

# Mining protein–protein interaction information on the internet<sup>☆</sup>

Hsi-Chieh Lee<sup>a,b</sup>, Szu-Wei Huang<sup>a</sup>, Eldon Y. Li<sup>c,d,\*</sup>

<sup>a</sup> Department of Information Management, Yuan Ze University, Taoyuan 320, Taiwan, ROC

<sup>b</sup> Department of Information Management, National Kinmen Institute of Technology, Jinning, Kinmen 892, Taiwan, ROC

<sup>c</sup> Department of Management Information Systems, National Chengchi University, Taipei 116, Taiwan, ROC

<sup>d</sup> Orfalea College of Business, California Polytechnic State University, San Luis Obispo, CA 93407, USA

## Abstract

In this study, a mining system is proposed for finding protein–protein interaction literatures from the databases on the Internet. In this system, we find out discriminating words for protein–protein interaction by way of statistics and the results from literatures. A threshold is also evaluated to check if a given literature is related to protein–protein interactions. In addition, a keypage-based search mechanism is used to find related papers for protein–protein interactions from a given document. To expand the search space and ensure better performance of the system, mechanisms for protein name identification and databases for protein names are also developed.

The system is designed with a web-based user interface and a job-dispatching kernel. Experiments are conducted and the results have been checked by a biomedical expert. The experimental results indicate that by using the proposed mining system, it is helpful for researchers to find out protein–protein literatures from the overwhelming piece of information available on the biomedical databases on the Internet.

© 2005 Published by Elsevier Ltd.

**Keywords:** Protein–protein interactions; Rule-based system; Data mining; Information retrieval; Keypage-based search

## 1. Introduction

In recent years, due to advances in information technology, more and more biomedical-related information is available electronically on the Internet. For example, the MEDLINE database contains over 12 million citations dating back to the mid-1960s. Therefore, it has become an important issue for mining valuable biomedical information from the literature (Valencia-García, Ruiz-Sánchez, Vicente, Fernández-Breis, & Martínez-Béjar, 2004; Wang, Kuo, Chen, Hsiao, & Tsai, 2005), especially information on the Internet (Hong & Han, 2002). Expert systems and data mining techniques have been used for years in medical diagnosis domain (Alonso, Caraçá-Valente, González, & Montes, 2002; Chou, Lee, Shao, & Chen, 2004).

In the post-genomic era, some scientists focus on finding meaningful information of DNA or try to use the information

of gene sequence in solving problems. However, the spirit of post-genomic era can be view broadly in three ways. The first one is the sequences from DNA level, and the second one is the Expressed Sequence Tag (EST) from RNA level. The last one is proteome from the protein level. People can use the analyzed information to understand the interaction between each other and discover the meaning behind it. In other words, after decoding the sequence, scientist can analyze the interaction between gene and protein, and understand the role the gene is playing on an organism. It has been shown that the protein and genomics would become the main issue in the post-genomic era (Eisenberg, Marcott, Xenarios, & Yeates, 2000).

Moreover, scientists try to understand the interaction and relation between proteins from biochemistry and gene-related angles. For example, the Database of Interaction Proteins (DIP) developed in UCLA (Xenarios, Rice, Salwinski, Baron, Marcotte and Eisenberg, 2000) has data about over 5900 proteins and 10,500 protein–protein interactions. Besides DIP, there exist many other databases with the collection of the data regarding protein function and pathway. However, if people want to know the relationship between proteins, they have to search different literatures and try to find some relationships. It is considered mission impossible to check on MEDLINE manually where there exist more than 15 million biomedical citations. It is time-consuming and ineffective. It would be helpful if the job can be processed automatically and

<sup>☆</sup> This research is supported partially under the grants of NSC91-2745-P-155-003 and NSC93-2745-E-155-010-URD, Taiwan.

\* Corresponding author. Address: Orfalea College of Business, California Polytechnic State University, San Luis Obispo, CA 93407, USA.

E-mail addresses: imhlee@saturn.yzu.edu.tw (H.-C. Lee), s907736@mail.yzu.edu.tw (S.-W. Huang), eli@calpoly.edu (E.Y. Li).

the database can be updated as soon as new literatures are available.

