

A robust method to detect structural and functional remote homologues

Ori Shachar¹ and Michal Linial^{2,3}

¹School of Computer Science and Engineering, ²Department of Biological Chemistry, Institute of Life Sciences, The Hebrew University, Jerusalem 91904, Israel

³Corresponding author

Michal Linial

Fax: +972-2-6586448

Phone: +972-2-6585425

E. mail: michall@cc.huji.ac.il

Abstract: 181 words

Tables: 3

Figures: 1

Pages: 20

Supplementary: <http://www.protonet.cs.huji.ac.il/homologues/>

Total words: ~5640

Running title

Remote homologues in protein families

Abstract

With currently available sequence data, it is feasible to conduct extensive comparisons among large sets of protein sequences. It is still a much more challenging task to partition the protein space into structurally and functionally related families solely based on sequence comparisons. The ProtoNet system automatically generates a tree-like classification of the whole protein space. It stands to reason that this classification reflects evolutionary relationships, both close and remote. In this paper we examine this hypothesis. We present a semi-automatic procedure that singles out certain inner nodes in the ProtoNet tree that should ideally correspond to structurally and functionally defined protein families. We compare the performance of this method against several expert systems. Some of the competing methods incorporate additional extraneous information on protein structure or on enzymatic activities. The ProtoNet-based method performs at least as well as any of the methods with which it was compared. This paper illustrates the ProtoNet-based method on several evolutionarily diverse families. Using this new method, an evolutionary divergence scheme can be proposed for a large number of structural and functional related superfamilies.

Keywords: Protein Family, Hierarchical classification, protein space, database

Introduction

The gap between the amount of sequences available and the number of homologous families known is still growing (Park et al., 1998; Russell et al., 1997). This paper addresses the problem of finding remote homologues in the sequence space and suggesting for any given protein sequence its relationship with other proteins, which basically form a structural or a functional homologous family.

To fully appreciate the complexity of the problem one needs to imagine the whole sequence space of biologically functional sequences, where the pairwise distance between every pair of sequences is exactly the pairwise score given by some chosen method for sequence similarity (i.e. FASTA or BLAST, with any chosen substitution matrix, Stark et al., 2003). Our task of defining functional families would sum up to drawing an outline around all proteins that belong to the same functional family and thus separating them from all sequences of other families. For single domain proteins that have evolved by divergent evolution, such a description is valid, as in the case of the globin and histone families. However, for multi-domain proteins that may have resulted from the fusion of domains and for those that have evolved by convergent evolution (Saier, 1996), such a simplified partition of the protein sequence space is incompatible. Such groups of proteins that are related by convergence in evolution or that are linked through shared domains are abundant in the protein space. They are especially abundant in complex proteomes. In reality, methods for classifications of such cases into functional or structural families are not reliable.

Most proteins, especially those derived from newly sequenced genomes, are not yet annotated and their evolutionary relations to other proteins are not evident. As a result, the task of expanding and generalizing the presently available classification schemes for new sequences is difficult and may result in a reduction of the quality of functional inference.

The basic principle for classifying a new sequence is to test its identity and similarity to all other sequences in the database. Should it have a sequence identity

of at least 30% with another sequence (throughout most of its length), a structural or a functional inference becomes valid. More accurately, functional annotation such as catalytic function of an enzyme requires at least 40% identity in amino acids (Rost, 2002), while reliable structural inference can be based on 25-30% sequence identity over at least 100 amino acids (Rost, 1999). Below these levels of pairwise sequence similarity, no unsupervised search method can suggest, with a reasonable confidence, the structural or functional relatedness to the closest homologue. While methods based on BLAST-based iterative profiles, Hidden Markov Models and methods for finding connections through intermediate sequences are superior in their sensitivity, they suffer from a higher level of false positives (Park et al., 1998). Specifically, PSI-BLAST, which is among the best search methods (Jones and Swindells, 2002), misses ~50% of all homologies when fixing the rate of false positives to 1/1000 (Karwath and King, 2002). In the case that the sequence similarity is below the level of secure inference (often referred to as the “twilight zone”), finding the protein family borders and the evolutionary remote homologues of a given protein is desirable. One should bear in mind that naïve sequence-based search methods that are based on similarity scoring are often inappropriate for multi-domain proteins and for proteins that are a result of convergent evolution.

The inherent limitation of detecting remote homologues was estimated based on a structural benchmark (Miller et al., 1999). Several methodologies for navigating in the protein space for a better detection of remote homologues were proposed. These include sequence hopping (Holm and Sander, 1997), transitivity of homology (Brenner et al., 1998), intermediate sequences (Park et al., 1998), tuned iterative-profile (Schaffer et al., 2001), profile-profile (Yona and Levitt, 2002), jumping alignments (Spang et al., 2002), string-based kernels (Leslie et al., 2002), and some hybrid methods (Jaroszewski et al., 2000).

