

Clustering FunFams using sequence embeddings improves EC purity

Maria Littmann^{1,2,*}, Nicola Bordin³, Michael Heinzinger^{1,2}, Christine Orengo^{3,*} & Burkhard Rost^{1,4}

1 TUM (Technical University of Munich) Department of Informatics, Bioinformatics & Computational Biology - i12, Boltzmannstr. 3, 85748 Garching/Munich, Germany

2 TUM Graduate School, Center of Doctoral Studies in Informatics and its Applications (CeDoSIA), Boltzmannstr. 11, 85748 Garching, Germany

3 Institute of Structural and Molecular Biology, University College London, London WC1E 6BT, UK.

4 Institute for Advanced Study (TUM-IAS), Lichtenbergstr. 2a, 85748 Garching/Munich, Germany & TUM School of Life Sciences Weihenstephan (WZW), Alte Akademie 8, Freising, Germany

* Corresponding author: Maria Littmann: littmann@rostlab.org, <http://www.rostlab.org/>, Tel: +49-289-17-814 (email rost: assistant@rostlab.org); Christine Orengo: c.orengo@ucl.ac.uk

Abstract

Motivation: Classifying proteins into functional families can improve our understanding of a protein's function and can allow transferring annotations within the same family. Toward this end, functional families need to be "pure", i.e., contain only proteins with identical function. Functional Families (FunFams) cluster proteins within CATH superfamilies into such groups of proteins sharing function, based on differentially conserved residues. 11% of all FunFams (22,830 of 203,639) also contain EC annotations and of those, 7% (1,526 of 22,830) have at least two different EC annotations, i.e., inconsistent functional annotations.

Results: We propose an approach to further cluster FunFams into smaller and functionally more consistent sub-families by encoding their sequences through embeddings. These embeddings originate from deep learned language models (LMs) transferring the knowledge gained from predicting missing amino acids in a sequence (ProtBERT) and have been further optimized to distinguish between proteins belonging to the same or a different CATH superfamily (PB-Tucker). Using distances between sequences in embedding space and DBSCAN to cluster FunFams, as well as identify outlier sequences, resulted in twice as many more pure clusters per FunFam than for a random clustering. 52% of the impure FunFams were split into pure clusters, four times more than for random. While functional consistency was mainly measured using EC annotations, we observed similar results for binding annotations. Thus, we expect an increased purity also for other definitions of function. Our results can help generating FunFams; the resulting clusters with improved functional consistency can be used to infer annotations more reliably. We expect this approach to succeed equally for any other grouping of proteins by their phenotypes.

Availability: The source code and PB-Tucker embeddings are available via GitHub: <https://github.com/Rostlab/FunFamsClustering>

Key words: Functional families, protein function, CATH, EC numbers, unsupervised learning, contrastive learning, word embeddings, transfer learning

Abbreviations used: DBSCAN, density-based spatial clustering of applications with noise; **d**, dimensions; **EC**, Enzyme Commission; **FunFam**, functional family; **LM**, language model; **NLP**, natural language processing

Introduction

Knowledge about the function of a protein is crucial for a wide array of biomedical applications and the classification of protein sequences into functional families can help transferring annotations from a functional family to an uncharacterized protein. Functional families can also reveal insights into the evolution of function through sequence changes [1]. To gain meaningful insights into protein functionality through functional families, it is important that those families are consistent, i.e., only contain functionally similar proteins.

CATH FunFams [2, 3] provide a functional sub-classification of CATH superfamilies [4, 5]. Superfamilies are the last level (H) in the CATH hierarchy; they group sequences which are related by evolution, often referred to as homologs. However, proteins in one superfamily can still be functionally and structurally diverse. Functional families (FunFams) provide a further sub-classification of superfamilies into coherent subsets of proteins with the same function. FunFams can be used to predict function on a per-protein level as described through Gene Ontology (GO) terms [6, 7], to predict functional sites [8], or to improve binding residue predictions by combining predictions of one FunFam in a consensus prediction [9].

The Enzyme Commission number (EC number) [10] numerically classifies enzymatic functions based on the reactions they catalyze. It consists of four levels and each level provides a more specific description of function than the previous one. The function of two proteins is more similar, the more levels of their two EC numbers are identical, particularly, for the levels EC3 and EC4 which describe the chemical reaction and its substrate specificity.

For 22,830 FunFams (11% of all), annotations for EC numbers for all four levels are available at least for one member. By design, proteins from the same FunFam should share the same EC class (annotated up to level 4). However, 1,526 FunFams (7% of 22,830) accounting for 16% of all sequences in the 22,830 FunFams with EC annotations have more than one annotation, and 180 (1% of 22,830) accounting for 2% of the sequences have even four or more different annotations (Fig. S1 in Supporting Online Material (SOM)). Different EC annotations within one FunFam could originate from moonlighting enzymes, i.e., enzymes with multiple functions [11]. Assuming the moonlighting enzyme to have two EC numbers, only one would be inconsistent with the other FunFam members rendering that FunFam inconsistent. However, different EC annotations can also result from impurity, i.e., FunFams containing proteins with different functions. Splitting FunFams further could provide a more fine-grained and consistent set of functionally related proteins.

Over the last few years, novel representations (embeddings) for proteins have emerged from adapting language models (LMs) developed for natural language processing (NLP) to protein sequences [12-16]. These embeddings are learnt solely from protein sequences without any additional annotations (self-supervised) using either auto-regressive pre-training (predicting the next amino acid, given all previous amino acids in a sequence, e.g., ELMo [17] or GPT [18]) or masked language modeling (reconstructing corrupted amino acids from the sequence, e.g., BERT [19]). To do well in those tasks, the LM is forced to learn frequently co-occurring sequence patches as well as more complex protein features such as those underlying secondary structure formation [20] as it is not possible to learn all possible amino acid permutations over the large set of protein sequences used for training. Features learnt implicitly by these models can be transferred to any task requiring protein representations by extracting the hidden states of the LM for a given protein sequence (transfer learning). It was shown previously that those learnt representations – referred to as embeddings – capture

1 higher-level features of proteins, including aspects of protein function beyond what is available
2 through traditional comparisons using sequence similarity or homology-based inference [21].
3 Therefore, we hypothesized that this orthogonal perspective – using embedding rather than
4 sequence space to transfer annotations – might help to find functionally consistent sub-groups
5 within protein families built using sequence similarity.