Generally, mining the literatures of protein–protein interactions requires natural language processing. The literature discussing protein–protein interactions does not contain a language that a computer can understand. As a result, there are two typical approaches in solving the problem. The first approach is transferring the format into a way that computer can understand by natural language processing. For example, in Ono, Hishigaki, Tanigami, and Takagi (2001), they brought up an idea to extract biomedical-related information with two steps. The first step is to scan the full document with a protein name dictionary. The second step is to extract content related to protein–protein interactions by predefined rules. The second approach is to extract biomedical-related information using statistics. The most typical way of statistics is calculating the frequency of words. In Marcotte, Xenarios, and Eisenberg (2001), they used statistics to find 83 words as discriminative words to check whether a paper is discussing protein–protein interaction.

In this study, we integrate both approaches mentioned above. The natural language processing techniques and statistics were utilized to process the biomedical-related literatures. Besides checking whether a paper is discussing protein–protein interaction, we also find the probable protein names in the target document and help people in finding other related literatures. The proposed mining system is described in detail in Section 2 and the experimental results are illustrated subsequently.

## 2. The proposed system

The literature mining procedures of the proposed mining system is illustrated in Fig. 1. The system starts by feeding a target document (e.g. a protein–protein interaction paper) to the system. The first step of the mining system is to calculate

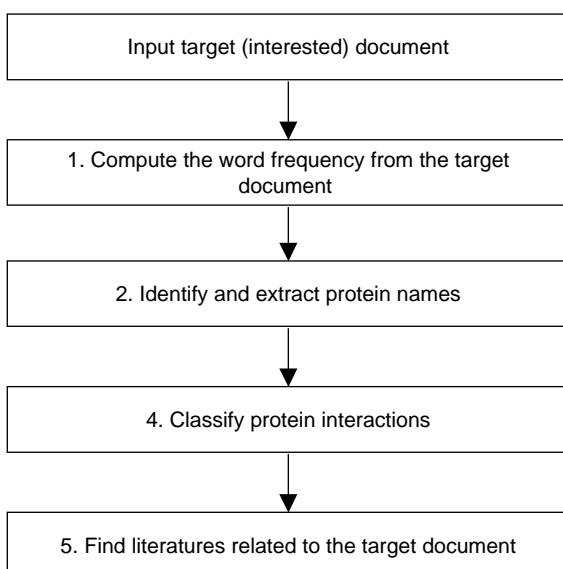


Fig. 1. The literature mining procedures of the proposed system.

the frequency of words and decide whether the target document can be classified as a paper related to protein–protein interaction. Next, the mining system will identify the potential protein names in the target document. The third step is to find out the category of interaction of the target document according to the frequencies of words in it. The final step is to compose keywords automatically and send the query to the PubMed. Once the query results are ready, the system will then retrieve the related literatures from the Internet (<http://www.ncbi.nlm.nih.gov/PubMed/>). The detailed description of each step is shown below.

**Step 1.** Compute the frequency of the word tokens in the target (interested) document. In the proposed mining system, we used 20 words that have the most discriminate capability. To ensure that they are the most discriminate words, we retrieved 2203 abstracts from Database of Interacting Proteins (<http://dip.doe-mbi.ucla.edu/>) (Xenarios et al., 2000)—a well-known database dedicated for data about protein interaction. We calculated the frequency of the words in these documents and took them as the most discriminating words (stop words were excluded). When the researcher inputs the target document, the system will calculate the frequencies of the 20 words. If the frequencies are high enough, this target document is considered as discussing protein interactions by the proposed system. In a different study, Marcotte et al. (2001) found statistics of 65,807 literatures and derived 83 discriminate words.

**Step 2.** *Identify and extract protein names.* In biology domain, it is not uncommon to have new findings everyday. When there is a new finding, the researcher will give it a new name to discriminate it from others. As a result, new protein names are created rapidly which makes it difficult to identify all the protein names. In general, a protein name can be classified into one of the following types of word.

- Single words with upper case letters, numerical figures, and non-alphabetical letters. Mostly derived from gene name.
- Compound words with upper case letters, numerical letters, and non-alphabetical letters.
- Single word with only lower case letters.

The first two types of protein names are more commonly found in research papers, while protein names of type 3 are rarely found (Fukuda, Tsunoda, Tamura, & Takagi, 1998).