Herein we suggest that a construction of the protein sequence space into a connected graph that captures the trace of evolutionary relatedness is useful in finding remote homologues. This method achieves satisfactory results that are comparable with results that rely on additional structural and functional knowledge. We suggest that the scaffold of the protein space as presented by ProtoNet (Sasson

et al., 2003) can be used to infer structural and functional information for remote homologues with high reliability. The success in inferring structural and functional relatedness for selected families by the ProtoNet navigation method is discussed.

Methodology

Methods for detecting remote homologues

A comparative study that tested the success in detecting remote homologues by various published strategies had been published (Holm, 1998). The list of protein homologues was manually selected to represent both structurally and functionally (i.e. similar enzymatic activities) related families of remote homology. The proteins in the test set were characterized by their very low pairwise sequence similarity.

For each homologous family, multiple strategies were applied and combined with the goal of unifying the proteins that belong to the same evolutionary-related group. The strategies that were used included advanced sequence-based search methods, including FASTA walk (Pearson, 1994), PSI-BLAST (Altschul et al., 1997), HMMer (Karplus et al., 1998), MAST (Bailey and Gribskov, 1998), and PROBE (Giles, 1992). Each of the methods initiated a search from a query protein from the group of proteins, and the number of proteins from the group that were identified by the search as remote homologues was noted (Table 1). We used this set as a benchmark for testing the ability of the ProtoNet cluster-map of the protein space to perform such a unification of protein families. We evaluated and quantified the degree of specificity and sensitivity in detecting those families.

ProtoNet graph

We base our study on the database of ProtoNet 2.4 (<http://www.protonet.cs.huji.ac.il>, Sasson et al., 2003). ProtoNet is an automatically generated hierarchical classification of all protein sequences from Swissprot (release 40.28 with 114,033 sequences). Each sequence is a vertex in a graph, and each pair of sequences defines an edge whose weight is the average pairwise BLAST E-score between those two sequences. A pairwise merging process is conducted as

part of an agglomerative clustering, where the initial sequences are used as singleton clusters. The clustering process follows an order of similarity, where each merge is the most significant in terms of the lowest average E-score between two groups of sequences. The main principle applied in ProtoNet is the extensive use of restricted transitivity based on hidden intermediate sequences, coupled with a statistical and biological validation of the hierarchy of the resulting graph. The resulting structure of the graph allows one to navigate through it and to evaluate the capacity of the ProtoNet algorithm to detect remote homologues (see examples in Sasson et al., 2002).

Results and Discussion

Unifying protein families in ProtoNet

ProtoNet is a bottom-up tree construction in which biologically relevant connections between proteins are evident (Sasson et al., 2002). In the initial steps of the clustering process connections between closely related sequences are made, but at more advanced stages of the clustering process more questionable connections are included; at the root of the tree, all proteins are interconnected (except for a few hundred singletons). The hierarchical structure of the graph that is built by ProtoNet enables one to follow the connections between different groups of sequences when considering the progress of the clustering process. Intuitively, one can view it as a reflection of the evolutionary diverging process. As the clustering process advances, similar sequences aggregate to form groups that then unite to form sub-families. As the clustering algorithm progresses, the families are enlarged, thus forming superfamilies and making remote homologues prevalent. As the number of pairwise connections becomes sparse during the progress of the clustering algorithm (and the statistical significance of the pairwise connections are extremely low. i.e. BLAST E-score as low as 100), the representation of ProtoNet resembles a graph in which the edges connect groups of sequences (clusters) rather than individual sequences.

The properties of the graph in view of its biological content was evaluated by comparing each cluster to other classification methods such as InterPro (Mulder et al., 2002), SCOP (Lo Conte et al., 2000), ENZYMES (Bairoch, 2000), and other external annotation sources. To evaluate the quality of ProtoNet clusters for structural and functional properties, we scored all *best clusters*. We define such a score as a *correspondence score* (CS). The CS for a certain cluster and a given keyword (i.e. *Enolase N-terminal domain-like* family, in SCOP) measures the correlation between the cluster and that keyword, using the intersect-union ratio.

$CS(\text{cluster } C \text{ for keyword } K) = |c \cap k| / |c \cup k| = TP / (TP + FP + FN)$, where: c is the set of annotated proteins in cluster C , k is the set of proteins annotated with K , TP , FP , FN stand for true positives, false positives, and false negatives, respectively. TP = the number of proteins in cluster C that have keyword annotation K ; FP = the number of annotated proteins in cluster C that do not have keyword annotation K ; FN = the number of proteins not in cluster C that have keyword annotation K . The cluster receiving the maximal CS for keyword K is considered the cluster that best represents K within the ProtoNet tree, and is called the *best cluster*. The score for a given cluster on keyword K ranges from 0 (no correspondence) to 1 (maximal correspondence to the keyword, the cluster contains exactly all of the proteins with keyword K).

