# A novel methodology on distributed representations of proteins using their interacting ligands

Hakime Öztürk<sup>1</sup>, Elif Ozkirimli<sup>2</sup> and Arzucan Özgür<sup>1</sup>

Departments of { 1 Computer Engineering, 2 Chemical Engineering }, Boğaziçi University, İstanbul, Turkey

### **1. ABSTRACT**

The effective representation of proteins is a crucial task. Related proteins usually bind to similar ligands. Chemical characteristics of ligands are known to capture the functional and mechanistic properties of proteins suggesting that a ligand based approach can be utilized in protein representation.

### **3. RESULTS**

**Table 1.** Distribution of families and super-families in A-50 dataset before and afterligand-interaction based filtering

Dataset	Num. Seq.	Super- families	Families
Refore	10816	1080	2100

- Total 10816 proteins, only (1639) 15% of them interacts at least a ligand.
- 64% of all interacting proteins have ligands fewer than 200 and 0.6% of all proteins are with single ligands.



We propose a novel method to compute similarity of proteins by describing them based on their ligands' SMILES.

### 2. METHODS

#### **COLLECTING INTERACTING LIGANDS OF PROTEINS**



After	1639	425	652

The mean number of the interacting ligands is 1791.







**n** words

For each chemical word that is extracted from ligand SMILES, a real-valued vector (embedding) is learned from a large training set [2]



SMILESVec (word) SMILESVec (char)

SMILESVec (word) SMILESVec (char)

The clustering is completed with TransClust [3] algorithm following similar pipeline to [4].

## 4. CONCLUSION

Using SMILESVec, we were able to define proteins based on their interacting ligands even in the absence of sequence or structure information.

✤ We showed that ligand-based protein representation, which uses only SMILES strings of the ligands that proteins bind to, performs as well as protein-sequence based representation methods in protein clustering.

Ligand-based protein description can be applied to different bioinformatics problems such as prediction of new protein-ligand interactions and protein function annotation.

**5. References** 







TUBITAK-BIDEB 2211-E and Bogaziçi University BAP are gratefully acknowledged.