Another difficulty of identifying a protein name is the inconsistency of naming convention. It is also not uncommon to find the same protein being expressed differently by different researchers, depending on the researchers' styles. For example, the term 'epidermal growth factor receptor', some researchers describe it directly while others use abbreviations as 'EGF receptor' or 'EGFR'. That is, even the abbreviation might be expressed differently. The following example demonstrates some variation of expression for a protein function (Fukuda et al., 1998).

- the Ras guanine nucleotide exchange factor Sos
- the Ras guanine nucleotide releasing protein Sos
- the Ras exchanger Sos

As mentioned above, protein names and their expressions are mainly decided by researchers' style. Therefore, it becomes a challenge to find them in the literature. The general way of finding protein names is constructing a dictionary of protein names manually and uses the dictionary for pattern matching (Thomas, Milward, Ouzounis, Pulman, & Carroll, 2000). Another way of finding the protein names is to check the words that are used often around the literature. In addition, it is also applicable to analyze the characteristics of words and try to find protein names in the literature. In the proposed system, we integrate the three approaches mentioned above for identifying the protein names in the literature. The detailed process of extracting protein names is shown in Fig. 2 and the discussion follows.

1. To construct a dictionary of protein names for the system, we retrieve the protein names from Protein Name Abbreviation Dictionary (<http://pnad.ontology.ims.u-tokyo.ac.jp/search/php/search.php>) and then add other known protein names manually.
2. Using pattern-matching approach to compare every word in the literature and extract the protein names from the literature.
3. A set of word patterns is commonly used for recognition of protein–protein interaction, namely, 'interact', 'bind', 'associate', and 'complex'. Unfortunately, they might appear in different forms. As a result, we try to extract these word patterns using regular expression and the rule shown below. In addition, the Brill POS tagger package (Brill, 1994) is used to analyze the target text.

*Rule:* If the sentence matches the following part-of-speech pattern as indicated by regular expression, it can be extracted.

Table 1  
Definition of symbols

Symbol	Definition
,	Comma
:	Colon
;	Semi-colon
CC	Coordinating conjunction
DT	Determiner
IN	Preposition or subordinating conjunction
JJ	Adjective
NN	Noun, singular or mass
NNP	Proper noun, singular
NNS	Noun, plural
P(1/2)	Phrase
P(3/4/5)	Phrase without verb
VB(1/2)	Verb
VBN	Verb, past participle
VBZ	Verb, third person singular present

`interact.* | complex.* | bind.* | associate.*`

When getting the sentence from the target document, the Brill POS tagger extracts the words which are tagged as NN or NNP (see Table 1). Once the sentence contains word patterns of protein–protein interaction, we treat the words that are tagged as NN or NNP in the sentence to be potential protein name. To further eliminate words that are not related, for example, words like 'domain', 'function', the Porter's Stemming Algorithm (Porter, 1980) is used. This algorithm processes suffix of words. If the word is plurality or other type, it will be transformed back to the original word. We will try to compare the words extracted to those on WordNet and try to eliminate words that are not interested. The Brill POS tagger definition of symbols is shown in Table 1.

4. In this step, we adopt Fukuda's (1998) approach to extracting protein names with the following characteristics.
  - (a) Based on characteristics of protein name, extract the words with upper cases, numerical figures, and/or special symbols.
  - (b) Use the rules listed below to filter out improbable words.
    - Exclude words whose length is more than nine characters and consists of '-' and lower cases.
    - Exclude words in which more than half of its character string consists of special symbols.
    - Exclude words related to numbers such as units. Eight words (aa, AA, fold, bp, nM, microM, %, UV) are registered as units.
    - Exclude words that agree to the reference template prepared beforehand.
  - (c) If extracting words are in common use, they can be filtered out by words in WordNet.
  - (d) The words that are left over are the probable protein names.
5. Integrating the probable protein names found using three different approaches: comparing with database, extracting by regular expression, and extracting by characteristics of protein names.