For annotation keywords related to structural properties, we used all keywords based on SCOP (fold, superfamily, family and domain levels) and for a functional view, we used the ENZYME (4 levels of EC hierarchy) and GO (in 3 categories - molecular function, cellular process and cellular localization) annotations. A complete list of all annotations with the associated ProtoNet *best clusters* is available (www.protonet.cs.huji.ac.il/best_cluster/).

The results from such analysis confirm that clusters that perfectly match functional and structural annotations as defined by the above listed sources are abundant within the ProtoNet graph. Representative results are shown in Table 1. The results for the SCOP family level (contains 1313 keywords) is $CS=0.84$ (with a higher score for proteins of a single domain, see discussion). By limiting the analysis for clusters size ranges from 20 to 1000 proteins, the CS is somewhat higher (0.86). The

average CS for best clusters for the E.C. ENZYME 4th level in hierarchy (contains 2157 keywords) is quite high, but drops significantly for the 3rd digit of the EC classification.

Properties	# keywords	# best clusters	CS ^a	Specificity ^b	Sensitivity ^c
SCOP Family	1313				
All		1283	0.84	0.96	0.88
20-1000		968	0.86	0.98	0.88
SCOP Superfamily	872				
All		854	0.77	0.95	0.79
20-1000		607	0.78	0.97	0.80
Enzyme EC. 4 th digit	2157				
All		2157	0.75	0.88	0.84
20-1000		1581	0.76	0.88	0.85
Enzyme EC. 3 rd digit	206				
All		206	0.50	0.87	0.55
20-1000		150	0.50	0.89	0.55

Table 1. A global analysis for the correspondence of ProtoNet clusters with structural and functional properties. ^aCorrespondence score. ^bSpecificity = TP / (TP+FP). ^cSensitivity = TP / (TP+FN). See text for definitions.

Results from the *best cluster* view point to a relatively narrow interval along the clustering process where the connection between protein families starts to dominate. Thus, two protein families that share the same cluster in the cluster-map described at that interval may be viewed as remote homologues. This last principle guided us when we tested for the ability of ProtoNet system to suggest remote homologues.

To compare the ability of ProtoNet to detect remote homologues to that of other available search methods, we choose to test a manually selected set of remote homologues as presented in a comparative survey (Holm, 1998). In Table 2, examples of proteins that share structural and functional relatedness are listed. We duplicated the comparative analysis to include ProtoNet method as an additional, testable methodology in view of the previously described five state of the art methods for remote homologues detection.

The remote family representatives are marked by the characteristics of the relevant protein functional family (i.e., Winged HTH DNA-binding domains, Table 2). We began each search with the query protein listed in Table 2, and ‘climbed’ the

cluster hierarchy until reaching the *largest* cluster in which the target proteins that are included in the defined functional family are encountered. The process is aborted and no success is recorded when the cluster that includes the query and the target proteins annotate less than 85% of the proteins in the cluster (denoted as TP - true positives; un-annotated proteins are not considered for such procedure). The number of protein families and representatives that were united under that search was reported. In the case that the selected largest cluster is already contaminated by unrelated identifiers or keywords we marked it (by minus sign, Table 2) to indicate a failure in the purity specificity (purity) of the cluster.