6 Here, we proposed a clustering approach to identify clusters in FunFams that are more
7 consistent in terms of shared functionality. To this end, shared functionality was defined as
8 sharing the same EC annotation up to the fourth level (i.e., completely identical EC numbers).
9 We represented protein sequences as embeddings, i.e., fixed-size vectors derived from pre-
10 trained LMs. We used the LM ProtBERT [15] to retrieve the initial embeddings, and applied
11 contrastive learning to map them onto a new embedding space where proteins within one
12 CATH superfamily were closer together than proteins from different superfamilies. The
13 resulting embeddings are called PB-Tucker. Clustering was then performed based on the
14 Euclidean distances between those embeddings using DBSCAN [22]. Within each FunFam
15 DBSCAN identified clusters as dense regions in which all sequences were close to each other
16 in embedding space; it classified proteins as outliers if they were not close to other sequences
17 in the FunFam. That allowed the identification of (i) a more fine-grained clustering of the
18 FunFams, and (ii) single sequences which might have been falsely assigned to this FunFam.
19 Analyzing whether or not embedding-based clustering reduced the number of different EC
20 annotations in a FunFam allowed validating our new approach.

21 22 23 **Methods**

24 **FunFams dataset.** The current version of CATH (v4.3) holds 4,328 superfamilies split into
25 212,872 FunFams. The FunFams generation process, albeit changing through time, consists
26 of various steps, starting with the clustering of all sequences within a CATH superfamily at
27 90% sequence similarity, encoding these clusters in Hidden Markov Models and creating a
28 relationship tree between all clusters using GeMMA [23] and HHsuite [24]. Subsequently,
29 CATH-FunFHMMer [7] is applied to transverse the tree and GroupSim [25] conservation
30 patterns are employed to merge or cut the tree branches to obtain the largest possible
31 alignment that is functionally pure. CATH FunFams have higher functional purity than CATH
32 superfamilies and conserved residues are enriched in functional sites [7].

33
34 **EC annotations and EC purity.** EC annotations for the FunFams dataset were obtained using
35 the UniProt [26] SPARQL API and cross-assigned to all UniProt IDs available within the
36 FunFams. Since proteins in the same FunFam are assumed to share a function, we expect all
37 proteins in one FunFam to have the same EC number(s). If not, the FunFam is considered
38 *impure*, i.e., it contains sequences which do not belong to this functional family. Impurity can
39 naively be defined as any FunFam with more than one EC number. However, some proteins
40 are annotated with multiple EC numbers. These proteins might actually execute multiple
41 functions (moonlighting) [11] or annotations might be wrong. Such an impurity is not caused
42 by an error in the creation of the FunFams and cannot be removed by further clustering the
43 FunFams. Instead, the naïve definition of impurity considers FunFams with one protein with
44 two different EC numbers as impure even if all other proteins share the same two EC numbers,
45 i.e., the annotations would be consistent and therefore, the FunFam should be considered as

1 pure. Consequently, we refined the definition considering FunFams as impure if one or more
2 relatives were annotated to additional EC numbers different to the other family members. We
3 only considered EC annotations with all four levels; all others were treated like those without
4 annotation.

5
6 **Protein representation.** We used ProtBERT [15] to create fixed-length vector
7 representations, i.e., vectors with the same number of dimensions irrespective of protein
8 length. ProtBERT uses the architecture of the LM BERT [19] which applies a stack of self-
9 attention [27] layers for masked language modeling (Supporting Online Material Section 1
10 (SOM_1) for details). Fixed-length vectors were derived by averaging over the representations
11 of each amino acid extracted from its last layer. This simple global average pooling provides
12 an effective baseline [12, 15, 20]. In the following, *ProtBERT* refers to this representation.

13 In order to capture the CATH hierarchy more explicitly, ProtBERT representations were
14 mapped to a new embedding space via contrastive learning. While supervised learning
15 requires phrasing a prediction task as e.g., a classification or regression task, contrastive
16 learning only requires some notion of similarity between samples. This similarity is used to
17 learn a new vector space that clusters similar samples while dissimilar items are separated.
18 Similarity can be defined between sample triplets potentially capturing their triangular relation;
19 in this case an *anchor* sample is given together with a *positive* and *negative* sample with the
20 positive being more similar to the anchor than the negative. The network then learns to push
21 anchor and positive toward each other while pushing anchor and negative apart. While
22 mapping a CATH-like hierarchy onto supervised classification is challenging, using a hierarchy
23 to define relative similarity between triplets is straightforward as anchor and positive only need
24 to share one level more in the hierarchy than anchor and negative. Toward this end, ProtBERT
25 representations were projected in two steps from 1024-dimensions (1024-d) to 128-d using
26 CATH v4.3 [3] for training the two-layer neural network (details in SOM_1). In the following, we
27 call these new 128-d embeddings *PB-Tucker* (Heinzinger et al., unpublished). PB-Tucker has
28 been trained to differentiate CATH superfamilies and seemed to better capture functional
29 relationships between proteins in one superfamily than the original ProtBERT (data not shown;
30 SOM_1 for more details).

31
32 **Clustering.** Representing sequences as PB-Tucker embeddings, we calculated the Euclidean
33 distance between all sequences within one FunFam. The distance d between two embeddings
34 x and y was defined as:

$$d(x, y) = \sqrt{\sum_{i=1}^{1024} (x_i - y_i)^2} \quad (\text{Eqn. 1})$$

37 Based on these distances, we clustered all sequences within one FunFam using the
38 implementation of DBSCAN [22] in scikit-learn [28] For a set of data points, DBSCAN identifies
39 dense regions, i.e., regions of points that are close to each other, and classifies these regions
40 as clusters. Data points not close to enough other data points are classified as outliers.
41 DBSCAN is based on the identification of *core points* that seed a cluster; all points within a
42 certain distance of the core point are added to this cluster. Two free parameters were
43 optimized: (1) The number of neighbors n (including the point itself) a point needs to have to

1 become “core point”; n implicitly controls the size and number of clusters, and (2) the distance
2 cutoff θ . Data points A and B are considered close, if $d(A,B) < \theta$. For our application, DBSCAN
3 has two major advantages: (1) The number of clusters does not have to be set *a priori*, and (2)
4 clustering and outlier detection are simultaneous.

