

## Kabul Edilmiş Araştırma Makalesi (Düzenlenmemiş Sürüm)

## Accepted Research Article (Uncorrected Version)

### Makale Başlığı / Title

Taxonomic diversity-based domain interaction prediction

Taksonomik çeşitlilik tabanlı protein altünite etkileşim tahmini

### Yazarlar / Authors

Erdem TÜRK<sup>1\*</sup>, Barış Ethem SÜZEK<sup>2</sup>

### Referans No / Reference No

PAJES-18828

### DOI

10.5505/pajes.2018.18828

Bu PDF dosyası yukarıda bilgileri verilen kabul edilmiş araştırma makalesini içermektedir. Sayfa düzeni, dizgileme ve son inceleme işlemleri henüz tamamlanmamış olduğundan, bu düzenlenmemiş sürüm bazı üretim ve dizgi hataları içerebilir.

This PDF file contains the accepted research article whose information given above. Since copyediting, typesetting and final review processes are not completed yet, this uncorrected version may include some production and typesetting errors.



(Ş ' « a « © ¥ œ · Ÿ Ÿ z̄ j̄ ī Ÿ · Ÿ « μ © š ¥ a · ¥ a ° j̄ ® š α  
(š š - « a « © ¥ š · Ç j̄ - ¥ ® « ° ¥ · ¥ ¥ š · š ° š ¥ > ± š æ ¥ a © j̄ ° š

Erdem T & Z̄ · · š ® á , : 112 Ÿ □ j̄ © · ·

<sup>1,2</sup>Department of Computer Engineering, Mugla Sıtkı Kocman University, Mugla, Turkey  
erdemturk@mu.edu.tr, barissuzek@mu.edu.tr

Received 19.04.2018; Accepted 05.11.2018  
\*Corresponding author

doi: 10.5505/pajes.2018.18828  
Research Article

J. 1

Ş @ « · j̄ ¥ a · ş ş ° · a j̄ ¥ š ° ¥ j̄ ī ¥ © · ī ¥ a ¥ a ·  
~ @ « ° j̄ ¥ a · ī @ ¥ a · C « a š · ¥ μ « a j̄ · 2 j̄ μ š  
ş Y á @ Y á @ ! ¥ @ @ ( @ j̄ ī ī evrsm̄ prensipine  
Ÿ š μ š a s a · , ¥ @ · ī · s □ © ¥ a · μ á ° a j̄ ©  
ī ° © j̄ š · Ÿ Ç @ a · s · á @ á · s a · ¥ š ¥ · ~ @ « ° j̄  
Ç j̄ ¥ · ¥ · Ÿ Ö j̄ · 2 j̄ ī 2 @ ¥ ± · Ç š · Ç Ç š Ÿ š  
? j̄ ¥ · ¥ · Ÿ Ö j̄ · s μ š · á · š © « ° j̄ ¥ a ·  
ş Y · s a Y á @ á · s a · ! ¥ @ @ « ( @ j̄ ī ī · s a · s a · μ á °  
ó a j̄ @ © j̄ s · j̄ μ ¥ Ÿ · ( · , · Ł \$ · ¥ š ¥ · ~ @ « ° j̄  
ó a œ j̄ · , · ± · š @ Ÿ š · Ÿ š □ š · μ š · j̄ s · Ÿ · Ÿ ·  
μ j̄ @ ¥ a j̄ · , · ¥ · ī · ī · Ÿ · ī Ÿ · Ÿ · j̄ · Ÿ · ¥ a  
j̄ š · j̄ μ j̄ @ j̄ š · ~ @ « ° j̄ ¥ a · □ q @ ¥ « E j̄ š ·  
ş @ · © ş · á a (á · , · š Ç · š Ö œ j̄ · ¥ š · j̄ · Ÿ j̄ a j̄ μ ·  
~ @ « ° j̄ ¥ a · Ç j̄ · Ÿ · Ÿ · j̄ · Ÿ · Ÿ · s · ° · a · ¥ · j̄  
« · ° š μ š a · ° š · E j̄ · j̄ · « · ± · ± @ ± · © ±  
¥ Ç j̄ @ j̄ a · , · ¥ @ á š š @ š š Ÿ j̄ Ö š @ · j̄ a · Ÿ ¥ @ ¥ · )  
Ÿ ± μ š @ · á · á š · 2 j̄ · ī · ó Ÿ E · · · š ·  
~ j̄ @ Ç « @ @ s a · á a á a · ! ī · 2 j̄ · & · // · š ·  
ş ş @ · á · s · á @ á · s a · s a · ¥ Ç ¥ a · ī ·  
¥ Ç j̄ @ j̄ a · , · s μ š · s a μ š · ¥ · š ± · s a · á · Ÿ á  
Ÿ ± μ š @ · á · á š · 2 j̄ · ī · ó Ÿ E · · · š · E ó · j̄ @ Ÿ  
« · ~ · š · Ÿ š · ó Ÿ E · · · Ö · k̄ ± Ç ¥ a · « · š ·  
E ó · j̄ @ @ ¥ · « · Ÿ ± Ö ± · ~ j̄ @ Ç « @ @ s a · ~ · š ·  
alters · ¥ Ç · ¥ @ · μ á œ š · ( · , · Ł \$ · ¥ a · )  
Ÿ š □ š · á · « a · s a · s @ š · s @ š · s μ š a · Ÿ « Ö ±  
bir araya getiren bir metá ī · · š □ © ¥ a · μ á ° a j̄ ©  
ş « a ± ± a Ÿ š · E · Ç · , · ¥ @ · ş « a ± © š · E j̄ ·