Functional family	Query sequence	A	B	C	D	E	PN	Cluster (# of proteins)	Remarks
Winged HTH DNA-binding domains (1opc, 1aoy, 1dprB, 1hst, 1ecl, 1bgw, 1smtA, 1lea, 1fokA)	1. ompr_ecoli (1opc)	2	1	1	1	2	2-	226420 (1088)	low TP, Detected: 1b9w
	2. argr_ecoli	1	1	1	1	1	1		
	3. 1lea	/	/	/	/	/	3	226170 (348)	TP = 1 for SF Detected: 1smtA, 1dprB
Nucleic-acid binding proteins with OB-fold 1. S1 domains, translation initiation factor IF1 (1ah9) 2. Cold shock protein (1mjc) (1asyA, 1lylA, 1cuk, 3ullA, 1cmkA, 1pfsA)	1. 1ah9	1	1	1	1	1	1		
	2. cspa_ecoli	1	1	1	1	0	1		
	3. 1asy	/	/	/	/	/	2	210265 (143)	Detected: 1lyl
1. Pertusis toxin 2. Aerolysin (domain) 3. LINK domains 4. CTL domain	1. tox3_borpe	1	1	0	1	0	1		
	2. lems_human	Ex	1	0	1	1	1		
1. Chromo domain 2. DNA-binding proteins (1sap)	mod3_human	1	1	1	1	1	1		
	1sap	1	1	1	1	1	1		
1. Mite allergen 2. Immunoglobulin-like	def2_derfa	1	1	1	1	1	1		
Urease **	ure1_kleae	4 [+6]	7	10	1	3	10	225659 (253)	TP = ~86%
1. PAPS reductase 2. ATP sulphurylase 3. N-type ATP Ppases (1nsyA, 1gpmA)	1. mt16_yeast	1	3	2	2	3	3-	226597 (841)	low TP for all three; high TP for 1 and 2 in cluster 212193.
	2.nade_bacsu	1	3	3	1	1			
	3. cysd_ecoli	3	3	3	3	3			
1. HIT (4rhn) 2. Ap ₄ A hydrolase 3. GalT (1hxpA)	1. fhit_human	2 [+1]	3	3	3	1	2	221320 (48)	Could not find Ap ₄ A proteins through search. High TP
	2. gal7_ecoli	1	1	3	1	1			
1. Covalent NAD-binding oxidases (1ahu) 2. UDP-N-acetylenolpyruvoylglucosamine (2mbr)	1. 1ahu	1	2	2	0	1	2	225895 (79)	High TP
	2. murb_ecoli	1 [+1]	1	1	2	0			
Coenzyme A transferase 1. Glutaconate 2. Succinyl, acetate, 3-oxoadipate, butyrate, etc	1. X81440	2	2	2	2	2	/	225850 (51)	X81440 not present in Swissprot. The 2 nd query unified this family
	2. atoa_ecoli	2	2	1	1	1	2 (*)		
1. Dioxygenase (1han) 2. Glyoxylase (1froA)	1. bphc_burce	1	2	2	1	2	2	222839 (70)	TP=1 for SF Including also 1byl from this SF
	2. igul_human	1	2	2	2	1			

Table 2. *Comparing the success in detecting remote homologues by alternative methods.* Functional families that unify two or more sequence families were combined as in (Holm, 1998). The success in detection is marked by the number of successes in view of the predetermined representatives for remote homology. Alternative methods for remote homologue detection were applied over single query proteins and an enumeration of the sequence families that were recognized by each of the methods was recorded (a result of 1 means the query is recovered; zero indicates failure to identify it). A search in ProtoNet tree was performed in the same way by climbing the cluster hierarchy starting from a query protein and determining the cluster in which the protein identifier of the target proteins reached > 85% TP. The methods applied are abbreviated as follows: A. FASTA walk; B. PSI-BLAST; C. PROBE; D. GRIBSKOV and E. HMMer (according to Holm, 1998), PN. ProtoNet. The mark (-) indicates contamination by unrelated structural families. The mark (/) indicates a query examined solely by ProtoNet. Ex – the number of protein exploded. (*) The entry X81440 could not be matched in ProtoNet's proteins, so we indicated the smallest cluster that unites both families of that query. (**) Discussed in detail in the text.

Evaluating the success in detecting remote homologues

Inspection of the results from Table 2 indicates that our simple procedure using the scaffold of ProtoNet clusters performs as well as the best method out of the five that were tested for almost every query. Note that the benchmark was completely independent of our work and was selected to evaluate the performance of the available five independent methods for detecting remote homologues. The different methods are based on exhaustive dynamic programming (FASTA), Iterative profile search (PSI-BLAST, PROBE) and traditional profile-based searches (GRIBSKOV, HMMer). For details on the listed methods, see (Holm, 1998).

A strong characteristic of our method is its symmetry in the results between query and target protein. This is simply due to the symmetry of the search up a cluster tree – the cluster that unites a query **A** with a target **B** is the same cluster that unites **B** with **A** when **B** is the query and **A** is the target. This symmetry property implies *consistency* of the search by the ProtoNet method. It can be appreciated from Table 1 in that the other methods do not perform symmetrically.

Other main characteristics of our method are the *simplicity* and *speed* of the search, which requires only basic ProtoNet navigation skills of the user, and its *robustness*, as the cluster map platform already contains all sequences from Swissprot (ver. 40.28). Every sequence can be referred to as a query or target sequence for such a search.

While performing our searches we were able to find in some functional families more examples of remote homologues that were not queried in the original benchmark (marked in bold, Table 2). These examples were found as a byproduct of the search in ProtoNet, and were added as additional representatives to enrich the list in Table 2.

The protein family of the urease-related hydrolases serves to illustrate the traces of diversifying proteins throughout evolution. Several methods were applied to unify the members of the urease family (Heger and Holm, 2003). All methods used structural alignment, sequence alignment and profile-based search (Holm and Sander, 1997). The entire group, as defined in the original study (Holm and Sander, 1997), contains 12 protein families that share structurally similar active-site architecture. Of those proteins, 9 are members of the *metallo-dependent hydrolases* SCOP superfamily, 2 are members of a related superfamily named *composite domain of metallo-dependent hydrolases*, an additional one belongs to a different superfamily (*arylphosphatase*). Of the 12 hydrolases that are in the DB, ProtoNet defined 10 of them in a relatively pure cluster (A225659, 253 proteins). This cluster best represents the *metallo-dependent hydrolases* superfamily hierarchy.