5 If not stated otherwise, we used the default $n=5$ although it has been suggested to use
6 values between $n=D+1$ and $n=2*D-1$ where D is the number of dimensions [29] With $d=128$ for
7 the PB-Tucker embeddings that implies $n=255$. Since FunFams vary in size, n might be
8 adjusted to that size. For five superfamilies, we tested, in addition to $n=5$, $n=129$, $n=255$ as
9 fixed neighborhood sizes, as well as $n=0.05*|F|$, $n=0.1*|F|$, $n=0.2*|F|$ ($|F|$ =number of
10 sequences in FunFam) as variable neighborhood sizes dependent of the size of the FunFam.

11 Observing differences in the distances between the members of different superfamilies
12 (Fig. S2), it appeared best to choose superfamily-specific values for θ . Initially, we wanted to
13 determine a distance threshold reflecting the expected distance between any two members of
14 the same FunFam. However, large distances between members in one FunFam might reveal
15 impurity rather than a generic width of a family. Instead, we computed the median over those
16 distances for all FunFams in one superfamily and used this value for each FunFam. This way,
17 the value still reflects the expected distance between two members of a FunFam, but the effect
18 of large distances due to impurity should be averaged out by considering all FunFams in a
19 superfamily. In detail, for each member in each FunFam in a superfamily, we calculated its
20 average distance to all other members of that FunFam (distance distribution for five
21 superfamilies in Fig. S2). Given the distribution of these average sequence distances, we
22 chose the median distance as θ , i.e., we chose a distance cutoff so that 50% of all sequences
23 in a superfamily were on average within a distance of θ to all other sequences in the same
24 FunFam. Decreasing θ raises outliers and yields smaller clusters while increasing θ reduces
25 outliers and yields larger clusters.

26
27 **Measuring purity of clustered FunFams.** To estimate whether the clustering of an impure
28 FunFam led to more consistent sub-families, we calculated the percentage of pure clusters.
29 Clusters with no EC annotation were excluded. For each FunFam, we calculated the clusters
30 with one single EC annotation as percentage of all clusters with EC annotations and defined
31 this measure as the purity of a FunFam (Eqn. 2). We then defined the percentage of completely
32 pure FunFams as the percentage of FunFams with a purity of 100.

$$33 \quad Purity(F) = \frac{\#clusters_{pure}}{\#clusters_{with\ ECs}} \cdot 100 \quad (\text{Eqn. 2})$$

34
35 We also calculated the purity of a FunFam in terms of its size, i.e., the number of sequences
36 contained in it:

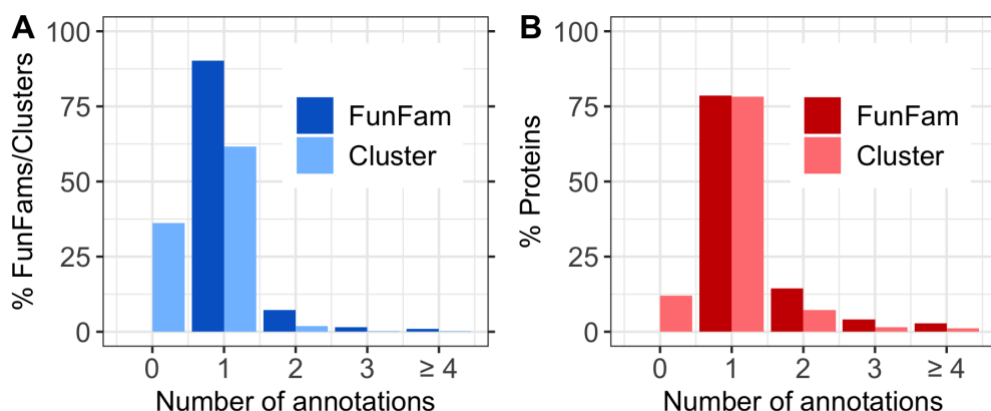
$$37 \quad Purity_{sequences}(F) = \frac{\#sequences\ in\ pure\ clusters}{\#sequences\ in\ clusters\ with\ ECs} \cdot 100 \quad (\text{Eqn. 3})$$

38
39 **Confidence intervals (CIs).** 95% symmetric confidence intervals (CIs) were calculated from
40 1, 000 bootstrap samples with replacement to indicate the spread of data and certainty of
41 average values.
42
43

1
2 **Final dataset.** To construct the dataset used in this analysis, we extracted all superfamilies
3 with at least one impure FunFam, i.e., at least one FunFam with more than one EC annotation.
4 Since embeddings could only be computed for continuous sequences, we excluded
5 sequences with multiple segments. After this removal, some FunFams became orphans
6 (single member) and were also excluded. This led to a final dataset of 458 superfamilies
7 (10.6% of all superfamilies) with 110,876 FunFams (52.1%) and 13,011 (6.1%) with EC
8 annotations. Those 13,011 FunFams accounted for 20% of all proteins in the FunFams
9 (1,669,245 sequences). All FunFams in a superfamily were used to determine a reasonable
10 distance cutoff for clustering while clustering was only performed for FunFams with EC
11 annotations. FunFams without EC annotations could have been clustered, too. However, since
12 EC annotations served as criterion for evaluation, only FunFams with such annotations were
13 clustered to save computational time and hence energy.

15 Results & Discussion

16 **Embedding clusters increased EC purity.** We began with 13,011 FunFams (6% of all) with
17 at least one EC annotation. Of these, 1,273 (10%) contained more than one EC annotation
18 (impure FunFams). Applying DBSCAN to all EC annotated FunFams, we split these into
19 26,464 clusters (21,546 for pure and 4,918 for impure FunFams). On average, 4.5% (95%
20 confidence interval (CI): [4.4%; 4.6%]) of the sequences in a FunFam were classified as
21 outliers (Table S1 in Supplementary Online Material (SOM)). 63% of the DBSCAN clusters
22 contained proteins with EC annotations; only 4% of those contained more than one EC
23 annotation (compared to 10% in all FunFams; Fig. 1A). Only 10% of all proteins (155,044 of
24 1,593,567) belonged to clusters with more than one EC annotation compared to 21% (356,565
25 of 1,668,273) for FunFams (Fig. 1B). Consequently, a larger fraction of clusters was pure (i.e.,
26 contained one EC annotation) than of FunFams both in terms of numbers of clusters and
27 numbers of proteins (Fig. 1).