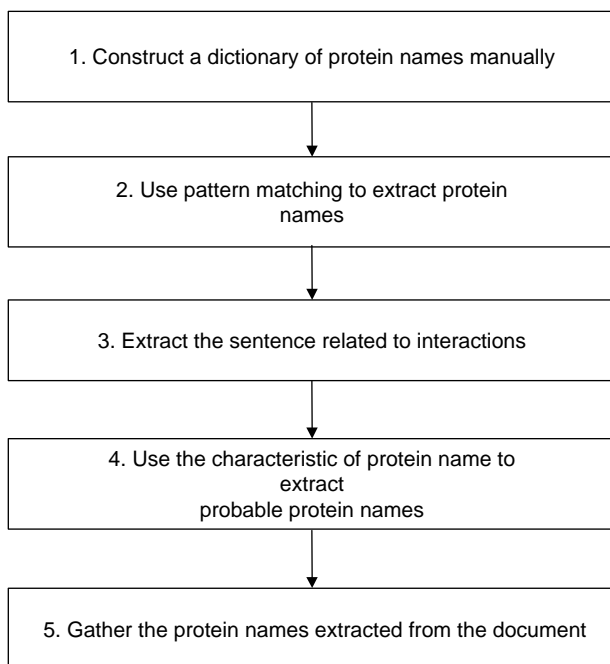


Fig. 2. Identification and extraction of protein names from a document.

*Step 3. Classification of protein interactions.* After extracting the protein names, we further try to classify the interaction discussed in the target document into a category. Based on the words commonly used to describe interaction and by checking the words that are discriminated by the statistics, four words (interact, associate, complex, bind) were used. According to the frequencies of the four words, the highest one will be considered as the category of interaction that is discussed in the target document.

*Step 4. Finding related literature of the target document.* Once the protein names and the category of interaction are identified, they will be used as input for search on the search engine of PubMed (URL: <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=PubMed>) for retrieving the candidate research papers related to the target document which has been entered into the proposed mining system by the researcher. The results from PubMed are URLs of candidate research papers. The system will then retrieve the papers from the Internet and then extracting features from each of them and calculate the similarity with the target document using SimNet's kernel function (Lee, Dagli, Ercal, & Ozbayoglu, 1995). SimNet is a neuro-fuzzy system that integrates the neural network architecture and fuzzy theory. The idea of SimNet was first introduced for character recognition and has later been used for speaker identification. It combines a neural network structure with the subethood concept of fuzzy logic to produce a rapid data clustering system that works similar to Adaptive Resonance Theory and Self-Organizing Maps. It has two neural network architectures (Lee et al., 1995), namely, a two-layer unsupervised learning model and a three-layer supervised learning model. In this study, only the similarity measurement equation is used to compare the similarity of documents. We use the algorithm to calculate similarity of documents and get the most related documents. The similarity measurement equation of SimNet is shown below:

$$MD(I, W) = \sqrt{\frac{(\sum_i \min(I_i, W_i))^2}{(\sum_i W_i)(\sum_i I_i) + \varepsilon}}$$

where  $\varepsilon$  is a non-zero minimal.

This procedure turns the keyword-based search into a keypage-based search that reduces the amount of un-related information significantly.

### 3. Experimental results

In order to decide whether the target document fell into the category that discussed protein–protein interaction, we started by finding the 20 most discriminating words. We downloaded 2203 abstracts from the DIP and then calculated the word frequencies. The average number of words in an abstract was 196 words. The top 20 most frequently used words were found as shown in Table 2. As we could see in Table 2, the word patterns, {interact, binding, complex} that were commonly used in protein–protein interaction papers were also included.

Since all the literatures in DIP were considered to have the topic regarding the protein–protein interaction. Therefore, we

Table 2  
The 20 most discriminating words

Discriminating word	Frequency	Discriminating word	Frequency
Complex/complexes	1845	Binding	1221
Interact/interaction/interactions/interacts	2449	Component	249
Two-hybrid	448	Required	626
With	3997	Suggesting	256
Protein/proteins	4434	From	907
Function	607	Demonstrate	251
Domain	1039	Kinase	803
		Essential	444

tried to analyze the behavior of the literatures in DIP by calculating the Matching Degree of each abstract in DIP with the overall frequency vector as shown in Table 2. The Matching Degree (MD) is the kernel function in SimNet (Lee et al., 1995). The similarity measurement equation MD is shown below:

$$MD(I, W) = \sqrt{\frac{(\sum_i \min(I_i, W_i))^2}{(\sum_i W_i)(\sum_i I_i) + \varepsilon}}$$

where  $\varepsilon$  is a non-zero minimal

In this equation, the vector ' $I$ ' means the word frequency that is shown in Table 2 and notice that it has been normalized. The vector ' $W$ ' means the frequency of 20 discriminate words of each abstract found in DIP. The detailed calculating steps are demonstrated in Fig. 3.