Inspecting this cluster obtained in ProtoNet shows that it includes all previously described urease-related superfamily proteins, as well as additional yet undefined related members as shown in Table 2. The different representatives in this large family all share an enzymatic activity marked by EC 3.5.-.-, with the exception of proteins having a different catalytic activity (EC 3.8.-.-). Interestingly, the evolutionary and functional relatedness between those two apparently different activities was confirmed (Sadowsky et al., 1998). Additionally, not only is the high coverage of the urease-related superfamily (in cluster A225659) validated, but

monitoring the descendant clusters in the cluster-map provides a good description of the partition of this functional family into more refined subfamilies. Relationships between the subfamilies can be proposed and the divergence process of the family is traceable.

Enzyme name	# of proteins in the cluster / # in DB	E.C.
(1) Adenosine deaminase	28 / 28	3.5.4.4
(2) Adenine deaminase,	6 / 6	3.5.4.2
(3) Allantoinase	6 / 6	3.5.2.5
(4) AMP deaminase	11 / 11	3.5.4.6
(5) N-isopropylammelide isopropylamidohydrolase	1 / 1	3.5.99.4
(6) Atrazine chlorohydrolase	1 / 1	3.8.1.8
(7) Hydroxydechloroatrazine ethylammonohydrolase	1 / 1	3.5.99.3
(8) Cytosine deaminase	1 / 3	3.5.4.1
(9) Dihydropyrimidinase 2	10 / 10	3.5.2.2
(10) Guanine deaminase	7 / 8	3.5.4.3
(11) Imidazolonepropionase	23 / 23	3.5.2.7
(12) N-acetylglucosamine-6-phosphate deacetylase	6 / 6	3.5.1.25
(13) N-acyl-D-aspartate deacylase	1 / 1	3.5.1.83
(14) D-aminoacylase	1 / 1	3.5.1.81
(15) N-acyl-D-glutamate deacylase	1 / 1	3.5.1.82
(16) Dihydroorotase (DHOase)	49 / 49	3.5.2.3
(17) Urease (Urea amidohydrolase)	78 / 78	3.5.1.5

Table 3. *The enzymatic groups that are associated with ProtoNet cluster A225659. The annotations for 17 different groups are extracted from the ENZYME database. Note that for almost all groups the coverage is complete, with no additional proteins carrying the same enzymatic activity reported in the database used.*

Fig. 1 shows the pairwise similarity within proteins of the cluster A225659 (253 proteins) as a matrix reflecting all-against-all pairwise similarity BLAST E-scores. Darkest color indicates a maximal BLAST E-score (E-score = 0). The gradient in color intensity indicates a decrease in the similarity score. White indicates an E-score that is worse than 100. The domination of the white color reflects the very low pairwise similarity among proteins in the cluster. Among all possible pairwise connections, ~70% are worse than E-score 100, indicating the sequence remoteness resulting from the BLAST search. The order of the proteins in the matrix reflects the construction of the ProtoNet tree for this sub-graph (starting at the bottom right towards top left). An interactive illustration (color-coded) of the weak connections between all 253 proteins in this cluster is provided at www.protonet.cs.huji.ac.il/homologues/.