28
29
30
31 **Fig. 1: EC purity for FunFams and embedding clusters.** This analysis considered 13,011 FunFams
32 with EC annotations. Panel A shows the distribution of all families (FunFam/Clusters), i.e., the
33 percentage of FunFams and embedding-based clusters with n EC annotations ($n \geq 1$ for FunFams and
34 $n \geq 0$ for new clusters; note: bars left and right of integer values n , not separated by a white space denote
35 n annotations). Panel B shows the distribution of all proteins, i.e., the percentage of proteins in families

(FunFam/Cluster) with n EC annotations. This number does not reveal how many proteins have an EC annotation. Of the 13,011 FunFams, 10% were impure, i.e., they contained more than one EC annotation (100-value for dark blue bar at 1 in panel A), and 21% of all proteins were part of these impure FunFams. After embedding-based split of FunFams, 64% (16,906) of the resulting clusters contain ECs (100-light blue bar at 0) and 4% (606) of those 16,906 were annotated to more than one EC accounting for 11% of proteins in clusters with ECs.

To further understand the extent to which the clustered FunFams provide a functionally more consistent subset, we determined for each impure FunFam, the fraction of clusters that were pure (Methods). To begin: 37% of all clusters had no EC annotated proteins and were excluded from further analysis. Of the remaining 16,906 clusters (63%), 22% were impure, i.e., contained more than one EC annotation. On average, 63% (CI: [60%; 66%]) of the clusters for a FunFam were pure (Fig. 2; dashed blue line) accounting for 58% (CI: [55%; 61%]) of all proteins (Fig. 2; dotted red line). 52% of all impure FunFams were split into completely pure clusters, i.e., for every other FunFam, the embedding-split clustered into functionally consistent sub-families (Fig. 2, right most blue point “100% Pure Clusters”) accounting for 38% of all proteins (Fig. 2, right most red point). This measure gave conservative estimates as it only considered completely pure clusters, ignoring improvements through reduction of EC annotations, e.g., when a group had originally $m+1$ annotations and the clustering improved to m , this improvement was ignored for all $m>1$.

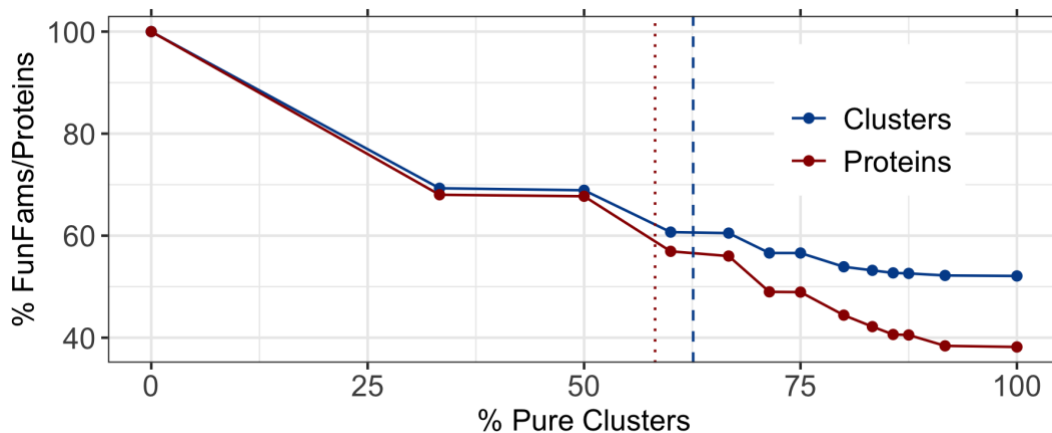
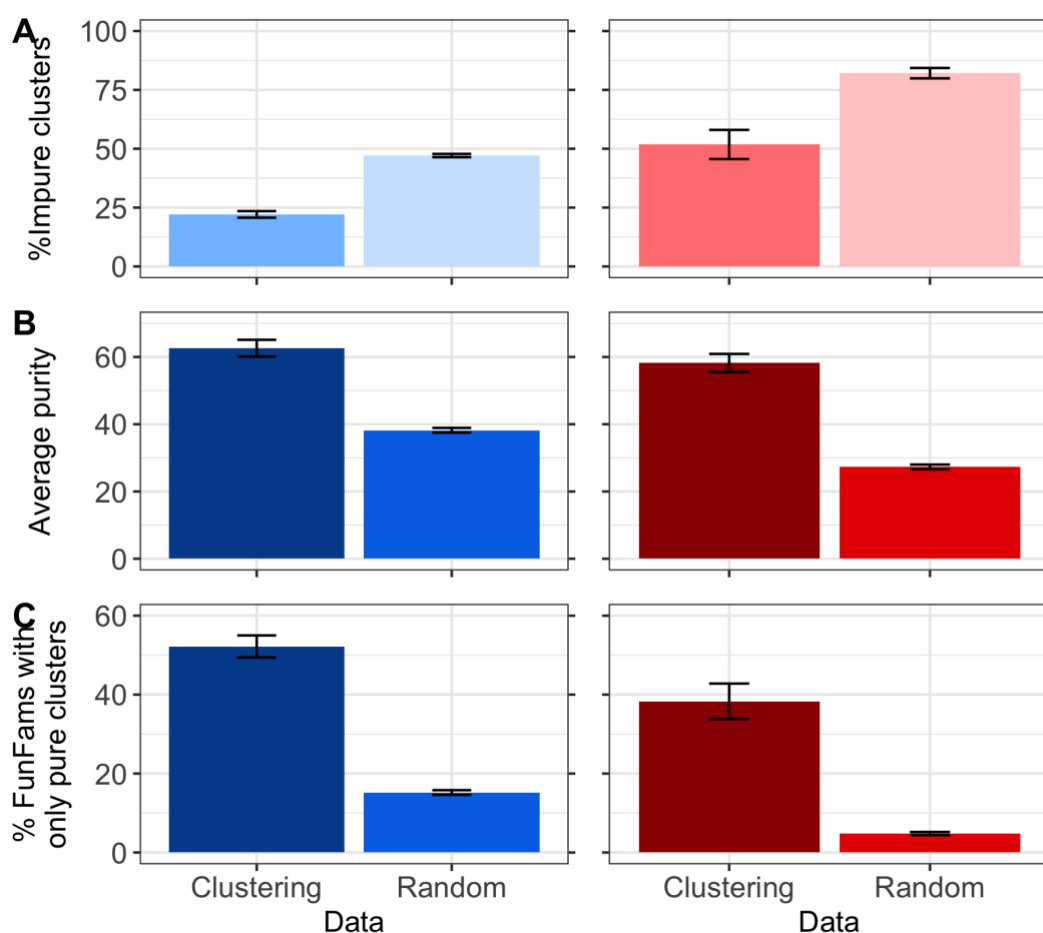


Fig. 2: Fraction of pure clusters for impure FunFams. We show the percentages of all FunFams (blue line) or of all proteins (red line) in FunFams at levels of increasing cluster purity (Eqns. 2, 3). On average, 63% of clusters for a FunFam were pure (dashed blue line) accounting for 58% of the proteins (dotted red line). 52% of impure FunFams were split only into pure clusters (right most blue point) accounting for 38% of the proteins.