Anahtar kelimeler: Ş @ « ° j̄ ¥ a · ş ş ° · a j̄ ¥ š ° ¥ j̄ ī · Ÿ · Ÿ ·  
birlikte evrimi, Protein fonksiyon analizi

Abstract

Identification of protein domain interactions (DDIs) is  
MirrorTree is a DDI prediction method that is based on the  
evolution. However, this method is sensitive to  
taxonomic diversity and evolutionary span within the two  
homolog sets compared to prediction methods work, we propose a  
MirrorTree based DDI prediction method, namely Taxonomic Div  
based Domain Interaction Prediction (TAXDIP). TAXDIP improves the  
MirrorTree method by adding a sampling step that  
representation of higher taxonomic ranks (e.g. family over sp  
in two protein homolog sets prior to their comparison. This  
step ensures increased evolutionary span within protein homologs.  
TAXDIP is first assessed using a set containing 6,514  
(interacting) domain pairs and a negative (non-interacting) set of equal  
size containing randomly generated domain pairs with no  
interactions. TAXDIP achieved 71.0% sensitivity and 63.0% spec  
on this set. Next, a benchmark containing 500 interacting and  
non-interacting domain pairs is used to compare the performance of  
TAXDIP against DDI prediction methods ME and RDF. TAXDIP  
outperforms ME and RDF. Its specificity remained better than ME. In  
conclusion, TAXDIP, with its performance, is a viable alternative to  
existing DDI prediction methods. In addition, the  
predictions are overlapping with, and furthermore, complementing other  
prediction methods, TAXDIP has a strong position in becoming  
a meta-DDI prediction method that combines multiple methods to  
consensus prediction.

Keywords: Protein domain interactions, Protein - evolution, Protein functional analysis.

### 1 Introduction

Proteins physically interact through domains located on them and these interactions are critical for their molecular functions and involvement in biological processes. Hence, identification of domain-domain interactions (DDIs) is a step towards identification of protein-protein interactions (PPIs) that helps in understanding biological processes. Several in vivo (e.g., yeast hybrid) and in vitro (e.g., x-ray crystallography, protein microarrays) PPI identification techniques exist, they require costly and time-consuming experiments. Compared to in vitro and in vivo techniques, computational (in silico) methods are fast and inexpensive alternatives. Computational DDI and PPI prediction methods with various algorithmic approaches are available. Spinzak et al. developed the Association Method (AM) to predict DDIs using the frequency of observed protein interactions. Alternatively, Deng and colleagues proposed a method based on an expectation maximization (EM) algorithm.

