoon S. Deep learning in bioinformatics. *Brief Bioinform* 2017;18:851–69.
99. Perera P, Patel VM. Learning deep features for one-class classification. *IEEE Trans Image Process* 2019;28:5450–63.
100. Greener JG, Kandathil SM, Jones DT. Deep learning extends de novo protein modelling coverage of genomes using iteratively predicted structural constraints. *Nat Commun* 2019;10:1–13.
101. Zhang C, Song D, Huang C et al. Heterogeneous graph neural network. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2019, p. 793–803.
102. Hanson J, Paliwal K, Litfin T, et al. Improving prediction of protein secondary structure, backbone angles, solvent accessibility and contact numbers by using predicted contact maps and an ensemble of recurrent and residual convolutional neural networks. *Bioinformatics* 2019;35:2403–10.
103. Li F, Fan C, Marquez-Lago TT, et al. PRISMOID: a comprehensive 3D structure database for post-translational modifications and mutations with functional impact. *Brief Bioinformatics* 2019. doi: 10.1093/bib/bbz050.
104. Hong J, Luo Y, Mou M, et al. Convolutional neural network-based annotation of bacterial type IV secretion system effectors with enhanced accuracy and reduced false discovery. *Brief Bioinformatics* 2019. doi: 10.1093/bib/bbz120.



105. Hong J, Luo Y, Zhang Y, et al. Protein functional annotation of simultaneously improved stability, accuracy and false discovery rate achieved by a sequence-based deep learning. *Brief Bioinformatics* 2019. doi: [10.1093/bib/bbz081](https://doi.org/10.1093/bib/bbz081).
106. Tang J, Wang Y, Fu J, et al. A critical assessment of the feature selection methods used for biomarker discovery in current metaproteomics studies. *Brief Bioinformatics* 2019. doi: [10.1093/bib/bbz061](https://doi.org/10.1093/bib/bbz061).
107. Lian X, Yang S, Li H, et al. Machine-learning-based predictor of human-bacteria protein-protein interactions by incorporating comprehensive host-network properties. *J Proteome Res* 2019;18:2195–205.
108. Yang S, Fu C, Lian X, et al. Understanding human-virus protein-protein interactions using a human protein complex-based analysis framework. *MSystems* 2019;4:e00303–18.