Improving EC purity without over-splitting. While splitting impure FunFams through embedding-based clustering based clearly improved the EC purity of these FunFams, we wanted to avoid over-splitting. Trivially, the more and smaller clusters a FunFam is split into, the more likely it is that these clusters are pure. In the non-sense extreme case of having N clusters for N sequences (each sequence a cluster), all clusters are trivially pure. One constraint to avoid generating too many clusters (over-splitting) is to do substantially better than by randomly splitting into the same number of clusters. We computed the random

1 clustering using the same cluster sizes and outlier numbers as realized by the embedding-
2 based clustering. Embedding-based clustering outperformed random (Fig. 3): More than twice
3 as many clusters were impure for random than for embedding-based clustering (Fig. 3A,
4 $47\pm 1\%$ vs $22\pm 1\%$); the average purity of a FunFam was almost two times higher for
5 embedding-based clustering than for random (Fig. 3B, $63\pm 3\%$ vs $38\pm 5\%$), and 3.5 times more
6 FunFams were split into exclusively pure clusters by the embedding-based clustering (Fig. 3C,
7 $52\pm 3\%$ vs $15\pm 1\%$). This corresponded to 4.8% (CI: [4.4%; 5.2%]) of all proteins clustered into
8 pure clusters at random compared to 38% (CI: [33%; 43%]) of all proteins for embedding-
9 based clustering, i.e., an over 7-fold increase (Fig. 3C, red bars).

10



11

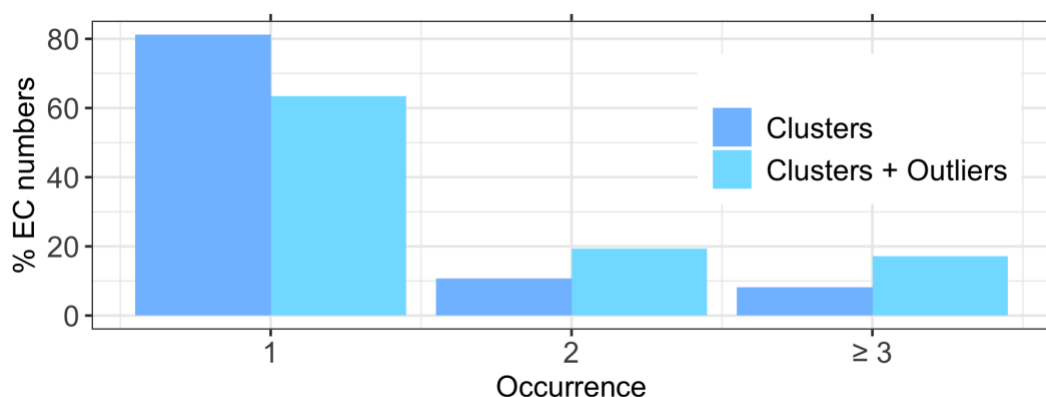
12

13 **Fig. 3: Embedding-based clusters improve EC purity over random.** Random clusters were
14 computed using the same cluster size and outlier number realized by the embedding-based clustering,
15 but the FunFam members were randomly assigned to one of these clusters or were classified as outliers.
16 Plots on the left (blue colors) show percentages of FunFams/Clusters, plots on the right show
17 percentages of proteins. **A.** The fraction of impure clusters was higher for the random clustering than
18 for our clustering (29% vs 12%). **B.** Through DBSCAN embedding-based clustering, each impure
19 FunFam was, on average, split into 63% pure clusters while for the random clustering, the average
20 purity was only 38%. **C.** More than half of all FunFams (53%) were split only into pure clusters for
21 embedding-based clustering but only 15% for a random clustering. Error bars indicate symmetric 95%
22 confidence intervals.

23

1 An ideal split of impure FunFams generates clusters defined by single EC numbers, i.e., all
2 cluster members share the same EC annotation and all proteins with the same EC annotation
3 end up in the same cluster. Ignoring the latter leads to over-splitting. For the embedding-based
4 clustering, 81% of the ECs occurred in one cluster (Fig. 4). However, some of the outliers had
5 EC annotations. When also counting those (as single member clusters), the percentage of EC-
6 exclusive clusters dropped to 63% (Fig. 4). These results suggested the embedding-based
7 clustering to have largely avoided over-splitting. Nevertheless, 8% of all experimentally known
8 EC numbers were annotated to proteins from at least three different clusters (17% if including
9 outliers; Fig. 4) and some (10%) of the outliers shared the EC number with the cluster from
10 which they had been removed. This might indicate over-splitting or suggest a more fine-grained
11 functional distinction between those proteins than is captured in the fourth EC level.

12



13

14

15 **Fig. 4: Most EC numbers only occur in one cluster.** For each EC number in a FunFam, we counted
16 the number of embedding-based clusters in which it occurred to gauge potential over-splitting. 81% of
17 the ECs only occurred in one cluster (darker bars). If we considered outliers as clusters with one
18 member, this number dropped to 63%. These results suggest that the clustering did not over-split the
19 FunFams and that functionally related proteins ended up in the same cluster.

20

21 If the increased purity through clustering had been a random effect, the embedding space
22 clustering would be EC-independent. If so, we expect no difference in the distributions of
23 embedding distances between pure and impure FunFams, and a similar number of clusters
24 and outliers. However, pure FunFams were, on average, split into only two clusters, while
25 impure FunFams were split into four clusters (Table S1). The number of clusters can, thus,
26 indicate whether a FunFam is impure or not, i.e., if a FunFam is split into many clusters, it
27 should be considered for further manual inspection to establish whether all proteins were
28 correctly assigned to this functional family (more details in SOM_2.1).