The experiment found the Matching Degrees (MDs) of the 2203 abstracts from DIP. The experimental results shown in Fig. 4 reveal that the values of MD were in the range between 0.02 and 0.74 and the average value of MD was 0.457. If the threshold was set to 0.02, all the literatures would be accepted as papers of protein–protein interaction and the recall rate would be maximized. However, the precision rate did not appear to be good. When we set the threshold to the average value 0.457, 45% of the literatures would be accepted and 55% would be dropped.

To generalize the results for the MD threshold and to find the precision and recall rate with different thresholds, we randomly chose 535 literatures from DIP that represented the research papers of protein–protein interaction. Moreover, we used the four keywords, 'interact', 'associate', 'bind', and 'complex' as inputs to MEDLINE. There were 30,277 abstracts obtained by the MEDLINE search using these keywords as the MeSH terms. With a random selection and a duplication–elimination process, 10,531 abstracts were retrieved from MEDLINE. The MDs were calculated from a total of 11,066 abstracts, in which, 535 abstracts derived from DIP and 10,531 abstracts derived from MEDLINE. The detailed calculating steps were demonstrated in Fig. 3 and the MDs of these abstracts were shown in Fig. 5. It could be observed that the test results of the 11,066 abstracts shown in Fig. 5 looked similar to the test results of the 2203 abstracts in Fig. 4.

To test if the MD thresholds could be applied to those literatures that were not included in the set of the 11,066

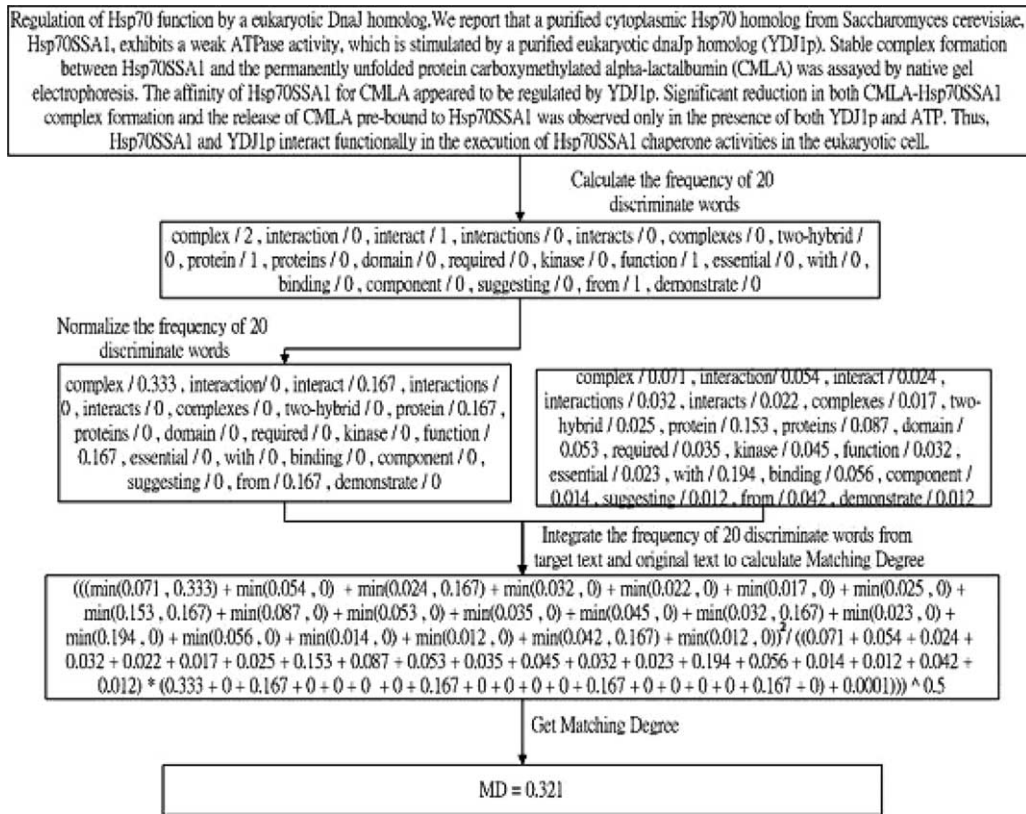


Fig. 3. The steps of calculating the Matching Degrees.

abstracts, a new experiment was conducted and the results for the precision and recall rate were verified by an expert of protein chemistry. This expert was a physician from a local hospital who held a PhD degree from an institution in the US and specialized in Protein Chemistry, Matrix Biology, and Glycobiology. A total of 740 abstracts were used for calculating the Matching Degree. Out of the 740 abstracts, 100 of them were selected randomly from DIP, and the rest of them were selected randomly from MEDLINE that were considered not literatures of protein–protein interaction. The experimental results of the 740 abstracts were shown in Fig. 6.