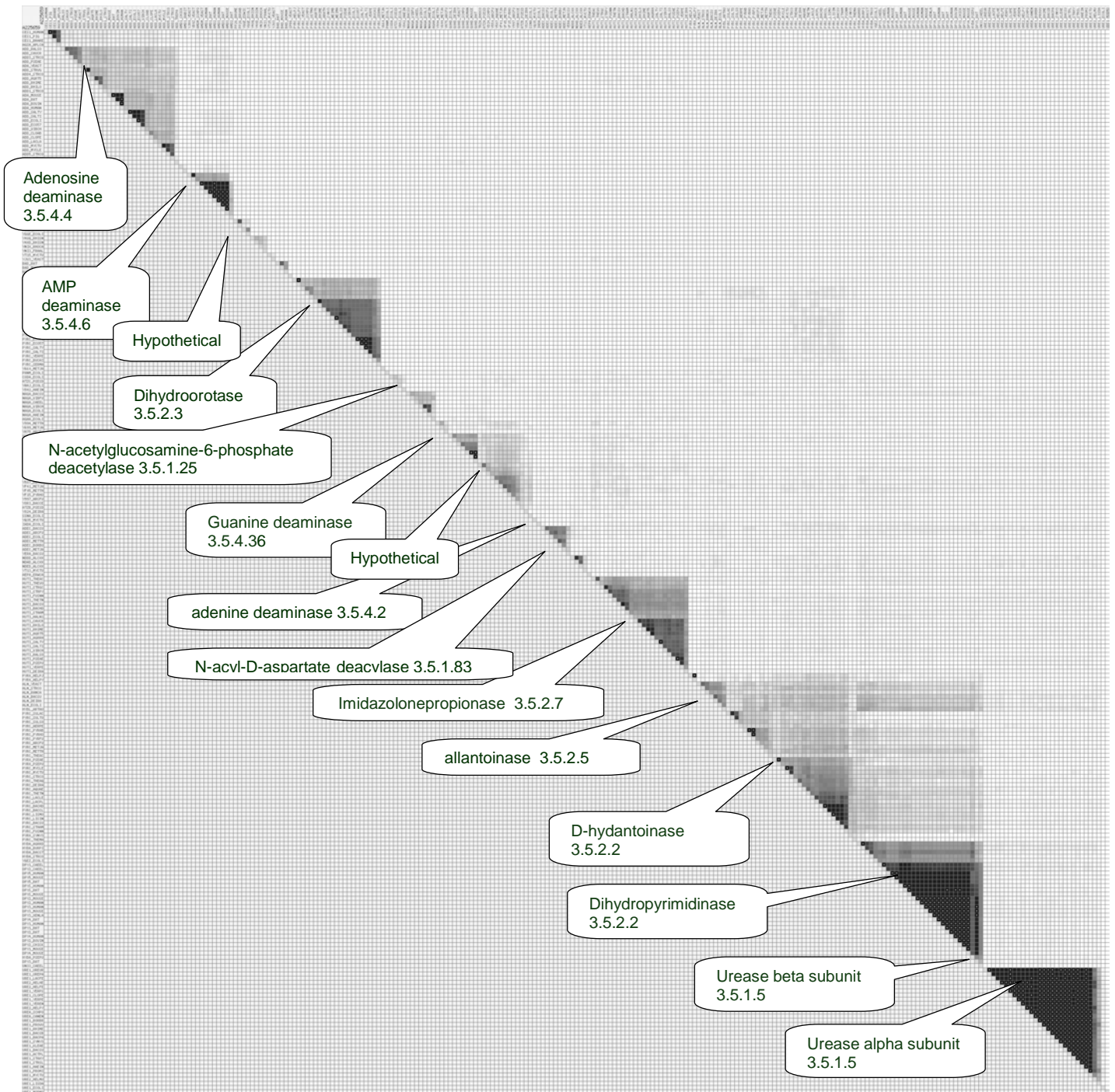


Figure 1. All-against-all BLAST E-score in a sub-tree of ProtoNet for urease related proteins. Representative enzymatic activity of some of the proteins is shown in the callouts. The matrix shows all 253 proteins in cluster A225659. Black indicates maximal BLAST E-score (E-score=0). Dark grey indicates very significant pairwise similarity. A gradient of black to white parallels a decrease in the similarity score.

In addition to the examples described in Table 3 and in Fig. 1, additional remote homologues are easily detected. For example, the proteins that share the histone core proteins H2A, H2B, H3 and H4 are linked together in a sub-tree of ProtoNet with a variety of transcription factors, such as the transcriptional initiator TFIIB, TFIID, TATA box, CAAC box binding protein and more. Despite extremely low sequence similarity (cluster A224843, 284 proteins), all these proteins share an identical structural fold (Baxevanis and Landsman, 1998).

In another very interesting case, similarity between alcohol dehydrogenases and crystalline of the eye and some proteins in synaptic vesicles in the nerve system has been proposed (Linial and Levius, 1993). Such homology is apparent at the sequence similarity level only through some weak connecting intermediate sequences. Proteins of the alcohol-dehydrogenases (also polyol-, threonine-, archaeon glucose-dehydrogenases), eye lens zeta-crystallins, E. coli quinone oxidoreductase, synaptic vesicle VAT-1 proteins, enoyl reductases of mammalian fatty acid and yeast erythronolide synthases share the same structural fold with only very few conserved amino acids throughout this diverged superfamily (Persson et al., 1994). This superfamily was termed medium-chain dehydrogenases/reductases (MDR). A cluster in ProtoNet with 226 proteins (cluster A220748) combines all the MDR proteins that carry very different functional role but nevertheless, share a high degree of structural and biophysical properties. By using the interactive visualization tools within ProtoNet (<http://www.protonet.cs.huji.ac.il/>), the connectivity among the proteins of the MDR superfamily is shown. Note that the intermediate sequences that connect specific families are traceable by inspecting the pairwise matrix (as in Fig. 1 and see www.protonet.cs.huji.ac.il/homologues/).