employs a maximum likelihood estimation (MLE) to find DDIs from protein interactions. Jothi et al. proposed a machine learning based method, namely, Relative Evolution of Domain Pairs (RCEDP) that uses sequence evolution to predict the domain pairs which are most likely to interact for a given protein. Gonzalez and Liao [4] employed a support vector machine (SVM) based computational method, interaction profile hidden Markov models (ipHMM) in order to predict DDIs. Furthermore, machine learning based algorithms (RDF [5], phylogenetic profiling [6] and probabilistic network models [7, 8]) have been proposed. In addition to these individual prediction methods in the literature, there are meta-DDI/PPI prediction methods. Rastan et al. proposed an integrated DDI and protein interaction analysis system (IPDI) which combines several prediction methods in order to build a consensus prediction score. The basis of the MirrorTree method is evolution; interacting proteins evolve simultaneously and likely to have similar phylogenetic trees. Briefly, the steps for MirrorTree based methods are 1) finding orthologs of two domains (or proteins)

that are candidates for interaction, 2) computing multiple sequence alignments for orthologs from organisms common to both domains, 3) computing two similarity matrices (one for each domain) using pairwise distances between orthologs, and finally 4) compute correlation coefficient between these similarity matrices to assess evolution and predict DDI (or PPI). The MirrorTree based methods are used extensively to predict DDI and PPI [3, 10-16]. Craig and Lia [10] used support vector machines calculated from phylogenetic trees in order to predict PPIs. Gertz et al. [11] used the Metropolis Monte Carlo optimization algorithm to detect possible matches between two distance matrices. Goh and colleagues used evolution for quantifying the behavior of the evolutionary histories of protein ligands and their receptors [12]. Jothi et al. [14] introduced a Monte Carlo search based method to detect interacting protein pairs. Pazos and Vaz [15, 16] employed a MirrorTree based approach for finding physical interactions between protein pairs and identifying the most probable sequence regions for these interactions. Although the MirrorTree is a competent approach for DDI prediction, they are sensitive to taxonomic diversity and evolutionary span of constructed matrices [17].

In this paper, we propose an improved MirrorTree based DDI prediction method, namely Taxonomic Diversity based Domain Interaction Prediction (TAXDIP), that uses representative homologs for taxonomic ranks higher than species (e.g., genus, family etc.) and consequently ensuring higher taxonomic span within similarity matrices prior to computation of correlation coefficient. TAXDIP is based on MirrorTree method but provides a solution to taxonomic diversity and evolutionary span problem of this method by adding a new taxonomy based sampling before computing similarity matrices. This extra step favor representation of higher taxonomic ranks (e.g. family) over lower ranks, therefore, provides increased evolutionary span.

The rest of the paper is organized as follows. Section 2 describes the datasets used. Section 3 describes our method, namely Taxonomic Diversity based Domain Interaction Prediction (TAXDIP), and steps towards its development including data preparation and experiments conducted to identify its parameters. Section 4 reports on the performance of our method and how it compares against available computational DDI prediction methods. Finally, Section 5 presents our conclusions.

## 2 Materials

### 2.1 Taxonomy dataset

The NCBI Taxonomy database [18] is one of the central sources that provides the species names and taxonomic lineage data for all known organisms. The NCBI Taxonomy database is downloaded from NCBI FTP site (<ftp://ftp.ncbi.nlm.nih.gov/pub/taxonomy/Species>, genus, family, order and class level taxa for about 1.49 million possible organisms are computed by traversing their lineages. These organisms represent 1.16 million species, 78.45 thousand genera, 8.30 thousand families, 381 orders and 304 classes.

### 2.2 Protein domain and multiple sequence alignment dataset

Pfam [19] is a protein family database which contains family annotations and multiple sequence alignments generated using

Hidden Markov models. The Pfam release 30.0 is used and contains a total of 16,306 families that covers 17.7 million protein sequences. The multiple sequence alignments (MSA) for protein domains are downloaded using Pfam website (<http://pfam.xfam.org>). The average number of members in a Pfam family for the release used is 2,667.