The experimental results for the 740 abstracts were transformed into the precision-recall chart as shown in Fig. 7.

When the MD threshold was set to 0.35, 87% abstracts would be correctly classified as papers of protein–protein interaction. However, the precision rate was not as high.

There was a tradeoff between the recall rate and the precision as shown in Fig. 7. In other words, to get higher precision, the MD threshold must be higher than 0.45. Conversely, to get higher recall rate, the MD threshold must be lower than 0.35. The selection of MD threshold mainly depended on the purpose of using the proposed mining system. For example, to retrieve more protein–protein interaction papers that were related to the target document, a lower threshold was preferred. In contrast, higher MD threshold was suggested if less but accurate results were demanded. Note that

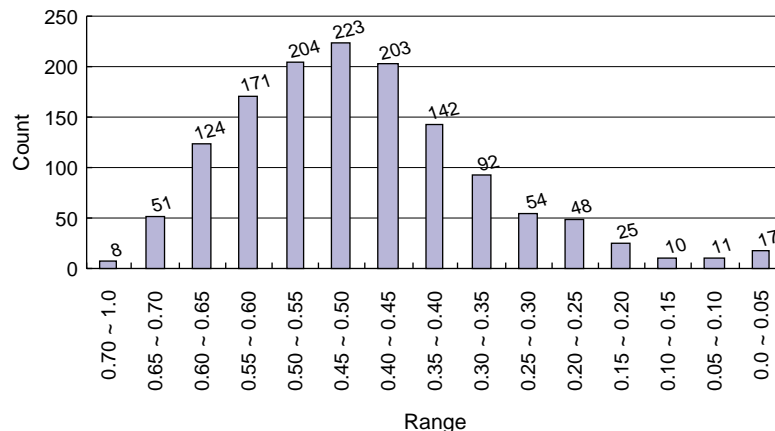


Fig. 4. The Matching Degree (MD) of literatures in DIP.

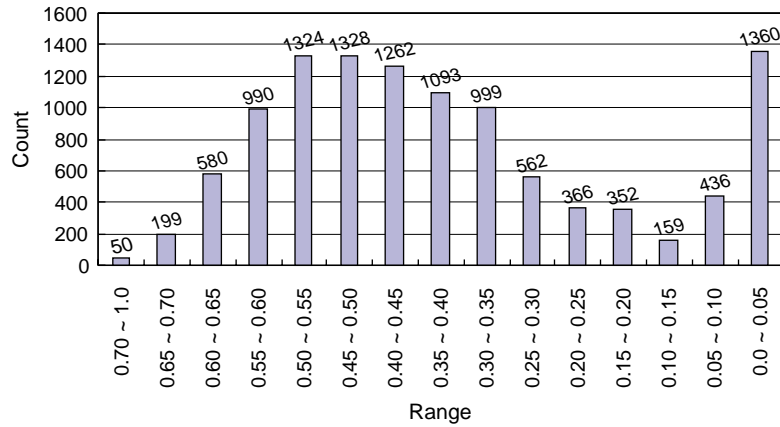


Fig. 5. The Matching Degree (MD) of 11066 abstracts.