The ProtoNet tree allows for the testing of the connectivity among proteins, thus detecting remote homologues (Fig. 2). However, our examples are based on analyzing proteins from the Swissprot database (~114.000 proteins). As the number of protein sequences currently available is at least 10 times higher, the relevance and the robustness of our method towards larger database should be re-evaluated. We developed an extended version of ProtoNet that includes over one million proteins (ProtoNet 3.0), based on Swissprot 41.12 combined with TrEMBL 24.8.

Many of the proteins in this combined dataset are still poorly annotated. In ProtoNet 3.0, nodes in the tree that are candidates for representing clusters of remote homologues are eventually much larger. We tested the robustness of the results in the extended version of ProtoNet. In general, despite ~9 fold increase in the database size and a reduced quantity and quality of the associated annotations, the performance of the system is only slightly reduced. We illustrate it for the urease superfamily that was discussed above. A cluster representing the urease-related superfamily (A293766, ProtoNet 3.0) is composed of 1508 proteins that include 16 out of the 17 groups (as in Table 2) with a very high sensitivity (i.e., collecting all proteins from the database that belong to this enzymatic group). One additional urease related enzymatic group (N-acyl-D-amino acid deacylase) that was missed from the Swissprot based cluster (Table 2) is detected in cluster A293766. This example illustrates the potential of the ProtoNet tree to detect remote homologues even in a relatively poorly annotated database. It is worth noticing that among the 1508 proteins, most of them (58%) have no Prosite annotation and 21.1% are marked as hypothetical proteins. Thus, the inference of functional groups is a direct byproduct of navigating the ProtoNet tree.

Careful inspection of the families with remote homologues showed that the performance of ProtoNet is very satisfactory for single domain proteins (as in the case of histones, urease and MDR proteins; a much lower success can be obtained for proteins that are multi-domains. The relatively low scores in the *best cluster* analysis (www.protonet.cs.huji.ac.il/best_cluster/) best reflect such instances. While ProtoNet is based on a whole protein clustering, and remote homology (especially from a structural perspective) is often associated with a domain level, it is appropriate to use this methodology with a critical objective.

In summary, we offer the use of ProtoNet's cluster-map as a platform for searching for remote homologues and for identification of weak connections between functional families. We propose a simple search procedure for finding remote homologues for a query sequence, and compare it to a benchmark of performance for different methods that rely on information beyond the elementary pairwise similarity. We show that in almost all cases, our search method performs

as good as the best method. Our method is very fast and consistent due to the fact that it relies on a pre-determined scaffold of the protein space that is in accord with the process of evolutionary diversification.

Partitioning the protein space into homologous families is motivated by Structural and Functional Genomics initiatives. As seen from the benchmark and the additional examples, the identification of remote homologues for a given sequence extends our ability in predicting structural relatedness and in many instances also the functional characteristics. We propose a robust method for finding such remote homologues and defining structural and functional superfamilies.

Acknowledgements

We would like to thank the ProtoNet team and especially Uri Inbar and Hillel Fleischer for excellent database management. We thank Noam Kaplan and Menachem Fromer for a critical reading and useful suggestions. This study is partially supported by the Israeli Ministry of Defense and the NIH support for the CESH Structural Genomics Consortium. We thank the Sudarsky Center for Computational Biology in the Hebrew University for a fellowship support (O.S).

References

- Altschul, S. F., Madden, T. L., Schaffer, A. A., Zhang, J., Zhang, Z., Miller, W., and Lipman, D. J. (1997). Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 25, 3389-3402.
- Bailey, T. L., and Gribskov, M. (1998). Methods and statistics for combining motif match scores. *J Comput Biol* 5, 211-221.
- Bairoch, A. (2000). The ENZYME database in 2000. *Nucleic Acids Res* 28, 304-305.
- Bateman, A., Birney, E., Durbin, R., Eddy, S. R., Howe, K. L., and Sonnhammer, E. L. (2000). The Pfam protein families database. *Nucleic Acids Res* 28, 263-266.
- Baxeavanis, A. D., and Landsman, D. (1998). Histone Sequence Database: new histone fold family members. *Nucleic Acids Res* 26, 372-375.