### 2.3 Protein domain-domain interaction dataset

DOMINE [20] contains a total of 6,634 DDIs are either extracted from Protein Data Bank (PDB) entries or predicted by at least one of the several computational methods including ME [21], RCDP [3], P-value [8], Interdom [22], DPEA [23], PE [24], GPE [25], DIPD [26], RDFF [5], K-GIDDI [27], Insite [28], and DomainGA [29] and DIMA [6]. The latest version of DOMINE has 26,219 unique DDIs involving 5,410 unique active Pfam domains. The same DDI may be predicted by multiple methods. Table 1 provides the number of DDIs predicted by individual computational methods in the DOMINE. A subset of DOMINE that contains 6,634 DDIs extracted from high resolution 3D structures is Pfam [30] and 3did [31], is referred to as DOMINE gold set. The main reason for our use of DOMINE gold set is testing TAXDIP algorithm with experimental data by avoiding potentially false DDI predictions made by other in silico methods.

Table 1: The number of domain-domain interactions in DOMINE predicted by individual computational methods.

Method	Number of DDIs
ME	2391
RCDP	960
P-value	596
Interdom	2768
DPEA	1812
PE	2588
GPE	1563
DIPD	2157
RDFF	2475
K-GIDDI	386
Insite	2408
DomainGA	459
DIMA	8012

## 3 Method

### 3.1 Creation of positive and negative sets

To assess the performance of our DDI prediction method, positive (interacting) and negative (non-interacting) domain pair sets are created. The positive set is based on the DOMINE gold set (see section 2.3) that contains 6,634 domain pairs. A total of 6,514 DDIs positive domain pairs are obtained after removal of 120 domain pairs due to missing data; in 81 pairs at least one of the Pfam domains are no longer in Pfam database and in 39 pairs Pfam multiple sequence alignments are not

available. A negative set of equal size (6,514 DDIs) is generated using randomly selected Pfam domain pairs that have no known or predicted interactions reported in the DOMINE database.

### 3.2 Generation of similarity matrices

The similarity matrices for Pfam domain pairs in positive and negative sets are computed using the following steps:

- i. For each Pfam domain pair, multiple sequence alignments are retrieved from Pfam web site
- ii. For each taxonomic rank from species to class (the taxonomic ranks above the class level are not examined since the number of taxa common to both domains becomes insignificant, e.g., 5 or less.)
  - a. For the selected taxonomic rank, find the list of taxa common to the proteins found in step (i).
  - b. Sample the intersecting taxa list from Step (a) using different sample sizes (10 to 100) and select one representative protein for each taxon. Figure 1 illustrates the impact of selection at different taxonomic ranks for one domain on an MSA based phylogenetic tree.
- iii. Compute similarity matrices based on pairwise distances between representative proteins from Step (b) using Coffee suite [32].

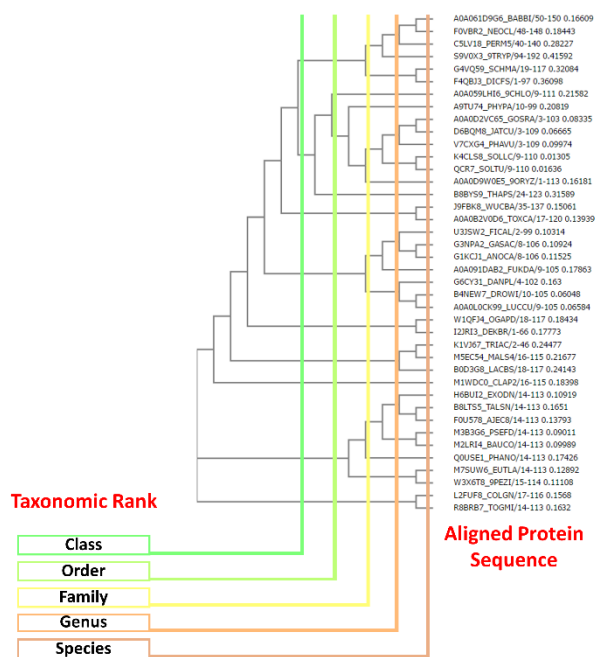


Figure 1: The impact of selection at different taxonomic ranks (species to class) for one domain illustrated on an MSA based phylogenetic tree.