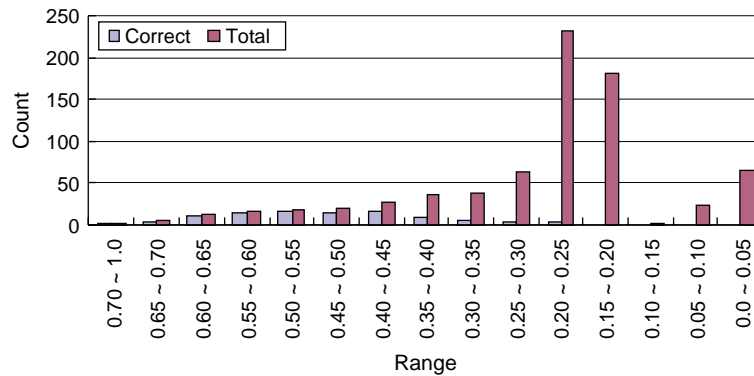


Fig. 6. The Matching Degree (MD) of 740 abstracts.

the mining job requested in this study was a non-trivial computation-intensive and communication-demanding task. The main reason was that the proposed system needed to analyze the target document in order to find the probable protein names and then send the query to the Internet search engines. In addition, each and every document found should also be retrieved from the Internet and the analysis process would be applied to each and every document again. Finally, computation for MDs between the target document and the documents retrieved from the Internet should be done to find the similarity. Fortunately, this time-consuming process could be resolved using a parallel and distributed environment proposed in [Chen & Lee \(2002\)](#).

#### 4. Managerial implications to bioinformatics

As human’s genetic disease has been proven to relate to gene variation, genomic sequence has been used for searching the genetic sequence of disease, researching the way of prevention, diagnosis, and treatment. It also helps people in understanding the relation between environment and evolution by comparing cross-gene of species. Besides analyzing the gene, the analysis of protein is another important approach. Moving the research from the one dimension of DNA into the three dimension of protein, people can understand the protein’s

structure and function. We can also use this information to develop new medicine or new prevention mechanism for diseases. During the course of this research, we must understand the interaction and relation between proteins in biochemistry and gene-related ways. In order to conserve

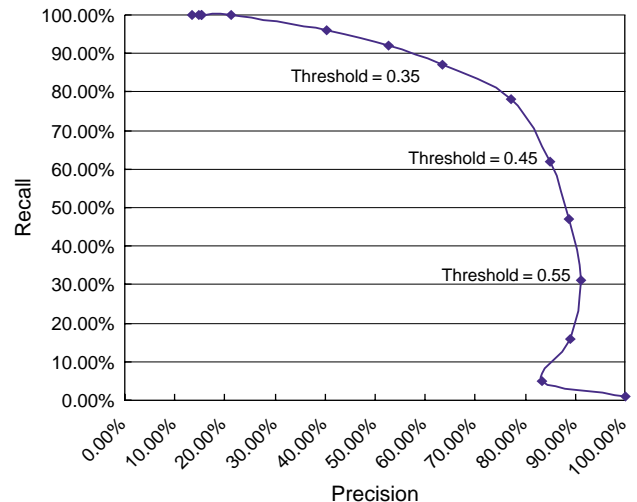


Fig. 7. Precision-recall chart with different thresholds.

the resources demanded by experimental studies, we can find the related experimental outcomes published in the literature. However, the information available on the extant literature is overwhelming. For example, just the MEDLINE (see <http://www.ncbi.nlm.nih.gov/Literature/>) alone contains over 15,000,000 biomedical journal citations. Likewise, Database of Interaction Proteins (Xenarios et al., 2000) alone contains over 5900 proteins and 10,500 interactions data about protein–protein interaction. It is impossible for a human being to search through all these pieces of data. Therefore, our proposed mining system will help researcher in finding related papers of protein interactions from the literatures available over the Internet, with a simple input of an interested document. The proposed mining system is expected to reduce the time and effort significantly for researchers in the stage of literature review.

## 5. Conclusions

In this study, a mining system has been proposed for finding protein–protein interaction literatures from the database on the Internet. Through this system, we found discriminating words for protein–protein interaction by way of statistics and results from literatures. A threshold was also evaluated using the MD function of SimNet to check if a given literature was related to protein–protein interactions. In addition, a keypage-based search mechanism was used to find papers related to protein–protein interactions from a given document. Moreover, to expand the search space and ensure better performance of the system, both mechanisms for protein name identification and databases for protein names were used. Experiments were conducted and the results were checked by an expert of protein chemistry.

The findings of the study suggest that to get precision as high as 90%, the MD threshold can be adjusted to a value higher than 0.45. In addition, to get recall rate as high as 90%, the MD threshold can be adjusted to be lower than 0.35. These indicate that by using the proposed mining system, it is helpful for researchers to find protein–protein literatures from the overwhelming amount of information available on the biomedical database over the Internet. Moreover, because of the nature of the architecture, the results can be derived in a reasonable time with the help of the parallel and distributed mechanism.