- Brenner, S. E., Chothia, C., and Hubbard, T. J. (1998). Assessing sequence comparison methods with reliable structurally identified distant evolutionary relationships. *Proceedings of the National Academy of Sciences of the United States of America* 95, 6073-6078.
- Giles, I. G. (1992). PROBE: a computer program to scan DNA sequence databases for the existence of potential probe sequences in DNA. *Biochem Soc Trans* 20, 292S.
- Heger, A., and Holm, L. (2003). Sensitive pattern discovery with 'fuzzy' alignments of distantly related proteins. *Bioinformatics* 19 *Suppl* 1, I130-I137.
- Holm, L. (1998). Unification of protein families. *Current Opinion in Structural Biology* 8, 372-379.
- Holm, L., and Sander, C. (1997). An evolutionary treasure: unification of a broad set of amidohydrolases related to urease. *Proteins* 28, 72-82.
- Jaroszewski, L., Rychlewski, L., and Godzik, A. (2000). Improving the quality of twilight-zone alignments. *Protein Sci* 9, 1487-1496.
- Jones, D. T., and Swindells, M. B. (2002). Getting the most from PSI-BLAST. *Trends Biochem Sci* 27, 161-164.
- Karplus, K., Barrett, C., and Hughey, R. (1998). Hidden Markov models for detecting remote protein homologies. *Bioinformatics* 14, 846-856.
- Karwath, A., and King, R. D. (2002). Homology Induction: the use of machine learning to improve sequence similarity searches. *BMC Bioinformatics* 3, 11.
- Leslie, C., Eskin, E., and Noble, W. S. (2002). The spectrum kernel: a string kernel for SVM protein classification. *Pac Symp Biocomput*, 564-575.
- Linial, M., and Levius, O. (1993). VAT-1 from Torpedo is a membranous homologue of zeta crystallin. *FEBS Lett* 315, 91-94.
- Lo Conte, L., Ailey, B., Hubbard, T. J., Brenner, S. E., Murzin, A. G., and Chothia, C. (2000). SCOP: a structural classification of proteins database. *Nucleic Acids Res* 28, 257-259.
- Mulder, N. J., Apweiler, R., Attwood, T. K., Bairoch, A., Bateman, A., Binns, D., Biswas, M., Bradley, P., Bork, P., Bucher, P., *et al.* (2002). InterPro: an integrated documentation resource for protein families, domains and functional sites. *Brief Bioinform* 3, 225-235.
- Muller, A., MacCallum, R. M. and Sternberg, M. J. (1999) Benchmarking PSI-BLAST in genome annotation. *J Mol Biol* 293, 1257-1271.
- Sternberg, M. J., Park, J., Karplus, K., Barrett, C., Hughey, R., Haussler, D., Hubbard, T., and Chothia, C. (1998). Sequence comparisons using multiple sequences detect three times as many remote homologues as pairwise methods. *J Mol Biol* 284, 1201-1210.
- Pearson, W. R. (1994). Using the FASTA program to search protein and DNA sequence databases. *Methods Mol Biol* 24, 307-331.

- Persson, B., Zigler, J. S., Jr., and Jornvall, H. (1994). A super-family of medium-chain dehydrogenases/reductases (MDR). Sub-lines including zeta-crystallin, alcohol and polyol dehydrogenases, quinone oxidoreductase enoyl reductases, VAT-1 and other proteins. *Eur J Biochem* 226, 15-22.
- Rost, B. (1999). Twilight zone of protein sequence alignments. *Protein Eng* 12, 85-94.
- Rost, B. (2002). Enzyme function less conserved than anticipated. *Journal of Molecular Biology* 318, 595-608.
- Russell, R. B., Saqi, M. A., Sayle, R. A., Bates, P. A., and Sternberg, M. J. (1997). Recognition of analogous and homologous protein folds: analysis of sequence and structure conservation. *J Mol Biol* 269, 423-439.
- Sadowsky, M. J., Tong, Z., de Souza, M., and Wackett, L. P. (1998). AtzC is a new member of the amidohydrolase protein superfamily and is homologous to other atrazine-metabolizing enzymes. *Journal of Bacteriology* 180, 152-158.
- Saier, M. H., Jr. (1996). Phylogenetic approaches to the identification and characterization of protein families and superfamilies. *Microb Comp Genomics* 1, 129-150.
- Sasson, O., Linial, N., and Linial, M. (2002). The metric space of proteins-comparative study of clustering algorithms. *Bioinformatics* 18 Suppl 1, S14-21.
- Sasson, O., Vaaknin, A., Fleischer, H., Portugaly, E., Bilu, Y., Linial, N., and Linial, M. (2003). ProtoNet: hierarchical classification of the protein space. *Nucleic Acids Res* 31, 348-352.
- Schaffer, A. A., Aravind, L., Madden, T. L., Shavirin, S., Spouge, J. L., Wolf, Y. I., Koonin, E. V., and Altschul, S. F. (2001). Improving the accuracy of PSI-BLAST protein database searches with composition-based statistics and other refinements. *Nucleic Acids Research* 29, 2994-3005.
- Spang, R., Rehmsmeier, M., and Stoye, J. (2002). A novel approach to remote homology detection: jumping alignments. *J Comput Biol* 9, 747-760.
- Stark, A., Sunyaev, S., and Russell, R. B. (2003). A model for statistical significance of local similarities in structure. *J Mol Biol* 326, 1307-1316.
- Yona, G., and Levitt, M. (2002). Within the twilight zone: a sensitive profile-profile comparison tool based on information theory. *J Mol Biol* 315, 1257-1275.