### 3.3 Computation of correlation coefficients

The coevolution of domain pairs is evaluated by based on the Pearson correlation coefficient (Equation 1) is given as:

where  $n$  represents the size of similarity matrices and  $X_i$  are the distances between proteins computed from MSAs. The value of  $r_{xy}$  ranges between -1.0 and +1.0. According to the MirrorTree method, higher correlation coefficients indicate a higher level of coevolution and, thus, an interaction between the domains. Figure 2 illustrates the similarity matrices and the correlation computed at taxonomic rank of family for two Pfam domains PF00005 (ABC binding domain of ABC

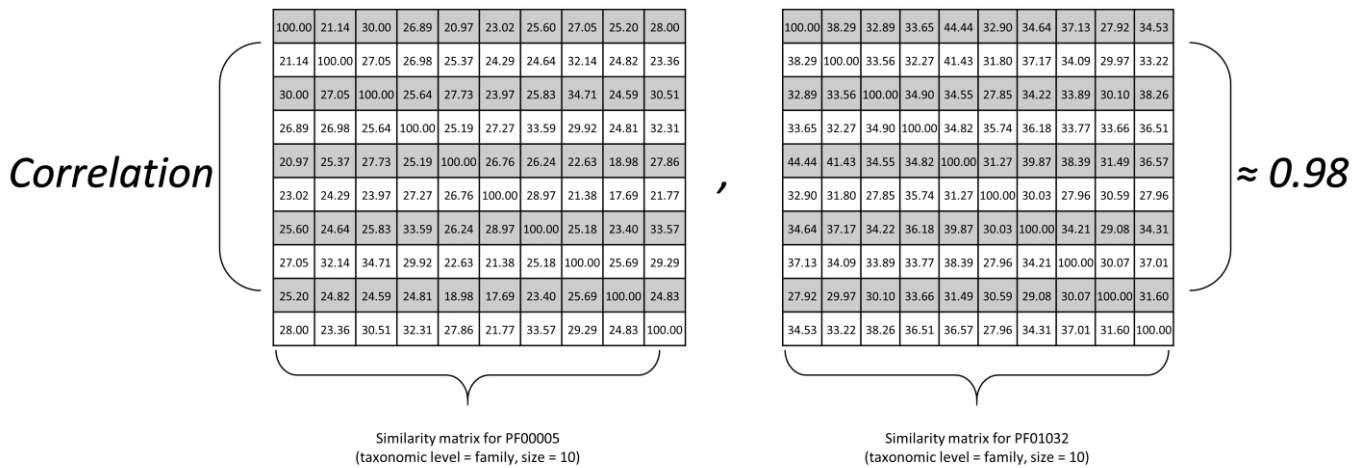


Figure 2: The similarity matrices computed at a family taxonomic rank level for two Pfam domains: PF00005 (ATPase ABC transporters) and PF01302 (Cytoskeleton-associated Protein Glycinerich domain), and their correlation. As per Mirror1 method a high correlation coefficient of 0.98 is considered as an indication for a domain-domain interaction between these two domains.

### 3.4 Proposed Taxonomic Diversity-based Domain Interaction Prediction (TAXDIP) algorithm

In the light of findings as reported in the Results section, TAXDIP algorithm is formulated as follows:

- i. Given a Pfam domain pair, the MSA of sequences used to generate domains are retrieved
- ii. For each taxonomic rank from family to species (higher to lower evolutionary span)
  - a. Find the taxa list common to the proteins containing the respective domains. If the intersection size is small (5 or less) proceed with the lower taxonomic rank.
  - b. Sample the common taxa list (40 is used as sample size) in Step (a) and select one representative protein for each taxon.
  - c. Generate two similarity matrices using proteins for both domains
  - d. Compute the Pearson correlation coefficient between the similarity matrices
- iii. Use the threshold of 0.5 on the Pearson correlation coefficient computed on the highest taxonomic rank in Step (ii) to decide on DDI. The thresholds based on the previous evolution based RCDIP method.