Finding related a literature of protein interactions from the Internet using simply an interested document is the major contribution of this study. It keeps the researchers from inputting many possible keywords and check hundreds and thousands of results returned from general search engines. However, there are two major concerns in this research which needs further investigation and study in the future. The first one is the prediction of protein names. As long as there exists no ‘golden’ standard for protein naming convention, there exists the problem of ‘guessing’ protein name available in the literatures. The second concern for implementing such

a mining system is the performance of the system which is usually measured by precision and recall rate. To measure the performance, one may build a benchmarking database in any target domain (e.g. DIP for protein interactions) which is a job requires considerably involvement of domain experts.

## Acknowledgements

The authors would like to thank Dr Wei-Chung Vivian Yang, a biomedical expert in the area of protein chemistry, matrix biology, and glycobiology from Taipei Medical University, for her professional assistance on verifying the results generated by the system.

## References

- Alonso, F., Caraça-Valente, J. P., González, A. L., & Montes, C. (2002). Combining expert knowledge and data mining in a medical diagnosis domain. *Expert Systems with Applications*, 23(4), 367–375.
- Brill, E. (1994). A report of recent progress in transformation-based error-driven learning. *Proceedings of the (ARPA-94) human language technology workshop*. NJ, USA: Princeton.
- Chen, P. W., & Lee, H. C. (2002). C<sup>2</sup>AS: An agent-based distributed and parallel processing virtual machine. *Seventh international conference on applications of high-performance computers in engineering (HPC 2002), September 23–25, Bologna, Italy*.
- Chou, S. M., Lee, T. S., Shao, Y. E., & Chen, I. F. (2004). Mining the breast cancer pattern using artificial neural networks and multivariate adaptive regression splines. *Expert Systems with Applications*, 27(1), 133–142.
- Eisenberg, D., Marcotte, E. M., Xenarios, I., & Yeates, T. O. (2000). Protein function in the post-genomic era. *Nature*, 405(6788), 823–826.
- Fukuda, K., Tsunoda, T., Tamura, A., & Takagi, T. (1998). Toward information extraction: Identifying protein names from biological papers. *Proceedings of the pacific symposium on biocomputing (PSB'98)* (pp. 707–718).
- Hong, T., & Han, I. (2002). Knowledge-based data mining of news information on the internet using cognitive maps and neural networks. *Expert Systems with Applications*, 23(1), 1–8.
- Lee, H. C., Dagli, C. H., Ercal, F., & Ozbayoglu, A. M. (1995). SimNet: A parallel neuro-fuzzy paradigm for data clustering. *OAI neural networks symposium and workshop (OAINN '95), Athens, Ohio, USA*.
- Marcotte, E. M., Xenarios, I., & Eisenberg, D. (2001). Mining literature for protein–protein interactions. *Bioinformatics*, 17(4), 359–363.
- Ono, T., Hishigaki, H., Tanigami, A., & Takagi, T. (2001). Automated extraction of information on protein–protein interactions from the biological literature. *Bioinformatics*, 17(4), 155–161.
- Porter, M. F. (1980). An algorithm for suffix stripping. *Program*, 14(3), 130–137.
- Thomas, J., Milward, D., Ouzounis, C., Pulman, S., & Carroll, M. (2000). Automatic extraction of protein interactions from scientific abstracts. *Pacific Symposium on Biocomputing*, 541–552.
- Valencia-García, R., Ruiz-Sánchez, J. M., Vicente, P. J. V., Fernández-Breis, J. T., & Martínez-Béjar, R. (2004). An incremental approach for discovering medical knowledge from texts. *Expert Systems with Applications*, 26(3), 291–299.
- Wang, H. C., Kuo, H. C., Chen, H. H., Hsiao, Y. Y., & Tsai, W. C. (2005). KSPF: Using gene sequence patterns and data mining for biological knowledge management. *Expert Systems with Applications*, 28(3), 537–545.
- Xenarios, I., Rice, D. W., Salwinski, L., Baron, M. K., Marcotte, E. M., & Eisenberg, D. (2000). DIP: The database of interacting proteins. *Nucleic Acids Research*, 28(1), 289–291.