29

30 **Different levels of EC annotations gave similar results.** Up to this point, we only
31 distinguished whether two proteins were annotated to the same or to different EC numbers,
32 ignoring that two proteins with ECs A.B.C.X and A.B.C.Y more likely have similar molecular
33 function than a pair with A.* and D.*. Pairs of the first type (difference only in 4th-EC level) will,
34 on average, be more sequence similar than pairs of the second type (different EC numbers at
35 the top level). Most impure FunFams were impure due to differences on the fourth level of EC
36 annotations (Fig. S4). Although we analyzed the clustering at higher levels of the EC

1 classification, the results were inconclusive, probably due to data sparsity (SOM_2.2 for more
2 details).

3
4 **Details of parameter choice mattered.** For a more detailed analysis of particular details of
5 our method, in particular, for the choice of embeddings and clustering parameters, we chose
6 five superfamilies with diverse properties (CATH identifiers: 3.40.50.150, 3.20.20.70,
7 3.40.47.10, 3.50.50.60, 1.10.630.10; SOM_3).

8 *More consistent clustering from PB-Tucker.* When using ProtBERT embeddings to
9 cluster the five chosen superfamilies, the number of clusters and outliers was smaller, but the
10 fraction of impure clusters was higher than when using the default PB-Tucker embeddings
11 (19% for ProtBERT vs 13% for “default”; Table S4). The average purity was also higher for
12 PB-Tucker (“default” = 59%) than for ProtBERT (51%) (Table S4). Thus, PB-Tucker appeared
13 superior in capturing functional differences, yielding a more fine-grained and pure clustering.

14 *Smaller distance thresholds led to smaller and purer clusters.* The distance threshold
15 q of DBSCAN defines whether or not two points are close enough to each other to be grouped.
16 For the default clustering, we chose the median distance between all proteins for each
17 superfamily (Methods). The observed distribution of distances (Fig. S2) suggested choosing
18 superfamily-specific thresholds. As expected, the smaller q , the more clusters and outliers will
19 result (“ q =1st quartile” vs “default”; Table S4). Largely due to splitting FunFams into more
20 clusters at smaller q , the resulting clusters were seemingly purer with only 4% impure clusters
21 (vs. 13%) and an average purity of 83% (vs. 59%) (Table S4). In contrast, larger q thresholds
22 (here the 3rd quartile) affected fewer, more impure clusters (Table S4). Thus, the choice of q
23 highly influences the clustering results. For some applications, lower values of q might be best
24 to obtain a large, highly consistent set of small sub-families that can serve e.g., as seed to
25 further extend those sub-families to larger functionally related families. Also, especially for
26 FunFams for which using a larger cutoff did not result in any clusters and only a small number
27 of outliers, decreasing the distance threshold can help to still identify which sequences might
28 cause impurity.

29 *Default neighborhood size resulted in best clustering.* DBSCAN forms clusters around
30 “core points” which are points with at least n neighbors. For the five superfamilies, we tested
31 fixed neighborhood sizes of $n \in \{5; 129; 255\}$ and variable neighborhood sizes dependent on the
32 size of the FunFam $n = x * |F|$, with $|F|$ as the number of proteins in a FunFam and $x \in [0.01; 0.1; 0.2]$
33 (Methods). While $n=129$ and $n=255$ were in the range of what is recommended for n [29], the
34 clustering was worse than for the default parameter ($n=5$) (Table S4). Specifically, the number
35 of outliers exploded for these large neighborhood sizes (Table S4, Fig. S5).

36 Since FunFams differ substantially in the number of proteins, we hypothesized that –
37 similar as for q – it could be reasonable to choose a different n for each FunFam. However,
38 this did not improve compared to the default clustering (Table S4); the default $n=5$ was a good
39 choice.

40 *No consistent influence of level of EC annotation.* Assessing how well our clustering approach
41 worked depending on the level on which EC annotations were different (E.g., for a given level,
42 for example EC level 3, we checked that annotations for level 1 and 2 were consistent) did not
43 reveal a consistent trend (Fig. S4). The results for the five chosen superfamilies were similar
44 (Fig. S6) underlining the more general findings that the level of EC annotation causing impurity

1 did not tremendously influence the results of the clustering (see SOM for a more detailed
2 analysis). The performance is rather impacted by other factors like the presence of
3 moonlighting proteins or missing annotations. We applied a rather conservative definition of
4 purity: If one protein is annotated to two EC numbers and another protein in the same cluster
5 is only annotated to one of those two, we considered this cluster impure. We would argue that
6 it makes sense to group all proteins with two annotations in one cluster and proteins with only
7 one annotation in another. However, those proteins also clearly share some function and
8 considering those FunFams or clusters as impure is probably too strict. Also, we cannot be
9 sure whether proteins with only one annotation are correctly annotated or are missing an
10 annotation. In general, missing annotations limited our approach. Many resulting clusters did
11 not contain any sequence with an EC annotation making it hard to assess whether those
12 clusters were pure or not. Assessing how well our clustering approach worked depending on
13 the level on which EC numbers differed in impure FunFams (e.g., for EC level 3, annotations
14 for levels 1 and 2 were consistent) did not reveal a consistent trend (Fig. S4). The results for
15 the five chosen superfamilies were similar (Fig. S6) underlining the more general findings that
16 the level of EC annotation causing impurity did not crucially affect the embedding-based
17 clustering (SOM_3.2). Instead, the performance was likely impacted more by other factors
18 such as the presence of moonlighting proteins or missing annotations. We applied a rather
19 conservative definition of purity: If one protein is annotated to two EC numbers and another
20 protein in the same cluster is only annotated to one of those two, we considered this cluster
21 impure. We would argue that it makes sense to group all proteins with two annotations in one
22 cluster and proteins with only one annotation in another. However, those proteins also clearly
23 share some function and considering those FunFams or clusters as impure is probably too
24 strict. Also, we cannot be sure whether proteins with only one annotation are correctly
25 annotated or are missing an annotation. In general, missing annotations limited our approach.
26 Many resulting clusters did not contain any protein with an experimental EC annotation making
27 it hard to assess whether those clusters were pure or not.

28
29 **Clustering increased purity of ligand binding.** Another way to assess the purity of molecular
30 protein function within a group of proteins is by comparing the extent to which they are similar
31 in terms of ligand-binding. We extracted bound ligands from BioLip [30] and only considered
32 annotations defined as the cognate ligand [31] (SOM_1.2). Of the 13,011 FunFams considered
33 so far, 950 (7%) contained any annotation about a ligand bound, and of those 950, 158 (17%)
34 were annotated with more than one different ligand. Embedding-based clustering split 33% of
35 these FunFams into clusters with only one type of ligand, i.e., “pure” clusters (compared to
36 52% for EC level 4) and an average purity of 36% (compared to 63% for ECs). Although ligand
37 annotations remained limited, these results confirmed that embedding-based clustering
38 increased functional purity of FunFams for an aspect of function not used during method
39 development.