### 3.5 Performance evaluation

The performance of DDI predictions is computed based on a confusion matrix (Figure 3). There are four different outcomes for the DDI prediction. The number of correct DDI predictions on the positive and negative sets correspond to True and False Positives (TP and FP), respectively. Similarly, the number of incorrect DDI predictions correspond to True and False Negatives (TN and FN).

Table 2 shows the general structure of a confusion matrix and the common metrics that can be calculated from it. For

comparison with other DDI prediction methods, mutually used metrics true positive rate (or sensitivity) and true negative rate (or specificity) are used.

The Matthews correlation coefficient (MCC) is another performance metric, given by Equation 2, which measures the quality of binary classification. It is a measure of the quality of binary classification, where +1, 0 and -1 correspond to the perfect prediction, random prediction and disagreement between actual and predicted class respectively.

$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (2)$$

In addition, the area under the receiver operating characteristic (ROC) curve (AUC); a commonly used method to reduce ROC performance to a single scalar value and a classifier performance metric, is computed.

Table 2: Confusion matrix and common performance metrics computed based on this matrix.

Actual class \ Predicted class	Positive	Negative
	Positive	True Positives (TP)
Negative	False Negatives (FN)	True Negatives (TN)
Total	P	N

## 4 Results

### 4.1 Evaluation of taxonomic rank and sample size parameters on MirrorTree-based domain-domain interaction prediction

The effect of taxonomic rank and sample size (number of taxa sampled from common taxa list at selected taxonomic rank) on MirrorTree method is assessed using positive rate (TPR) and true negative rate (TNR) metrics. Different correlation coefficient thresholds from 0.4 to 0.9 are tested. The best performance is obtained using a correlation coefficient of 0.5 and this corroborates with earlier MirrorTree method RCDP [3].

Figure 3 illustrates the true positive and negative rates for different taxonomy ranks and sample sizes using the correlation coefficient of 0.5. The MirrorTree method achieves best TPR/TNR ratios for sample sizes in 50-100 bracket. The best AUC value of 0.759 is the taxonomy rank family and the sample size 40.

The motivation behind the proposed TAXDIP method (see section 3.4) is based on findings from above results. In summary:

- ◀ A sample size corresponding to 10% of overall common taxa list size (mostly 40-50 bracket) results in optimal DDI predictions
- ◀ The taxonomic diversity can change by domain. For instance, in some cases, the great majority of domain family members are from the same species. Furthermore, the diversity becomes next to none above the taxonomic rank of family. Hence, although DDI predictions are furthermore accurate when higher taxonomic ranks (e.g. family) are used, using such ranks may not be possible due to the limited evolutionary span of the domain family.

### 4.2 Comparison of TAXDIP with other methods

Table 3 provides the number of DDIs in DOMINE gold set predicted by individual computational methods [21] has the highest number of DDI predictions in DOMINE gold set as shown in Table 2 and better performance results in Pfam [30] and 3did [31] databases.

based on DOMINE gold set. A wide range of TPR values (DPEA:23.63%, PE:29.63%, DIPD:29.76%, RCDP:52.13%, DIMA:54.00%, ME:55.00%, RDFF:79.78%) and TNR (ME:55.55%, REF:64.38%, DIMA:99.84%) have been reported in their respective manuscripts. RDFF [5] has the best reported TPR and TNR values (79.78%/64.38%).