Conclusions

FunFams [4, 5] provide a high-quality sub-classification of CATH superfamilies into families of functionally related proteins [3, 32]. However, some FunFams are impure and 7% of all FunFams with EC annotations contain at least two different ECs (Fig. S1). Here, we introduced a novel approach toward clustering proteins through embeddings derived from the LM ProtBERT [15] and further optimized to capture relationships between proteins within one CATH superfamily (called *PB-Tucker*). Similarity between embeddings can capture information different from what is captured by sequence similarity. In particular, it can reveal new functional relations between proteins [21]. Clustering all FunFams with more than one EC annotation (impure FunFams) using DBSCAN [22] reduced the percentage of impure clusters to 22% (95% confidence interval (CI): [21%, 23%]). An impure FunFam was on average split into 63% pure clusters (CI: [60%: 66%]) and more than half (53%, CI: [50%; 56%]) of all impure FunFams were split into fully pure sub-families (Fig. 2). This corresponded to a four-fold increase over random clustering (Fig. 3B). In terms of number of proteins (rather than number of clusters), the increase was almost ten-fold. Only 4.8% (CI: [4.4%; 5.2%]) of the proteins were in FunFams split into pure clusters for random while this number rose to 38% (CI: [33%; 43%]) for the *PB-Tucker* embedding-based clustering.

A more detailed analysis of five hand-picked superfamilies (Table S2) showed that the default choices for the DBSCAN parameters were reasonable (Figs. S4 & S5, Table S4), with the default $n=5$ to define the number of neighbors for a point to be considered a core point and the distance threshold q determined automatically based on the median distance between proteins within one FunFam.

Restricting the analysis to experimental EC annotations limited the validation of our approach to a small fraction (6.1%) of all FunFams and even for those FunFams, most EC annotations remain unknown. Nevertheless, we have shown that our approach could capture more fine-grained functional relationships and enabled splitting FunFams into more functionally consistent sub-families. Especially for FunFams without many known functional annotations, our clustering can be used to (i) investigate whether or not the family could be impure based on the number of clusters resulting from the embedding-based split, or (ii) more safely infer functional annotations between members of one functional cluster than between members of one FunFam. We presented evidence suggesting that the findings for EC annotations will hold for other aspects of protein function, e.g., for binding. While we only applied this approach to FunFams using embeddings optimized for CATH, this clustering could be applied to any database of functional families using a more generalized version of those embeddings.

Acknowledgements

Thanks to Tim Karl and Inga Weise (both TUM) for invaluable help with technical and administrative aspects of this work. We would like to acknowledge Ian Sillitoe (UCL) for helpful comments on EC data. Last, but not least, thanks to all maintainers of public databases and to all experimentalists who enabled this analysis by making their data publicly available.

This work was supported by the Bavarian Ministry of Education through funding to the TUM, by a grant from the Alexander von Humboldt foundation through the German Ministry for Research and Education (BMBF: Bundesministerium für Bildung und Forschung), and by two grants from BMBF (031L0168 and Program “Software Campus 2.0 (TUM): 01IS17049) as well as by a grant from Deutsche Forschungsgemeinschaft (DFG–GZ: RO1320/4–1). We gratefully acknowledge the support of NVIDIA Corporation with the donation of one Titan GPU used for this research. Nicola Bordin acknowledges financial support from the Biotechnology and Biological Sciences Research Council (UK) [BB/R009597/1].

References

- [1] S. S. Hannenhalli and R. B. Russell, “Analysis and prediction of functional sub-types from protein sequence alignments,” *J Mol Biol*, vol. 303, no. 1, pp. 61–76, Oct. 2000, doi: 10.1006/jmbi.2000.4036.
- [2] C. A. Orengo, A. D. Michie, S. Jones, D. T. Jones, M. B. Swindells, and J. M. Thornton, “CATH—a hierarchic classification of protein domain structures,” *Structure*, vol. 5, no. 8, pp. 1093–108, Aug. 1997, doi: 10.1016/s0969-2126(97)00260-8.
- [3] “CATH: Protein Structure Classification Database at UCL.” <https://www.cathdb.info/> (accessed Nov. 02, 2020).
- [4] I. Sillitoe *et al.*, “CATH: increased structural coverage of functional space,” *Nucleic Acids Res.*, no. gkaa1079, Nov. 2020, doi: 10.1093/nar/gkaa1079.
- [5] S. Das, D. Lee, I. Sillitoe, N. L. Dawson, J. G. Lees, and C. A. Orengo, “Functional classification of CATH superfamilies: a domain-based approach for protein function annotation,” *Bioinformatics*, vol. 32, no. 18, p. 2889, Sep. 2016, doi: 10.1093/bioinformatics/btw473.
- [6] I. Sillitoe *et al.*, “New functional families (FunFams) in CATH to improve the mapping of conserved functional sites to 3D structures,” *Nucleic Acids Res*, vol. 41, no. Database issue, pp. D490-8, Jan. 2013, doi: 10.1093/nar/gks1211.
- [7] N. Zhou *et al.*, “The CAFA challenge reports improved protein function prediction and new functional annotations for hundreds of genes through experimental screens,” *Genome Biol*, vol. 20, no. 1, p. 244, Nov. 2019, doi: 10.1186/s13059-019-1835-8.
- [8] S. Das *et al.*, “CATH FunFHMMer web server: protein functional annotations using functional family assignments,” *Nucleic Acids Res*, vol. 43, no. W1, pp. W148-53, Jul. 2015, doi: 10.1093/nar/gkv488.
- [9] S. Das, H. M. Scholes, and C. A. Orengo, “CATH functional families predict protein functional sites,” 2020, doi: 10.1101/2020.03.23.003012.
- [10] L. Scheibenreif, M. Littmann, C. Orengo, and B. Rost, “FunFam protein families improve residue level molecular function prediction,” *BMC Bioinformatics*, vol. 20, no. 1, p. 400, Jul. 2019, doi: 10.1186/s12859-019-2988-x.