Considering the real life situation where the negatives far outnumber the positive interactions, a pseudo-life situation is simulated by subsampling positive and negative DDI sets with a 1:10 ratio (100 positive and 1000 negative DDIs) from the complete DDI set and by assessing the performance of TAXDIP on this subset. These simulations are repeated 1000 times. TAXDIP shows rather consistent performance regardless of any bias, with mean TPR of 71.00% and TNR of 300% for all the simulations. Hence, the

real life situation has an insignificant impact on the overall performance of TAXDIP.

Table 3: The number of domain-domain interactions in DOMINE gold set predicted by individual computational methods.

Method	Number of DDIs
ME	1326
RCDP	144
P-value	63
Interdom	399
DPEA	247
PE	328
GPE	362
DIPD	588
RDFF	148
K-GIDDI	68
Insite	147
DomainGA	459
DIMA	106

Using the benchmark set with a total of 1,000 domain pairs (500 interacting, 500 non-interacting) is compared with two existing DDI prediction methods considering two criteria. The first criterion is the reported performance results of these methods in their respective manuscripts [5, 21] and the second one is their coverage within or awareness of DOMINE gold set. Table 3) ME [21] is selected as it has the highest overlap with DOMINE gold set and RDFF is with the best reported TPR/TNR (79.78%/64.38%). The benchmark set contains DDIs common to sets used by TAXDIP, RDFF and ME. The DDI sets for RDFF and ME are obtained using their manuscripts. TPR, TNR and MCC values are reported in Table 4 provides the benchmark results of these methods.

Table 4: Benchmark results of TAXDIP, ME, and RDFF

Method	TPR	TNR	MCC
TAXDIP	75.40%	52.51%	0.2867
ME	66.20%	90.60%	0.5878
RDFF	21.00%	49.00%	-0.3116

The DDI prediction algorithms typically used in combination to establish a level of confidence on predictions. Thus, to assess how TAXDIP corroborates with other methods, the DOMINE gold set is used to generate predictions and other prediction methods contributed to the DOMINE database, namely ME [21], RCDP [3], Pvalue [8], Interdom [22], DPEA [23], PE [24], GPE [25], DIPD [26], RDFF [5], K-GIDDI [27], Insite [28], and DomainGA [29] and DIMA [6]. Table 5 provides the percentage of DDIs predicted by both

Uncorrected Version

Figure3: True positive and true negative rates for different taxonomic levels and sample sizes using 0.5 as correlation threshold.

Insite	51.70%
DomainGA	67.65%
DIMA	80.19%

TAXDIP and another method among all the DOMINE gold set DDI predictions made this other method a percentage of 100% indicates all the DDIs predicted by this other method is also predicted by TAXDIP while a percentage of 0% indicates the predictions made by TAXDIP complements these other methods predictions by 100%.

Table5: The percentage of DDIs predicted by both TAXDIP and other prediction methods among all DDIs predicted by the respective method within DOMINE gold set.

Method	Percentage
ME	74.66%
RCDP	74.30%
P-value	66.66%
Interdom	71.00%
DPEA	71.42%
PE	75.30%
GPE	67.96%
DIPD	74.00%
RDFD	73.98%
K-GIDDI	79.41%

## 5 Discussion

Identification of interactions between proteins is a critical step in understanding biological processes. Computational methods to predict DDIs and, consequently, PPIs is a cost-effective and rapid way to complement experimental studies, especially, to identify and prioritize interacting protein candidates for experimental validation.

In this work, we proposed a new algorithm called TAXDIP that is based on the MirrorTree method. TAXDIP fixes the known span as previously reported [67] by introducing an effective taxonomic rank-based sampling prior to generation of similarity matrices and computation of correlation coefficients.

Based on the reported performance results of existing DDI prediction methods, TAXDIP predicts DDIs with the best sensitivity/specificity (71.00%/63.00%) almost any other method. The only exception to this is RIFE [5] that has a pretty small coverage within DOMINE gold set (148 DDIs). TAXDIP predictions show overlap with other prediction methods.

According to the benchmark set, we used to compare TAXDIP against RDFD and ME, TAXDIP outperformed RDFD with better sensitivity, specificity and MCC score. The difference is significantly