- 1 [11] E. C. Webb, *Enzyme nomenclature 1992: recommendations of the Nomenclature*
2 *Committee of the International Union of Biochemistry and Molecular Biology on the*
3 *nomenclature and classification of enzymes*. Academic Press, 1992.
- 4 [12] C. J. Jeffery, "Moonlighting proteins: old proteins learning new tricks," *Trends Genet*, vol.
5 19, no. 8, pp. 415–7, Aug. 2003, doi: 10.1016/S0168-9525(03)00167-7.
- 6 [13] M. Heinzinger *et al.*, "Modeling aspects of the language of life through transfer-learning
7 protein sequences," *BMC Bioinformatics*, vol. 20, no. 1, p. 723, Dec. 2019, doi:
8 10.1186/s12859-019-3220-8.
- 9 [14] R. Rao *et al.*, "Evaluating Protein Transfer Learning with TAPE," *ArXiv190608230 Cs Q-*
10 *Bio Stat*, Jun. 2019, Accessed: Nov. 03, 2020. [Online]. Available:
11 <http://arxiv.org/abs/1906.08230>.
- 12 [15] E. C. Alley, G. Khimulya, S. Biswas, M. AlQuraishi, and G. M. Church, "Unified rational
13 protein engineering with sequence-based deep representation learning," *Nat. Methods*,
14 vol. 16, no. 12, Art. no. 12, Dec. 2019, doi: 10.1038/s41592-019-0598-1.
- 15 [16] A. Elnaggar *et al.*, "ProtTrans: Towards Cracking the Language of Life's Code Through
16 Self-Supervised Deep Learning and High Performance Computing," *bioRxiv*, p.
17 2020.07.12.199554, Jul. 2020, doi: 10.1101/2020.07.12.199554.
- 18 [17] A. Madani *et al.*, "ProGen: Language Modeling for Protein Generation," *bioRxiv*, p.
19 2020.03.07.982272, Mar. 2020, doi: 10.1101/2020.03.07.982272.
- 20 [18] M. E. Peters *et al.*, "Deep contextualized word representations," *ArXiv180205365 Cs*,
21 Mar. 2018, Accessed: Nov. 02, 2020. [Online]. Available: <http://arxiv.org/abs/1802.05365>.
- 22 [19] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving Language
23 Understanding by Generative Pre-Training," p. 12.
- 24 [20] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep
25 Bidirectional Transformers for Language Understanding," *ArXiv181004805 Cs*, May
26 2019, Accessed: Nov. 02, 2020. [Online]. Available: <http://arxiv.org/abs/1810.04805>.
- 27 [21] A. Rives *et al.*, "Biological structure and function emerge from scaling unsupervised
28 learning to 250 million protein sequences," *bioRxiv*, p. 622803, Jan. 2020, doi:
29 10.1101/622803.
- 30 [22] M. Littmann, M. Heinzinger, C. Dallago, T. Olenyi, and B. Rost, "Embeddings from deep
31 learning transfer GO annotations beyond homology," *bioRxiv*, p. 2020.09.04.282814, Oct.
32 2020, doi: 10.1101/2020.09.04.282814.
- 33 [23] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering
34 clusters in large spatial databases with noise," 1996, vol. 96, pp. 226–231.
- 35 [24] D. A. Lee, R. Rentzsch, and C. Orengo, "GeMMA: functional subfamily classification
36 within superfamilies of predicted protein structural domains," *Nucleic Acids Res.*, vol. 38,
37 no. 3, pp. 720–737, Jan. 2010, doi: 10.1093/nar/gkp1049.
- 38 [25] M. Steinegger, M. Meier, M. Mirdita, H. Vöhringer, S. J. Haunsberger, and J. Söding, "HH-
39 suite3 for fast remote homology detection and deep protein annotation," *BMC*
40 *Bioinformatics*, vol. 20, no. 1, p. 473, Sep. 2019, doi: 10.1186/s12859-019-3019-7.
- 41 [26] J. A. Capra and M. Singh, "Characterization and prediction of residues determining
42 protein functional specificity," *Bioinformatics*, vol. 24, no. 13, pp. 1473–1480, Jul. 2008,
43 doi: 10.1093/bioinformatics/btn214.
- 44 [27] T. U. Consortium, "UniProt: a worldwide hub of protein knowledge," *Nucleic Acids Res.*,
45 vol. 47, no. D1, pp. D506–D515, Jan. 2019, doi: 10.1093/nar/gky1049.

- 1 [28] “UniProt.” <https://sparql.uniprot.org/> (accessed Nov. 12, 2020).
- 2 [29] D. Bahdanau, K. Cho, and Y. Bengio, “Neural Machine Translation by Jointly Learning to
3 Align and Translate,” *ArXiv14090473 Cs Stat*, May 2016, Accessed: Nov. 02, 2020.
4 [Online]. Available: <http://arxiv.org/abs/1409.0473>.
- 5 [30] D. Shen *et al.*, “Baseline Needs More Love: On Simple Word-Embedding-Based Models
6 and Associated Pooling Mechanisms,” *ArXiv180509843 Cs*, May 2018, Accessed: Nov.
7 02, 2020. [Online]. Available: <http://arxiv.org/abs/1805.09843>.
- 8 [31] F. Pedregosa *et al.*, “Scikit-learn: Machine Learning in Python,” *J. Mach. Learn. Res.*, vol.
9 12, pp. 2825–2830, 2011.
- 10 [32] J. Sander, M. Ester, H.-P. Kriegel, and X. Xu, “Density-Based Clustering in Spatial
11 Databases: The Algorithm GDBSCAN and Its Applications,” *Data Min. Knowl. Discov.*,
12 vol. 2, no. 2, pp. 169–194, Jun. 1998, doi: 10.1023/A:1009745219419.
- 13 [33] J. Yang, A. Roy, and Y. Zhang, “BioLiP: a semi-manually curated database for biologically
14 relevant ligand–protein interactions,” *Nucleic Acids Res.*, vol. 41, no. Database issue, pp.
15 D1096–D1103, Jan. 2013, doi: 10.1093/nar/gks966.
- 16 [34] J. D. Tyzack, L. Fernando, A. J. M. Ribeiro, N. Borkakoti, and J. M. Thornton, “Ranking
17 Enzyme Structures in the PDB by Bound Ligand Similarity to Biological Substrates,”
18 *Structure*, vol. 26, no. 4, pp. 565-571.e3, Apr. 2018, doi: 10.1016/j.str.2018.02.009.
- 19